# Expectation-enforcing strategies for repeated games

Nikos Dimou and Alex McAvoy*

**Abstract**

Originating in evolutionary game theory, the class of "zero-determinant" strategies enables a player to unilaterally enforce linear payoff relationships in simple repeated games. An upshot of this kind of payoff constraint is that it can shape the incentives for the opponent in a predetermined way. An example is when a player ensures that the agents get equal payoffs. While extensively studied in infinite-horizon games, extensions to discounted games, nonlinear payoff relationships, richer strategic environments, and behaviors with long memory remain incompletely understood. In this paper, we provide necessary and sufficient conditions for a player to enforce arbitrary payoff relationships (linear or nonlinear), in expectation, in discounted games. These conditions characterize precisely which payoff relationships are enforceable using strategies of arbitrary complexity. Our main result establishes that any such enforceable relationship can actually be implemented using a simple two-point reactive learning strategy, which conditions on the opponent's most recent action and the player's own previous mixed action, using information from only one round into the past. For additive payoff constraints, we show that enforcement is possible using even simpler (reactive) strategies that depend solely on the opponent's last move. In other words, this tractable class is universal within expectation-enforcing strategies. As examples, we apply these results to characterize extortionate, generous, equalizer, and fair strategies in the iterated prisoner's dilemma, asymmetric donation game, nonlinear donation game, and the hawk-dove game, identifying precisely when each class of strategy is enforceable and with what minimum discount factor.

## 1 Introduction

A foundational question in repeated games is how much an individual can shape long-run outcomes without imposing equilibrium discipline on others. Classical work, notably the Folk Theorem [1], characterizes the set of payoffs sustainable in equilibrium, provided the horizon of the game is sufficiently long. In contrast, far less is known about the unilateral problem: how much can a player control the space of possible outcomes, even when the opponent behaves arbitrarily? The space of feasible payoffs against a fixed strategy provides outcome-relevant information not always apparent from the behavioral mechanics of a strategy itself. Importantly, this perspective gives insights into the incentives an opponent faces, whether it be in an adaptive (evolutionary) setting or one of reinforcement learning. In fact, despite the rich space of subgame perfect equilibria in repeated games, reinforcement learning illustrates that gradual modifications of one's strategy based on self-interested objectives (i.e., caring about maximizing one's own payoff) leads to inefficient equilibria [2, 3], which raises a question about how to productively shape incentives.

---

*Please direct correspondence to A.M. (`alexmcavoy@gmail.com`).

Although we are ultimately interested in repeated games in a general setting, some of the motivation of our approach comes from evolutionary game theory, where there is a broad body of work on cooperation in simple repeated games like the prisoner's dilemma. The mechanism by which cooperation emerges in the repeated prisoner's dilemma ("direct reciprocity" [4]) is predicated on the fact that defection now can be punished in the future and cooperation now can be reciprocated in the future [5]. The underlying stage game has two actions, "cooperate" ($C$) and "defect" ($D$), and two agents facing one another receive payoffs according to the matrix

$$
\begin{array}{cc}
 & \begin{array}{cc} C & \quad D \end{array} \\
\begin{array}{c} C \\ D \end{array} & \left( \begin{array}{cc} R,\ R & S,\ T \\ T,\ S & P,\ P \end{array} \right),
\end{array}
\tag{1}
$$

where $T > R > P > S$. This payoff ranking ensures that defection is the dominant action, despite mutual cooperation yielding the socially optimal outcome.

For this kind of game, Press and Dyson [6] described a new class of "zero-determinant" (ZD) strategies that permit a striking level of control over expected linear payoff outcomes. In the infinitely repeated prisoner's dilemma, ZD strategies allow player $X$ to choose constants $\alpha$, $\beta$, and $\gamma$ and unilaterally enforce the equation $\alpha \pi_X + \beta \pi_Y + \gamma = 0$, where $\pi_X$ and $\pi_Y$ are the long-term average payoffs of $X$ and $Y$, respectively. This equation is "enforced" by $X$ in the sense that it holds regardless of the opponent's behavior in the repeated game. Of course, not all linear relationships are feasible, but several surprising classes of relationships are enforceable in the repeated prisoner's dilemma, including those in the family $\kappa - \pi_X = \chi (\kappa - \pi_Y)$ for $\chi \geqslant 1$ and $\kappa \in [P, R]$. Of particular interest are strategies enforcing $\pi_X - P = \chi (\pi_Y - P)$ for $\chi > 1$, which allow $X$ to effectively extort the opponent and claim an unfair share of payoffs beyond those of mutual defection. A concrete example of a ZD strategy is tit-for-tat (TFT) in the (infinite-horizon) iterated prisoner's dilemma [7, 8], which enforces equal payoffs ($\chi = 1$).

Since their discovery, ZD strategies have been investigated in numerous social dilemmas, including iterated public goods [9] and (nonlinear) donation games [10], and their properties have been studied in the contexts of discounted games [11], games with continuous action spaces [10], alternating games [12], evolutionary environments [13–15], multiplayer games [16], human experiments [17, 18], and stochastic, multi-state games [19]. However, the study of ZD strategies remains incomplete. Most work investigates linear payoff relationships in games where both players have access to only two actions. It also is not understood exactly what kinds of payoff constraints can be enforced, even in simple games.

Furthermore, ZD strategies are typically derived under the assumption that one is searching within the class of "memory-one" strategies, which specify the player's reaction depending on the outcome of the previous round only. This class of strategies strikes a balance between mathematical tractability and behavioral complexity, and many famous strategies (including tit-for-tat and win-stay, lose-shift [20]) condition on only the previous round's outcome. It has been shown that longer finite memory, which is computationally intractable [21], is not always needed against reasonable opponents, and strategies of reduced memory suffice to produce optimal payoff outcomes in some specific game structures [22, 23]. Even though longer-memory ZD strategies have been developed [24], any advantage they have over memory-one strategies remains unclear.

An additional constraint is imposed by discounting, which weights payoffs from future rounds less than current payoffs through a discount factor $\lambda \in [0, 1)$. Discounting reflects time

preferences or uncertainty about future interactions, and it restricts the set of enforceable payoff relationships: as $\lambda$ decreases (placing less weight on the future), fewer relationships can be enforced [11]. Given the restricted range of feasible ZD strategies in discounted games and the limited understanding of longer-memory properties, the following questions naturally arise:

- How much control does a player have against an adaptive opponent?

- Can a player enforce additional payoff relationships by extending memory?

This paper addresses both questions and provides a definitive answer to what is possible and what is implementable by a single player in repeated games. More precisely, we adopt the framework of "autocratic" strategies, which generalize ZD strategies and allow for nonlinear constraints on expected payoffs. We provide necessary and sufficient conditions of which payoff relationships can be enforced using strategies of arbitrary complexity, thus establishing a complete characterization on the level of control that a player is able to exert against any adaptive opponent. Our main contribution gives a (perhaps surprising) answer to the second question: extending memory beyond a simple reactive learning structure provides no additional power. In fact, any enforceable payoff relationship can be implemented using a two-point reactive learning strategy. Reactive learning strategies generalize memory-one strategies by conditioning on the opponent's most recent action and the player's own previous mixed action (rather than realized action), allowing the player to track their own randomization history. A two-point reactive learning strategy further simplifies this by mixing between just two (fixed) possible responses based on the opponent's last action. In addition, for additive objective functions, we show that enforcement is possible using even simpler reactive strategies that depend solely on the opponent's last action. These results demonstrate that reactive learning strategies are universal within the class of all strategies that endow a player with the ability to control expected payoff outcomes.

For any candidate payoff constraint, we give a concrete "next-round correction" condition that is both necessary and sufficient for enforcement. Furthermore, we provide an explicit formula for the minimum discount factor required to enforce a given payoff relationship. This result resolves a number of open problems. First, the problem of identifying the minimum discount factor in games with a more complex structure than that of the iterated prisoner's dilemma has appeared previously in several works, e.g., in [25]. Second, it answers an open question raised by Hilbe et al. [14] regarding the existence of autocratic strategies in discounted games, and extends prior work [11] by providing exact thresholds rather than just existence results.

In contrast to the approach that resulted in the discovery of ZD strategies [6], we do not rely on the standard method of studying the transition probabilities of memory-one strategies adopted by the focal player, as this technique turns out to be prohibitive when dealing with strategies of broad cognitive complexity. Instead, we explore the dynamics of the repeated game by identifying a mechanism that allows one to control the stochastic path of the game and correct past "suboptimal" behaviors or errors. Although we focus on discounted games, which are more realistic, we also extend our techniques to the infinite-horizon regime.

The main theoretical results are applied to several classical games. In the iterated prisoner's dilemma, we provide exact conditions for extortionate, generous, and equalizing strategies, computing the minimum discount factor required for each class. Our results verify and extend prior work [11, 26] by providing constructive formulas. The generality of our framework also enables important non-existence results. For instance, we prove that symmetric relationships (such as

3

$\pi_X = \pi_Y$ in the infinitely repeated prisoner's dilemma) cannot be enforced in properly discounted games of finite horizon, resolving the question of when fair strategies exist. This shows that tit-for-tat's ability to enforce equal payoffs is fundamentally limited to the infinite-horizon setting. Beyond symmetric, two-action games, we characterize autocratic strategies in a multi-action nonlinear donation game, an asymmetric donation game, and the hawk-dove game. The nonlinear donation game demonstrates that our framework extends naturally to nonlinear payoff relationships, where traditional ZD techniques fail. The asymmetric donation game illustrates the distinction between equality (enforcing $\pi_X = \pi_Y$) and fairness (proportional sharing based on action costs), showing that only the former can sometimes be enforced unilaterally.

A final important algorithmic consequence of our characterization is that verifying whether a given payoff relationship is enforceable, as well as computing the minimum discount factor and constructing an associated autocratic strategy, can be accomplished in polynomial time using linear programming. This stands in stark contrast to the computational intractability of analyzing general behavioral strategies [21, 27]. Our results thus provide both theoretical closure on the memory question of enforcing payoff constraints, as well as practical tools for computing autocratic strategies.

The remainder of this paper is organized as follows. In Section 2, we review the framework of repeated games and introduce the notion of autocratic strategies, establishing key background results including the pointwise and generalized "next-round correction" conditions. In Section 3, we present our main theoretical contributions: we prove that any enforceable payoff relationship can be implemented using a two-point reactive learning strategy (Theorem 1), characterize the minimum discount factor required for enforcement (Proposition 4), and show that additive payoff constraints admit even simpler reactive implementations (Theorem 2). We also extend our results to the undiscounted setting ($\lambda \to 1$). In Section 4, we establish computational and structural properties of autocratic strategies, including polynomial-time algorithms for verification and construction, and convexity properties of the space of enforceable relations. Finally, in Section 5, we apply our framework to four examples of social dilemmas [28]: the iterated prisoner's dilemma, a nonlinear donation game, an asymmetric donation game, and the hawk-dove game, providing complete characterizations of all linear payoff relationships that can be unilaterally enforced, as well as the minimum average time horizon needed to do so.

## 2 Background and auxiliary results

### 2.1 Repeated games and discounting

We consider repeated games between two players, $X$ and $Y$, with finite action spaces $S_X$ and $S_Y$, respectively. The players receive short-term payoffs via functions $u_X : S_X \times S_Y \to \mathbb{R}$ and $u_Y : S_X \times S_Y \to R$. In each round $t \in \{0, 1, 2, \dots\}$, players simultaneously choose actions $(s_X^t, s_Y^t) \in S_X \times S_Y$ from a distribution and receive stage-game payoffs $u_X\left(s_X^t, s_Y^t\right)$ and $u_Y\left(s_X^t, s_Y^t\right)$. The game continues with probability $\lambda \in [0, 1)$ after each round. The discounted payoff for player $X$ over a realization of action outcomes is

$$\pi_X = (1 - \lambda) \sum_{t=0}^{\infty} \lambda^t u_X\left(s_X^t, s_Y^t\right), \tag{2}$$

with a similar expression for player $Y$, with $u_Y$ replacing $u_X$. The factor $1 - \lambda$ normalizes payoffs to be comparable across different continuation probabilities.

4

## 2.2 Behavioral strategies and histories

Let $\mathcal{H}^T = (S_X \times S_Y)^T$ denote the set of all possible histories of length $T$, representing the sequence of action pairs played from round 0 through round $T-1$. We write $\mathcal{H} = \bigcup_{T \geqslant 0} \mathcal{H}^T$ for the set of all finite histories, where $\mathcal{H}^0 = \{\varnothing\}$ represents the null history at the start of the game.

Let $\Delta(S_X)$ and $\Delta(S_Y)$ denote the corresponding spaces of mixed actions (distributions over pure actions). A behavioral strategy for player $X$ is a map $\sigma_X : \mathcal{H} \to \Delta(S_X)$ that specifies a mixed action for each possible history. Given strategies $\sigma_X$ and $\sigma_Y$, we write $\mathbb{E}_{\sigma_X,\sigma_Y}[\cdot]$ for the expected value with respect to the probability distribution over infinite sequences of realized play induced by $\sigma_X$ and $\sigma_Y$. Notably, the mean long-term payoffs for $X$ and $Y$, respectively, are

$$\pi_X := \mathbb{E}_{\sigma_X,\sigma_Y}\left[(1-\lambda)\sum_{t=0}^{\infty} \lambda^t u_X\left(s_X^t, s_Y^t\right)\right]; \tag{3a}$$

$$\pi_Y := \mathbb{E}_{\sigma_X,\sigma_Y}\left[(1-\lambda)\sum_{t=0}^{\infty} \lambda^t u_Y\left(s_X^t, s_Y^t\right)\right]. \tag{3b}$$

A behavioral strategy is a "memory-one" strategy if its response depends only on the most recent round of play. Formally, a memory-one strategy for player $X$ consists of *(i)* an initial mixed action $\sigma_X^0 \in \Delta(S_X)$ and *(ii)* a response rule $\sigma_X : S_X \times S_Y \to \Delta(S_X)$ that specifies a mixed action $\sigma_X[s_X, s_Y] \in \Delta(S_X)$ for each action pair $(s_X, s_Y)$ in the previous round. Importantly, the response depends only on the realized actions $(s_X, s_Y)$, not on how those actions were generated through randomization. Let $\mathbf{Mem}_X^1$ denote the set of all memory-one strategies for player $X$.

## 2.3 Reactive learning strategies

Reactive learning strategies, introduced by McAvoy and Nowak [29], generalize memory-one strategies by allowing a player to condition on their own mixed action from the previous round.

**Definition 1.** A reactive learning strategy for $X$ consists of an initial action $\sigma_X^0 \in \Delta(S_X)$ and a response rule $\sigma_X : \Delta(S_X) \times S_Y \to \Delta(S_X)$.

Let $\mathbf{RL}_X$ denote the set of all reactive learning strategies for $X$. Every memory-one strategy naturally induces a reactive learning strategy via the canonical embedding

$$^* : \mathbf{Mem}_X^1 \longrightarrow \mathbf{RL}_X$$
$$: \left(\sigma_X^0, \sigma_X\right) \longmapsto \left(\sigma_X^0, \sigma_X^*\right), \tag{4}$$

where the initial actions are the same, and $\sigma_X^*[\tau_X, s_Y](\cdot) := \mathbb{E}_{s_X \sim \tau_X}[\sigma_X[s_X, s_Y](\cdot)]$ over $S_X$. Conversely, restricting a reactive learning strategy to Dirac measures recovers a memory-one response function, $\sigma_X[s_X, s_Y] := \sigma_X^*[\delta_{s_X}, s_Y]$, where $\delta_{s_X}$ is the point mass at $s_X$.

**Definition 2.** A sequence $\left\{\tau_X^t\right\}_{t=0}^{\infty} \subseteq \Delta(S_X)$ is a "chain" of mixed actions derived from $\left(\sigma_X^0, \sigma_X^*\right)$ if $\tau_X^0 = \sigma_X^0$ and, for every $t \geqslant 1$, there exists an action $s_Y^{t-1} \in S_Y$ such that $\tau_X^t = \sigma_X^*\left[\tau_X^{t-1}, s_Y^{t-1}\right]$.

5

## 2.4 Feasible and enforceable payoffs

The feasible region is the set of all payoff pairs $(\pi_Y, \pi_X)$ that can arise from some pair of strategies (note the unusual ordering, which we fix throughout). For a generic two-player game, this forms a convex subset of $\mathbb{R}^2$ of full dimension. When a player commits to a strategy, $\sigma_X$, the achievable payoffs as $\sigma_Y$ varies form a convex subset of the feasible region. For a generic strategy $\sigma_X$, characterizing the geometry of achievable payoff pairs $(\pi_Y, \pi_X)$ as $\sigma_Y$ ranges over all opponent strategies requires detailed knowledge of the repeated game dynamics. There is currently no simple way to determine this payoff region, apart from in special cases, such as two-action games with infinite horizon, where this region is the convex hull of at most 11 points [29].

## 2.5 Autocratic strategies

The following definition formalizes the notion of unilateral enforcement of expectations:

**Definition 3.** Let $\varphi : S_X \times S_Y \to \mathbb{R}$ be a fixed function on the space of joint actions. A strategy $\sigma_X : \mathcal{H} \to \Delta(S_X)$ is $(\varphi, \lambda)$-autocratic if, for all opponent strategies $\sigma_Y : \mathcal{H} \to \Delta(S_Y)$,

$$\mathbb{E}_{\sigma_X, \sigma_Y}\left[(1-\lambda)\sum_{t=0}^{\infty} \lambda^t \varphi\left(s_X^t, s_Y^t\right)\right] = 0. \tag{5}$$

Intuitively, an autocratic strategy unilaterally enforces a constraint on the expected value of $\varphi$, regardless of the opponent's behavior. When $\varphi(s_X, s_Y) = \alpha u_X(s_X, s_Y) + \beta u_Y(s_X, s_Y) + \gamma$, this corresponds to enforcing the linear payoff relationship $\alpha\pi_X + \beta\pi_Y + \gamma = 0$, which is exactly the motivation behind zero-determinant (ZD) strategies [6].

**Definition 4.** For $\varphi : S_X \times S_Y \to \mathbb{R}$ and $\lambda \in [0,1]$, we say that $\varphi \equiv 0$ is $\lambda$-enforceable if there exists a $(\varphi, \lambda)$-autocratic strategy $\sigma_X : \mathcal{H} \to \Delta(S_X)$. We say $\varphi \equiv 0$ is enforceable if it is $\lambda$-enforceable for some $\lambda \in [0,1]$.

## 2.6 The pointwise next-round correction condition

The following result from McAvoy and Hauert [10] provides a sufficient condition for a memory-one strategy to enforce a linear payoff relationship.

**Theorem 0** (McAvoy and Hauert [10]). Suppose that $\left(\sigma_X^0, \sigma_X[s_X, s_Y]\right)$ is a memory-one strategy for X. If there exists a function $\psi : S_X \to \mathbb{R}$ such that

$$\alpha u_X(s_X, s_Y) + \beta u_Y(s_X, s_Y) + \gamma = \psi(s_X) - \lambda\mathbb{E}_{s_X' \sim \sigma_X[s_X, s_Y]}\left[\psi\left(s_X'\right)\right] - (1-\lambda)\mathbb{E}_{s_X' \sim \sigma_X^0}\left[\psi\left(s_X'\right)\right] \tag{6}$$

holds for every $s_X \in S_X$ and $s_Y \in S_Y$, then $\left(\sigma_X^0, \sigma_X[s_X, s_Y]\right)$ enforces the linear payoff relationship

$$\alpha\pi_X + \beta\pi_Y + \gamma = 0 \tag{7}$$

against any behavioral strategy of player $Y$, including those with infinite memory.

We call **Eq. 6** the pointwise next-round correction condition, and we call $\psi$ the enforcement potential. This condition provides a local, action-by-action characterization that guarantees the global payoff constraint **Eq. 7**.
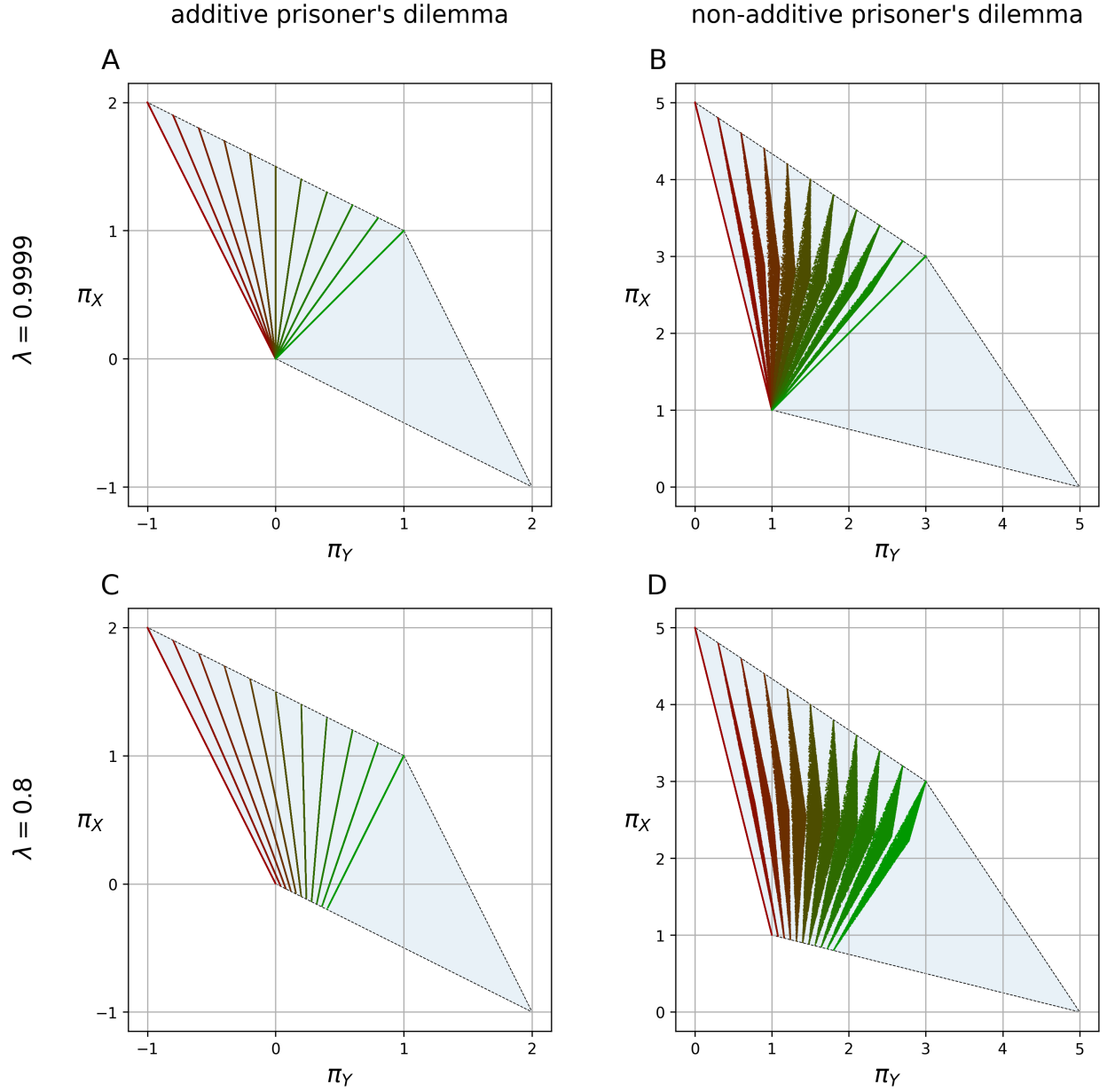
**Figure 1:** Payoff regions enforced when $X$ plays weighted averages of ALLD (red) and TFT (green) in repeated prisoner's dilemmas. Each colored region shows the payoff region obtained from the strategy $\sigma_X := (1-p)\,\text{ALLD} + p\,\text{TFT}$ played against $10^4$ randomly-chosen opposing strategies, for $p \in \{k/10\}_{k=0}^{10}$. (A,C) Additive prisoner's dilemma with $(R, S, T, P) = (1, -1, 2, 0)$: for $p \notin \{0, 1\}$, the strategy enforces a linear payoff relationship. (B,D) Non-additive prisoner's dilemma with $(R, S, T, P) = (3, 0, 5, 1)$: the strategy enforces a two-dimensional convex region. While line-enforcing strategies naturally arise in additive games through simple mixtures of well-known strategies, non-additive games require more sophisticated constructions. Panels A and B use $\lambda = 0.9999$ (a game with 10,000 rounds, on average, approximating an undiscounted game), and panels C and D use $\lambda = 0.8$ (a game with 5 rounds, on average).

**Example 1** (Tit-for-tat in the undiscounted prisoner's dilemma). Consider the prisoner's dilemma with payoff matrix **Eq. 1** in the undiscounted setting ($\lambda \to 1$). The well-known strategy of tit-for-tat (TFT) plays $C$ initially and then copies the opponent's previous action. That is, $\sigma_X^0 = C$ and $\sigma_X[s_X, s_Y] = s_Y$ for all $s_X, s_Y \in \{C, D\}$. It is known that TFT enforces the fair relationship $\pi_X = \pi_Y$, or equivalently, $\varphi \equiv 0$ where $\varphi(s_X, s_Y) = u_Y(s_X, s_Y) - u_X(s_X, s_Y)$ [6]. To verify this using the pointwise next-round correction condition (**Eq. 6**), we seek an enforcement potential $\psi : \{C, D\} \to \mathbb{R}$ such that

$$\varphi(s_X, s_Y) = \psi(s_X) - \psi(\sigma_X[s_X, s_Y]) \tag{8}$$

for all $s_X, s_Y \in \{C, D\}$. This equation simplifies to $\varphi(s_X, s_Y) = \psi(s_X) - \psi(s_Y)$ since TFT satisfies $\sigma_X[s_X, s_Y] = s_Y$. Taking $\psi(C) = 0$ and solving, we obtain $\psi(D) = T - S$. By Theorem 0 (extended to $\lambda = 1$), TFT enforces $\pi_X = \pi_Y$ against any opponent strategy. This example illustrates how the enforcement potential $\psi$ captures the "correction" that each action provides toward achieving the target payoff relationship.

A natural question is whether the pointwise next-round correction condition is also necessary for enforcing linear payoff relationships. The next example shows that it is not, in general:

**Example 2.** Consider a two-player, two-action game with payoff matrix

$$\begin{array}{c} \\ C \\ D \end{array} \begin{array}{cc} C & D \\ \left( \begin{array}{cc} -1 & -2 \\ +1 & +2 \end{array} \right). \end{array} \tag{9}$$

By playing $C$ and $D$ with equal probability in each round, player $X$ can ensure $\pi_X = 0$ regardless of $Y$'s strategy. Specifically, the constant mixed strategy with $\sigma_X^0(C) = \sigma_X^0(D) = 1/2$ and $\sigma_X[s_X, s_Y] = \sigma_X^0$ for all $s_X, s_Y \in \{C, D\}$ enforces $\pi_X = 0$ for all $\lambda \in [0, 1)$.

However, this strategy does not satisfy the pointwise next-round correction condition (**Eq. 6**) for every action pair $(s_X, s_Y)$ because $u_X(s_X, s_Y)$ depends on $s_Y$. To see this more explicitly, suppose there exists $\psi : S_X \to \mathbb{R}$ such that **Eq. 6** holds with $\alpha = 1$ and $\beta = \gamma = 0$. By scaling, we may assume $\psi(C) = 0$. The pointwise next-round correction condition then requires:

$$-1 = -\lambda \sigma_X[C, C](D) \psi(D) - (1 - \lambda) \sigma_X^0(D) \psi(D); \tag{10a}$$
$$-2 = -\lambda \sigma_X[C, D](D) \psi(D) - (1 - \lambda) \sigma_X^0(D) \psi(D); \tag{10b}$$
$$+1 = \psi(D) - \lambda \sigma_X[D, C](D) \psi(D) - (1 - \lambda) \sigma_X^0(D) \psi(D); \tag{10c}$$
$$+2 = \psi(D) - \lambda \sigma_X[D, D](D) \psi(D) - (1 - \lambda) \sigma_X^0(D) \psi(D). \tag{10d}$$

This system has no solution for $\psi(D)$, which shows that the pointwise next-round correction condition (**Eq. 6**) is sufficient but not necessary for enforcing linear payoff relationships.

## 2.7 The generalized next-round correction condition

We now extend the pointwise condition to reactive learning strategies. Suppose $\left(\sigma_X^0, \sigma_X[s_X, s_Y]\right)$ satisfies the pointwise next-round correction condition (**Eq. 6**) for some enforcement potential $\psi : S_X \to \mathbb{R}$. For $\tau_X \in \Delta(S_X)$, define the map $\Psi(\tau_X) := \mathbb{E}_{s_X \sim \tau_X}[\psi(s_X)]$ (which we also refer to

as an enforcement potential), and recall the induced reactive learning strategy, $\sigma_X^* [\tau_X, s_Y] (\cdot) :=$
$\mathbb{E}_{s_X \sim \tau_X} [\sigma_X [s_X, s_Y] (\cdot)]$.

Let $\mathcal{M}_X \subseteq \Delta(S_X)$ be the reachable set of mixed actions,

$$\mathcal{M}_X := \bigcap \left\{ \mathcal{M} \subseteq \Delta(S_X) \mid \sigma_X^0 \in \mathcal{M} \text{ and } \sigma_X^* [\tau_X, s_Y] \in \mathcal{M} \text{ for all } \tau_X \in \mathcal{M} \text{ and } s_Y \in S_Y \right\}. \quad (11)$$

That is, $\mathcal{M}_X$ is the smallest subset of $\Delta(S_X)$ containing the initial action $\sigma_X^0$ and closed under the response map $\tau_X \mapsto \sigma_X^* [\tau_X, s_Y]$ for every $s_Y \in S_Y$.

Let $\varphi(\tau_X, s_Y) := \mathbb{E}_{s_X \sim \tau_X} [\varphi(s_X, s_Y)]$ be the linear extension of $\varphi$ to mixed actions in the first coordinate. Taking expectations of **Eq. 6** with respect to $\tau_X \in \mathcal{M}_X$ yields the generalized next-round correction condition,

$$\varphi(\tau_X, s_Y) = \Psi(\tau_X) - \lambda \Psi(\sigma_X^* [\tau_X, s_Y]) - (1 - \lambda) \Psi(\sigma_X^0) \quad (12)$$

for every $\tau_X \in \mathcal{M}_X$ and $s_Y \in S_Y$, where $\varphi(s_X, s_Y) = \alpha u_X (s_X, s_Y) + \beta u_Y (s_X, s_Y) + \gamma$. Unlike the pointwise condition, we will show that the generalized next-round correction condition is both necessary and sufficient for $(\sigma_X^0, \sigma_X^*)$ to be autocratic, for general functions $\varphi$. We first need a technical lemma, which shows that autocratic strategies can be characterized by enforcement against deterministic sequences:

**Lemma 1.** A behavioral strategy $\sigma_X : \mathcal{H} \rightarrow \Delta(S_X)$ is $(\varphi, \lambda)$-autocratic if and only if it enforces $\varphi \equiv 0$ with discount factor $\lambda$ against all exogenous (deterministic) sequences $\{s_Y^t\}_{t=0}^\infty \subseteq S_Y$.

*Proof.* The "only if" direction follows immediately since exogenous sequences are behavioral strategies. For the converse, suppose $\sigma_X$ enforces $\varphi \equiv 0$ against all exogenous sequences. For any behavioral strategy $\sigma_Y : \mathcal{H} \rightarrow \Delta(S_Y)$, the tower property of total expectation gives

$$
\begin{aligned}
\mathbb{E}_{\sigma_X, \sigma_Y} \left[ \sum_{t=0}^\infty \lambda^t \varphi(s_X^t, s_Y^t) \right] &= \mathbb{E}_{\sigma_X, \sigma_Y} \left[ \mathbb{E}_{\sigma_X, \sigma_Y} \left[ \sum_{t=0}^\infty \lambda^t \varphi(s_X^t, s_Y^t) \mid (s_Y^t)_{t=0}^\infty \right] \right] \\
&= \mathbb{E}_{\sigma_X, \sigma_Y} \left[ \mathbb{E}_{\sigma_X, (s_Y^t)_{t=0}^\infty} \left[ \sum_{t=0}^\infty \lambda^t \varphi(s_X^t, s_Y^t) \mid (s_Y^t)_{t=0}^\infty \right] \right] \\
&= \mathbb{E}_{\sigma_X, \sigma_Y} [0] \\
&= 0, \quad (13)
\end{aligned}
$$

where the third equality comes from the hypothesis. $\square$

The generalized next-round correction condition provides a complete characterization of autocratic reactive learning strategies. Since **Eq. 6** implies **Eq. 12** by linearity of expectation, the following strengthens the main result of McAvoy and Hauert [10].

**Proposition 1.** Consider a function $\varphi : S_X \times S_Y \rightarrow \mathbb{R}$ and fix $\lambda \in [0, 1)$. If the generalized next-round correction condition holds for the reactive learning strategy $(\sigma_X^0, \sigma_X^*)$, then $(\sigma_X^0, \sigma_X^*)$ is a $(\varphi, \lambda)$-autocratic strategy.

*Proof.* By Lemma 1, it suffices to show that $(\sigma_X^0, \sigma_X^*)$ enforces $\varphi \equiv 0$ with discount factor $\lambda$ against all exogenous sequences. Consider an arbitrary sequence $\{s_Y^t\}_{t=0}^\infty \subseteq S_Y$. The strategy $(\sigma_X^0, \sigma_X^*)$

and this sequence generate a chain of mixed actions $\{\tau_X^t\}_{t=0}^{\infty}$ with $\tau_X^0 = \sigma_X^0$ and $\tau_X^{t+1} = \sigma_X^* \left[\tau_X^t, s_Y^t\right]$ for all $t \geqslant 0$. Let $\mu_X^t \in \Delta\left((\mathcal{M}_X \times S_Y)^t\right)$ be the probability distributions defined inductively by

$$\mu_X^0\left[\varnothing\right] := \sigma_X^*\left[\varnothing\right] \times \sigma_Y\left[\varnothing\right] = \sigma_X^0 \times \delta_{s_Y^0}; \tag{14a}$$

$$\mu_X^{t+1}\left(\left(\tau_X^0, s_Y^0\right), \ldots, \left(\tau_X^t, s_Y^t\right)\right) := \mu_X^t\left(\left(\tau_X^0, s_Y^0\right), \ldots, \left(\tau_X^{t-1}, s_Y^{t-1}\right)\right)$$
$$\times \delta_{\sigma_X^*\left[\tau_X^{t-1}, s_Y^{t-1}\right], \tau_X^t} \times \delta_{\sigma_Y\left[\tau_X^{t-1}, s_Y^{t-1}\right], s_Y^t}. \tag{14b}$$

From this we obtain the marginalized distributions $\{\nu_X^t\}_{t=0}^{\infty}$ over $\mathcal{M}_X \times S_Y$ defined by

$$\nu_X^t\left(A \times B\right) := \mu_X^t\left((\mathcal{M}_X \times S_Y)^{t-1} \times (A \times B)\right). \tag{15}$$

The generalized next-round correction condition (**Eq. 12**) gives

$$\mathbb{E}_{(\tau_X, s_Y) \sim \nu_X^t}\left[\varphi\left(\tau_X, s_Y\right)\right] = \mathbb{E}_{(\tau_X, s_Y) \sim \nu_X^t}\left[\Psi\left(\tau_X\right)\right] - \lambda \mathbb{E}_{(\tau_X, s_Y) \sim \nu_X^t}\left[\Psi\left(\sigma_X^*\left[\tau_X, s_Y\right]\right)\right] - (1 - \lambda)\Psi\left(\sigma_X^0\right). \tag{16}$$

Due to the deterministic nature of the exogenous sequence, by induction we find $\nu_X^{t+1}\left(\tau_X, s_Y\right) = \delta_{\left\{\tau_X = \tau_X^t, s_Y = s_Y^t\right\}}$ for every $(\tau_X, s_Y) \in \mathcal{M}_X \times S_Y$, where $\tau_X^t := \sigma_X^*\left[\tau_X^{t-1}, s_Y^{t-1}\right]$ for every $t \geqslant 1$. Thus,

$$\mathbb{E}_{(\tau_X, s_Y) \sim \nu_X^t}\left[\Psi\left(\sigma_X^*\left[\tau_X, s_Y\right]\right)\right] = \mathbb{E}_{(\tau_X, s_Y) \sim \nu_X^{t+1}}\left[\Psi\left(\tau_X\right)\right] \tag{17}$$

for every $t \geqslant 0$. From **Eq. 16** and **Eq. 17** we deduce the recursion

$$\mathbb{E}_{(\tau_X, s_Y) \sim \nu_X^t}\left[\varphi\left(\tau_X, s_Y\right)\right] = \mathbb{E}_{(\tau_X, s_Y) \sim \nu_X^t}\left[\Psi\left(\tau_X\right)\right] - \lambda \mathbb{E}_{(\tau_X, s_Y) \sim \nu_X^{t+1}}\left[\Psi\left(\tau_X\right)\right] - (1 - \lambda)\Psi\left(\sigma_X^0\right), \tag{18}$$

which gives a telescoping sum. The result follows by multiplying **Eq. 18** by $\lambda^t$ and summing over $t \geqslant 0$. $\qquad\square$

We next prove that the generalized next-round correction condition is also a necessary condition, and we consider the uniqueness of the enforcement potential $\Psi$.

**Proposition 2.** If $\left(\sigma_X^0, \sigma_X^*\right)$ is a reactive learning strategy that is $(\varphi, \lambda)$-autocratic, then there exists $\mathcal{M}_X \subseteq \Delta\left(S_X\right)$ and an enforcement potential $\Psi : \mathcal{M}_X \to \mathbb{R}$ such that

(a) $\lim_{t \to \infty} \lambda^t \Psi\left(\tau_X^t\right) = 0$ for every chain, $\{\tau_X^t\}_{t=0}^{\infty} \subseteq \mathcal{M}_X$, derived from $\left(\sigma_X^0, \sigma_X^*\right)$;

(b) The generalized next-round correction condition holds for all $\tau_X \in \mathcal{M}_X$ and $s_Y \in S_Y$.

(c) if $\widetilde{\mathcal{M}}_X \subseteq \Delta\left(S_X\right)$ and $\widetilde{\Psi} : \widetilde{\mathcal{M}}_X \to \mathbb{R}$ also satisfy $\lim_{t \to \infty} \lambda^t \widetilde{\Psi}\left(\widetilde{\tau}_X^t\right) = 0$ for chains $\{\widetilde{\tau}_X^t\}_{t=0}^{\infty} \subseteq \widetilde{\mathcal{M}}_X$, and **Eq. 12** holds for all $\tau_X \in \widetilde{\mathcal{M}}_X$ and $s_Y \in S_Y$, then $\mathcal{M}_X \subseteq \widetilde{\mathcal{M}}_X$ and $\widetilde{\Psi}|_{\mathcal{M}_X} - \Psi$ is constant on $\mathcal{M}_X$.

*Proof.* Let $\mathcal{M}_X$ be the set of all truncated chains of mixed actions derived from $\sigma_X^*$. In particular, $\mathcal{M}_X$ is the set of all actions $\tau_X \in \Delta\left(S_X\right)$ for which there exist sequences $\{\tau_X^0, \ldots, \tau_X^T\} \subseteq \Delta\left(S_X\right)$

and $\left\{s_Y^0,\ldots,s_Y^{T-1}\right\} \subseteq S_Y$ with $\tau_X^0 = \sigma_X^0$, $\tau_X^T = \tau_X$, and $\tau_X^t = \sigma_X^* \left[\tau_X^{t-1}, s_Y^{t-1}\right]$ for $1 \leqslant t \leqslant T$. For $\tau_X \in \mathcal{M}_X$, with sequences $\left\{\tau_X^0,\ldots,\tau_X^T\right\} \subseteq \Delta\left(S_X\right)$ and $\left\{s_Y^0,\ldots,s_Y^{T-1}\right\} \subseteq S_Y$ chosen as above, let

$$\Psi\left(\tau_X\right) := -\lambda^{-T} \sum_{t=0}^{T-1} \lambda^t \varphi\left(\tau_X^t, s_Y^t\right). \tag{19}$$

(We define $\Psi\left(\sigma_X^0\right) := 0$.) To see that **Eq. 19** gives a well-defined value of $\Psi$ at $\tau_X$, suppose that $\left\{\widetilde{\tau}_X^0,\ldots,\widetilde{\tau}_X^{\widetilde{T}}\right\} \subseteq \Delta\left(S_X\right)$ and $\left\{\widetilde{s}_Y^0,\ldots,\widetilde{s}_Y^{\widetilde{T}-1}\right\} \subseteq S_Y$ are also sequences with $\widetilde{\tau}_X^0 = \sigma_X^0$, $\widetilde{\tau}_X^{\widetilde{T}} = \tau_X$, and $\widetilde{\tau}_X^t = \sigma_X^* \left[\widetilde{\tau}_X^{t-1}, \widetilde{s}_Y^{t-1}\right]$ for $1 \leqslant t \leqslant \widetilde{T}$. For any $s_Y^* \in S_Y$, we can extend the sequences $\left\{\tau_X^t\right\}_{t=0}^T$ and $\left\{\widetilde{\tau}_X^t\right\}_{t=0}^{\widetilde{T}}$ as follows: for each $t \geqslant 1$, let $\tau_X^{T+t} := \sigma_X^* \left[\tau_X^{T+t-1}, s_Y^*\right]$ and $\widetilde{\tau}_X^{\widetilde{T}+t} := \sigma_X^* \left[\widetilde{\tau}_X^{\widetilde{T}+t-1}, s_Y^*\right]$. Since $\tau_X^T = \widetilde{\tau}_X^{\widetilde{T}} = \tau_X$, we have $\tau_X^{T+t} = \widetilde{\tau}_X^{\widetilde{T}+t}$ for every $t \geqslant 0$. Therefore, by the hypothesis,

$$\begin{aligned}
0 &= \sum_{t=0}^{T-1} \lambda^t \varphi\left(\tau_X^t, s_Y^t\right) + \sum_{t=T}^{\infty} \lambda^t \varphi\left(\tau_X^t, s_Y^*\right) \\
&= \sum_{t=0}^{\widetilde{T}-1} \lambda^t \varphi\left(\widetilde{\tau}_X^t, \widetilde{s}_Y^t\right) + \sum_{t=\widetilde{T}}^{\infty} \lambda^t \varphi\left(\widetilde{\tau}_X^t, s_Y^*\right) \\
&= \sum_{t=0}^{\widetilde{T}-1} \lambda^t \varphi\left(\widetilde{\tau}_X^t, \widetilde{s}_Y^t\right) + \lambda^{\widetilde{T}-T} \sum_{t=T}^{\infty} \lambda^t \varphi\left(\tau_X^t, s_Y^*\right).
\end{aligned} \tag{20}$$

It follows immediately that $\Psi$ is well-defined since, from **Eq. 20**, we obtain

$$-\lambda^{-T} \sum_{t=0}^{T-1} \lambda^t \varphi\left(\tau_X^t, s_Y^t\right) = -\lambda^{-\widetilde{T}} \sum_{t=0}^{\widetilde{T}-1} \lambda^t \varphi\left(\widetilde{\tau}_X^t, \widetilde{s}_Y^t\right). \tag{21}$$

That $\lim_{t \to \infty} \lambda^t \Psi\left(\tau_X^t\right) = 0$ for every chain of mixed actions $\left\{\tau_X^t\right\}_{t=0}^{\infty}$ derived from $\left(\sigma_X^0, \sigma_X^*\right)$ similarly follows from the fact that $\left(\sigma_X^0, \sigma_X^*\right)$ is $(\varphi, \lambda)$-autocratic, which establishes part *(a)*.

For part *(b)*, consider $\tau_X \in \mathcal{M}_X$. We may assume, without a loss of generality, that $\tau_X$ is not the initial action $\sigma_X^0$, since $\Psi\left(\sigma_X^0\right) = 0$. Then, as above, there exist $T \geqslant 1$ and sequences $\left\{\tau_X^0,\ldots,\tau_X^T\right\} \subseteq \Delta\left(S_X\right)$ and $\left\{s_Y^0,\ldots,s_Y^{T-1}\right\} \subseteq S_Y$ with $\tau_X^0 = \sigma_X^0$, $\tau_X^T = \tau_X$, and $\tau_X^t = \sigma_X^* \left[\tau_X^{t-1}, s_Y^{t-1}\right]$ for $1 \leqslant t \leqslant T$. Let $s_Y \in S_Y$, and set $\tau_X^{T+1} := \sigma_X^* \left[\tau_X, s_Y\right]$. By the definition of $\Psi$, we have

$$\begin{aligned}
\Psi\left(\tau_X\right) - \lambda \Psi\left(\sigma_X^* \left[\tau_X, s_Y\right]\right) &= \Psi\left(\tau_X^T\right) - \lambda \Psi\left(\tau_X^{T+1}\right) \\
&= -\lambda^{-T} \sum_{t=0}^{T-1} \lambda^t \varphi\left(\tau_X^t, s_Y^t\right) + \lambda \lambda^{-(T+1)} \sum_{t=0}^{T} \lambda^t \varphi\left(\tau_X^t, s_Y^t\right) \\
&= \varphi\left(\tau_X, s_Y\right).
\end{aligned} \tag{22}$$

The generalized next-round correction condition follows.

Since $\mathcal{M}_X$ is defined as the set of all chains of mixed actions derived from $\left(\sigma_X^0, \sigma_X^*\right)$, it follows that any other such $\widetilde{\mathcal{M}}_X$ must contain $\mathcal{M}_X$. With $\widehat{\Psi} := \widetilde{\Psi}|_{\mathcal{M}_X} - \widetilde{\Psi}\left(\sigma_X^0\right)$, we have $\widehat{\Psi}\left(\sigma_X^0\right) = 0$ and

$$\widehat{\Psi}\left(\sigma_X\left[\tau_X, s_Y\right]\right) = -\lambda^{-1} \varphi\left(\tau_X, s_Y\right) + \lambda^{-1} \widehat{\Psi}\left(\tau_X\right) \tag{23}$$

for every $\tau_X \in \mathcal{M}_X$ and $s_Y \in S_Y$. Using this recurrence, it follows by induction that $\widehat{\Psi} = \Psi$, and thus $\widetilde{\Psi}|_{\mathcal{M}_X} - \Psi = \widetilde{\Psi}\left(\sigma_X^0\right)$ on $\mathcal{M}_X$, which completes the proof of part *(c)*. $\qquad\square$

11

# 3 Autocratic strategies with short memory

For most of this section, we assume $\lambda \in [0, 1)$; that is, the game terminates with positive probability $1 - \lambda > 0$ after each round. Only in Section 3.3, which covers the undiscounted, infinite-horizon case, do we consider $\lambda \to 1$.

A basic question for implementing autocratic strategies in practice is: how much memory is required? While Definition 3 allows strategies with arbitrary memory, we show that strategies with shorter memory suffice. In this section, we prove that every autocratic strategy, regardless of its complexity, can be replaced by a two-point reactive learning strategy, which is one that mixes between just two fixed distributions (mixed actions).

We proceed in three steps. First, we show that autocratic strategies need only condition on the opponent's history (Proposition 3). Second, we establish that such strategies correspond to reactive learning strategies with a right-invariance property, which allows longer action histories to be "rolled up" into information that can be carried from round to round. Finally, we construct explicit two-point strategies that enforce any enforceable payoff relationship (Theorem 1).

**Proposition 3.** Suppose that $\sigma_X : \mathcal{H} \to \Delta(S_X)$ is $(\varphi, \lambda)$-autocratic. Then, with the opponent-action history space $\mathcal{H}_Y := \bigcup_{t \geqslant 0} S_Y^t$ (where we interpret $S_Y^0 = \{\varnothing\}$ as the "empty" history), there exists an opponent-conditioned strategy, $\widetilde{\sigma}_X : \mathcal{H}_Y \to \Delta(S_X)$, that is also $(\varphi, \lambda)$-autocratic. In particular, it always suffices for $X$ to condition on only the observed history of $Y$ alone.

*Proof.* Suppose that $\sigma_X : \mathcal{H} \to \Delta(S_X)$ is an autocratic strategy for $X$. If the opponent plays an exogenous sequence of pure actions, $\{s_Y^t\}_{t=0}^{\infty}$, which we denote by $\sigma_Y$, then

$$
\begin{aligned}
0 &= \mathbb{E}_{\sigma_X, \sigma_Y} \left[ \sum_{t=0}^{\infty} \lambda^t \varphi\left(s_X^t, s_Y^t\right) \right] \\
&= \sum_{t=0}^{\infty} \lambda^t \mathbb{E}_{\sigma_X, \sigma_Y} \left[ \varphi\left(\sigma_X\left[\left(s_X^0, s_Y^0\right), \ldots, \left(s_X^{t-1}, s_Y^{t-1}\right)\right], s_Y^t\right) \right] \\
&= \sum_{t=0}^{\infty} \lambda^t \varphi\left(\widetilde{\sigma}_X\left[s_Y^0, \ldots, s_Y^{t-1}\right], s_Y^t\right),
\end{aligned}
\tag{24}
$$

where $\widetilde{\sigma}_X : \mathcal{H}_Y \to \Delta(S_X)$ is defined by $\widetilde{\sigma}_X^0 := \sigma_X^0$ and

$$
\widetilde{\sigma}_X\left[s_Y^0, \ldots, s_Y^{t-1}\right]\left(s_X^t\right) := \mathbb{E}_{\left(s_X^0, \ldots, s_X^{t-1}\right) \in S_X^t} \left[ \sigma_X\left[\left(s_X^0, s_Y^0\right), \ldots, \left(s_X^{t-1}, s_Y^{t-1}\right)\right]\left(s_X^t\right) \right].
\tag{25}
$$

We simply do not track dependence on $X$'s actions, which we can do because $X$'s mixed actions determine their own history and we are taking expectations in the evaluation of autocratic strategies. We note that this strategy, while defined using exogenous sequences of opponent (pure) actions, can be used against any opponent, and it remains autocratic by Lemma 1, as desired. $\qquad\square$

Proposition 3 suggests an interesting connection between autocratic strategies and reactive learning strategies. Consider the map taking a reactive learning strategy, $(\sigma_X^0, \sigma_X^*)$, to a behavioral strategy, $\mathcal{U}(\sigma_X^0, \sigma_X^*) : \mathcal{H}_Y \to \Delta(S_X)$, with $\mathcal{U}(\sigma_X^0, \sigma_X^*)[\varnothing] = \sigma_X^0$ and $\mathcal{U}(\sigma_X^0, \sigma_X^*)[s_Y^0, \ldots, s_Y^t] = \tau_X^{t+1}$, where $\tau_X^{i+1} = \sigma_X^*[\tau_X^i, s_Y^i]$ for $i = 0, \ldots, t$ and $\tau_X^0 = \sigma_X^0$. Recall that $(\sigma_X^0, \sigma_X^*)$ is implicitly based on a subset $\mathcal{M}_X \subseteq \Delta(S_X)$, such that $\sigma_X^0 \in \mathcal{M}_X$ and $\sigma_X^*$ is a map from $\mathcal{M}_X \times S_Y$ to $\mathcal{M}_X$. The map $\mathcal{U}$ is not surjective because that would require that if two distinct histories prescribe the same

randomization for player $X$ at time $t$, then they must prescribe the same randomization for $X$ at all $T \geqslant t$ when $Y$ uses the same actions in both sequences thereafter. We can re-frame this problem slightly. We say that $\sigma_X$ is "right-invariant" if, whenever $\sigma_X[\alpha] = \sigma_X[\beta]$ for $\alpha, \beta \in \mathcal{H}_Y$, we have $\sigma_X[\alpha\gamma] = \sigma_X[\beta\gamma]$ for all $\gamma \in \mathcal{H}_Y$. One can check that if $\sigma_X : \mathcal{H}_Y \to \Delta(S_X)$ is in the image of $\mathcal{U}$, then $\sigma_X$ is right-invariant. Conversely, if $\sigma_X : \mathcal{H}_Y \to \Delta(S_X)$ is right-invariant, then $\sigma_X$ lifts to a canonical element in the domain of $\mathcal{U}$.

One of the goals of this paper is to show that every autocratic strategy, regardless of complexity, can be replaced by a reactive learning strategy, i.e., an element of the preimage of $\mathcal{U}$. In fact, we show that this can be done with a particularly simple and convenient reactive learning strategy.

## 3.1 Enforcing payoff constraints with short memory

The key to constructing simple autocratic strategies is identifying when two mixed actions can serve as the building blocks for enforcement. Lemma 2 ultimately provides necessary and sufficient conditions: two mixed actions $\tau_X^+$ and $\tau_X^-$ can enforce a constraint if their maximum and minimum payoffs satisfy certain inequalities that balance the discounting and allow for stable enforcement across all opponent responses.

Intuitively, these inequalities ensure that player $X$ can always find an appropriate mixture of $\tau_X^+$ and $\tau_X^-$ in response to any opponent action $s_Y$ such that the expected payoff remains on target. The "+" and "−" superscripts suggest their roles: $\tau_X^+$ typically yields higher values of $\varphi$ while $\tau_X^-$ yields lower values, and the strategy adjusts the mixture to maintain the target.

**Lemma 2.** Suppose that $\lambda \in [0, 1)$ and that $\tau_X^{\pm} \in \Delta(S_X)$ satisfy the inequalities

$$\min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) \geqslant (1 - \lambda) \max_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) + \lambda \max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right); \tag{26a}$$

$$\max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right) \leqslant \lambda \min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) + (1 - \lambda) \min_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right). \tag{26b}$$

Then, there exists a function $p^* : [0, 1] \times S_Y \to [0, 1]$ such that, with the reaction

$$\sigma_X^*\left[p\tau_X^+ + (1 - p)\tau_X^-, s_Y\right] = p^*[p, s_Y]\tau_X^+ + (1 - p^*[p, s_Y])\tau_X^-, \tag{27}$$

$X$ can enforce $\varphi \equiv K$ for any $K \in \left[\max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right), \min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right)\right]$.

The proof strategy is as follows. By Proposition 1, satisfying the generalized next-round correction condition (**Eq. 12**) with an appropriate enforcement potential $\Psi$ is sufficient to guarantee that a reactive learning strategy is autocratic. For a two-point strategy that mixes between $\tau_X^+$ and $\tau_X^-$, the map $\Psi$ is completely determined by just two values: $\psi\left(\tau_X^+\right)$ and $\psi\left(\tau_X^-\right)$. The strategy begins with an initial mixture $\sigma_X^0 = p_0\tau_X^+ + (1 - p_0)\tau_X^-$ for some $p_0 \in [0, 1]$, and responds to each opponent action $s_Y \in S_Y$ by playing $\sigma_X^*\left[p\tau_X^+ + (1 - p)\tau_X^-, s_Y\right] = p^*[p, s_Y]\tau_X^+ + (1 - p^*[p, s_Y])\tau_X^-$. The challenge is to choose $\psi\left(\tau_X^+\right)$, $\psi\left(\tau_X^-\right)$, and $p_0$ such that **Eq. 12** holds and all transition probabilities $p^*[p, s_Y]$ remain in $[0, 1]$ for every $p \in [0, 1]$ and $s_Y \in S_Y$. We construct these values explicitly below.

*Proof.* For two-point strategies, finding an enforcement potential $\psi$ satisfying the generalized next-round correction condition amounts to finding two constants, $\psi\left(\tau_X^+\right)$ and $\psi\left(\tau_X^-\right)$. Let

$$\psi\left(\tau_X^+\right) := \frac{1}{1 - \lambda} \min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right); \tag{28a}$$

13

$$\psi\left(\tau_X^-\right) := \frac{1}{1-\lambda} \max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right). \tag{28b}$$

From **Eq. 26**, we must have $\psi\left(\tau_X^+\right) \geqslant \psi\left(\tau_X^-\right)$. If $\psi\left(\tau_X^+\right) = \psi\left(\tau_X^-\right)$, then **Eq. 26** implies that $\varphi\left(\tau_X^+, s_Y\right) = \varphi\left(\tau_X^-, s_Y\right) = \psi\left(\tau_X^+\right)$ for every $s_Y \in S_Y$, in which case $\varphi \equiv \psi\left(\tau_X^+\right) = \psi\left(\tau_X^-\right)$ can be enforced by an unconditional strategy. Therefore, we may assume that $\psi\left(\tau_X^+\right) - \psi\left(\tau_X^-\right) > 0$.

Fix $K \in \left[\max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right), \min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right)\right]$ and consider the response function

$$p^*\left[p, s_Y\right] = \frac{K - p\varphi\left(\tau_X^+, s_Y\right) - (1-p)\,\varphi\left(\tau_X^-, s_Y\right)}{\lambda\left(\psi\left(\tau_X^+\right) - \psi\left(\tau_X^-\right)\right)} + \frac{p - (1-\lambda)\,p_0}{\lambda}. \tag{29}$$

To have $p^*\left[p, s_Y\right] \in [0,1]$ for all $p \in [0,1]$ and $s_Y \in S_Y$, necessary and sufficient conditions are

$$\frac{K - \min_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right)}{(1-\lambda)\left(\psi\left(\tau_X^+\right) - \psi\left(\tau_X^-\right)\right)} - \frac{\lambda}{1-\lambda} \leqslant p_0 \leqslant \frac{K - \max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right)}{(1-\lambda)\left(\psi\left(\tau_X^+\right) - \psi\left(\tau_X^-\right)\right)}; \tag{30a}$$

$$1 - \frac{\min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) - K}{(1-\lambda)\left(\psi\left(\tau_X^+\right) - \psi\left(\tau_X^-\right)\right)} \leqslant p_0 \leqslant \frac{1}{1-\lambda} - \frac{\max_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) - K}{(1-\lambda)\left(\psi\left(\tau_X^+\right) - \psi\left(\tau_X^-\right)\right)}. \tag{30b}$$

The fact that each of these two intervals is non-trivial follows from **Eq. 26**. Moreover, by the values given in **Eq. 28**, the two intervals intersect at a unique initial probability, which is

$$p_0 = \frac{K - \max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right)}{(1-\lambda)\left(\psi\left(\tau_X^+\right) - \psi\left(\tau_X^-\right)\right)} = \frac{K - \max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right)}{\min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) - \max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right)}. \tag{31}$$

With $p_0$ and $p^*$ well-defined and taking values in $[0,1]$, the generalized next-round correction condition is satisfied for $\varphi - K$, and thus $(p_0, p^*)$ allows $X$ to enforce $\varphi \equiv K$. $\qquad\square$

**Remark 1.** In the statement of Lemma 2, we implicitly assume that the mixing probability, $p$, can be determined by the value of $p\tau_X^+ + (1-p)\,\tau_X^-$. If $p, q \in [0,1]$ are mixing probabilities satisfying $p\tau_X^+ + (1-p)\,\tau_X^- = q\tau_X^+ + (1-q)\,\tau_X^-$, then $(p-q)\,\tau_X^+ = (p-q)\,\tau_X^-$. If $\tau_X^+ \neq \tau_X^-$, then $p = q$. If $\tau_X^+ = \tau_X^-$, then the Lemma holds trivially.

We are especially interested in enforcing relationships of the form $\varphi \equiv 0$, and there is no loss of generality in setting $K = 0$ since we can absorb this constant into $\varphi$, if necessary.

The following result is an immediate consequence of Lemma 2:

**Corollary 1.** If $\max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right) \leqslant 0 \leqslant \min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right)$ and the hypotheses of Lemma 2 hold for $\tau_X^+$ and $\tau_X^-$, then $X$ can enforce $\varphi \equiv 0$ (using a two-point reactive learning strategy).

Motivated by this result and our focus on enforcing $\varphi \equiv 0$, we define the sets

$$\Phi_X^+ := \left\{\tau_X \in \Delta\left(S_X\right) \mid \min_{s_Y \in S_Y} \varphi\left(\tau_X, s_Y\right) \geqslant 0\right\}; \tag{32a}$$

$$\Phi_X^- := \left\{\tau_X \in \Delta\left(S_X\right) \mid \max_{s_Y \in S_Y} \varphi\left(\tau_X, s_Y\right) \leqslant 0\right\}. \tag{32b}$$

In general, either or both of these sets can be empty.

We are now in a position to state and prove our main theoretical result:

**Theorem 1.** Suppose that $\sigma_X : \mathcal{H} \to \Delta(S_X)$ is a $(\varphi, \lambda)$-autocratic strategy of arbitrary memory. Then, there exists a two-point reactive learning strategy that is also $(\varphi, \lambda)$-autocratic.

*Proof.* If $\lambda = 0$, then conditioning is irrelevant and only the initial mixed action matters, so the result is trivial. Therefore, we assume that $\lambda > 0$ going forward. By Proposition 3, we may assume that $\sigma_X$ is a map from $\mathcal{H}_Y = \bigcup_{t \geqslant 0} S_Y^t$ to $\Delta(S_X)$. Consider the map defined by

$$\Theta : \mathcal{H}_Y \longrightarrow \mathbb{R}$$

$$: \left( s_Y^0, \ldots, s_Y^{T-1} \right) \longmapsto -\lambda^{-T} \sum_{t=0}^{T-1} \lambda^t \varphi \left( \sigma_X \left[ s_Y^0, \ldots, s_Y^{t-1} \right], s_Y^t \right), \tag{33}$$

where $\Theta(\varnothing) := 0$. From the definition of $\Theta$, we see that for all $h \in \mathcal{H}_Y$ and $s_Y \in S_Y$,

$$\varphi(\sigma_X[h], s_Y) = \Theta(h) - \lambda \Theta(h, s_Y). \tag{34}$$

From this equation, we also see that for every $h \in \mathcal{H}_Y$,

$$\max_{s_Y \in S_Y} \varphi(\sigma_X[h], s_Y) + \lambda \inf_{h \in \mathcal{H}_Y \setminus \{\varnothing\}} \Theta(h)$$

$$\leqslant \Theta(h) \leqslant \min_{s_Y \in S_Y} \varphi(\sigma_X[h], s_Y) + \lambda \sup_{h \in \mathcal{H}_Y \setminus \{\varnothing\}} \Theta(h). \tag{35}$$

To see that these bounds are finite, we note that for any $h \in \mathcal{H}_Y$ and $s_Y \in S_Y$,

$$\Theta(h) = \sum_{t=T}^{\infty} \lambda^{t-T} \varphi \left( \sigma_X \left[ h, \underbrace{s_Y, \ldots, s_Y}_{t - T \text{ times}} \right], s_Y \right), \tag{36}$$

since $\sigma_X$ is $(\varphi, \lambda)$-autocratic. This equation gives $\|\Theta\|_\infty \leqslant \frac{1}{1-\lambda} \|\varphi\|_\infty < \infty$ since $\varphi$ is bounded.

Fix two sequences of histories, $\{h_n^+\}_{n=0}^{\infty}$ and $\{h_n^-\}_{n=0}^{\infty}$, such that $\{\Theta(h_n^+)\}_{n=0}^{\infty}$ converges monotonically to $\sup_{h \in \mathcal{H}_Y \setminus \{\varnothing\}} \Theta(h)$ and $\{\Theta(h_n^-)\}_{n=0}^{\infty}$ converges monotonically to $\inf_{h \in \mathcal{H}_Y \setminus \{\varnothing\}} \Theta(h)$. For every fixed $\varepsilon > 0$, there must then exist $N_\varepsilon \geqslant 0$ such that, whenever $n \geqslant N_\varepsilon$,

$$\Theta(h_n^+) > \sup_{h \in \mathcal{H}_Y \setminus \{\varnothing\}} \Theta(h) - \frac{1}{2\lambda} \varepsilon; \tag{37a}$$

$$\Theta(h_n^-) < \inf_{h \in \mathcal{H}_Y \setminus \{\varnothing\}} \Theta(h) + \frac{1}{2\lambda} \varepsilon. \tag{37b}$$

We note that if $\sup_{h \in \mathcal{H}_Y \setminus \{\varnothing\}} \Theta(h) < 0$, then **Eq. 37a** holds for $h_n^+ = \varnothing$ for all $n \geqslant 0$. Similarly, if $\inf_{h \in \mathcal{H}_Y \setminus \{\varnothing\}} \Theta(h) > 0$, then **Eq. 37b** holds when $h_n^- = \varnothing$ for all $n \geqslant 0$. Consider now the pair,

$$\left( \tau_{X,n}^+, \tau_{X,n}^- \right) := \begin{cases} (\sigma_X[\varnothing], \sigma_X[h_n^-]) & \sup_{h \in \mathcal{H}_Y \setminus \{\varnothing\}} \Theta(h) < 0, \\[2mm] (\sigma_X[h_n^+], \sigma_X[\varnothing]) & \inf_{h \in \mathcal{H}_Y \setminus \{\varnothing\}} \Theta(h) > 0, \\[2mm] (\sigma_X[h_n^+], \sigma_X[h_n^-]) & \text{otherwise.} \end{cases} \tag{38}$$

15

Since $\Delta\left(S_X\right)$ is compact, by passing to subsequences of histories if necessary, we may assume that $\left\{\tau_{X,n}^+\right\}_{n=0}^\infty$ and $\left\{\tau_{X,n}^-\right\}_{n=0}^\infty$ are convergent sequences in $\Delta\left(S_X\right)$, with limits $\tau_X^+$ and $\tau_X^-$, respectively. To complete the proof, using Lemma 2 (Corollary 1), we show that **Eq. 26** holds and that $\tau_X^+ \in \Phi_X^+$ and $\tau_X^- \in \Phi_X^-$. To see that **Eq. 26** holds, we fix $\varepsilon > 0$ and let $n \geqslant N_\varepsilon$ and note that

$$
\begin{aligned}
(1-\lambda) & \left(\sup_{h\in\mathcal{H}_Y\backslash\{\varnothing\}} \Theta\left(h\right) - \inf_{h\in\mathcal{H}_Y\backslash\{\varnothing\}} \Theta\left(h\right)\right) \\
& < \Theta\left(h_n^+\right) - \Theta\left(h_n^-\right) + \frac{1}{\lambda}\varepsilon - \lambda\left(\sup_{h\in\mathcal{H}_Y\backslash\{\varnothing\}} \Theta\left(h\right) - \inf_{h\in\mathcal{H}_Y\backslash\{\varnothing\}} \Theta\left(h\right)\right) \quad \textbf{(Eq. 37)} \\
& \leqslant \min_{s_Y\in S_Y} \varphi\left(\tau_{X,n}^+,s_Y\right) - \max_{s_Y\in S_Y} \varphi\left(\tau_{X,n}^-,s_Y\right) + \frac{1}{\lambda}\varepsilon. \quad \textbf{(Eq. 35)}
\end{aligned} \tag{39}
$$

From this inequality, we can conclude that

$$
\begin{aligned}
(1-\lambda)\max_{s_Y\in S_Y}\varphi & \left(\tau_{X,n}^+,s_Y\right) + \lambda\max_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^-,s_Y\right) \\
& \leqslant (1-\lambda)\min_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^+,s_Y\right) + \lambda\max_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^-,s_Y\right) \\
& \quad + (1-\lambda)\lambda\left(\sup_{h\in\mathcal{H}_Y\backslash\{\varnothing\}}\Theta\left(h\right) - \inf_{h\in\mathcal{H}_Y\backslash\{\varnothing\}}\Theta\left(h\right)\right) \quad \textbf{(Eq. 35)} \\
& < (1-\lambda)\min_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^+,s_Y\right) + \lambda\max_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^-,s_Y\right) \\
& \quad + \lambda\left(\min_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^+,s_Y\right) - \max_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^-,s_Y\right) + \frac{1}{\lambda}\varepsilon\right) \quad \textbf{(Eq. 39)} \\
& = \min_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^+,s_Y\right) + \varepsilon;
\end{aligned} \tag{40a}
$$

$$
\begin{aligned}
\lambda\min_{s_Y\in S_Y}\varphi & \left(\tau_{X,n}^+,s_Y\right) + (1-\lambda)\min_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^-,s_Y\right) \\
& \geqslant \lambda\min_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^+,s_Y\right) + (1-\lambda)\max_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^-,s_Y\right) \\
& \quad - (1-\lambda)\lambda\left(\sup_{h\in\mathcal{H}_Y\backslash\{\varnothing\}}\Theta\left(h\right) - \inf_{h\in\mathcal{H}_Y\backslash\{\varnothing\}}\Theta\left(h\right)\right) \quad \textbf{(Eq. 35)} \\
& > \lambda\min_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^+,s_Y\right) + (1-\lambda)\max_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^-,s_Y\right) \\
& \quad - \lambda\left(\min_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^+,s_Y\right) - \max_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^-,s_Y\right) + \frac{1}{\lambda}\varepsilon\right) \quad \textbf{(Eq. 39)} \\
& = \max_{s_Y\in S_Y}\varphi\left(\tau_{X,n}^-,s_Y\right) - \varepsilon.
\end{aligned} \tag{40b}
$$

Since $\varepsilon > 0$ was arbitrary, it follows that **Eq. 26** holds for the pair $\left(\tau_X^+,\tau_X^-\right)$.

What remains to be shown is that $\tau_X^+ \in \Phi_X^+$ and $\tau_X^- \in \Phi_X^-$, which we establish in three cases:

*(i)* If $\sup_{h\in\mathcal{H}_Y\backslash\{\varnothing\}} \Theta\left(h\right) < 0$, then, by **Eq. 35**,

$$
\min_{s_Y\in S_Y} \varphi\left(\sigma_X\left[\varnothing\right],s_Y\right) \geqslant -\lambda\sup_{h\in\mathcal{H}_Y\backslash\{\varnothing\}}\Theta\left(h\right) > 0, \tag{41}
$$

and thus $\min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) > 0$. For $n \geqslant N_\varepsilon$, **Eq. 35** and **Eq. 37** give

$$\max_{s_Y \in S_Y} \varphi\left(\tau_{X,n}^-, s_Y\right) < (1-\lambda) \inf_{h \in \mathcal{H}_Y \setminus \{\varnothing\}} \Theta\left(h\right) + \frac{1}{2\lambda}\varepsilon, \tag{42}$$

and thus $\max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right) < 0$ in the limit, since $\varepsilon > 0$ was arbitrary.

(ii) If $\inf_{h \in \mathcal{H}_Y \setminus \{\varnothing\}} \Theta\left(h\right) > 0$, then, by **Eq. 35**,

$$\max_{s_Y \in S_Y} \varphi\left(\sigma_X\left[\varnothing\right], s_Y\right) \leqslant -\lambda \inf_{h \in \mathcal{H}_Y \setminus \{\varnothing\}} \Theta\left(h\right) < 0, \tag{43}$$

which gives $\max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right) < 0$. For $n \geqslant N_\varepsilon$, **Eq. 35** and **Eq. 37** give

$$\min_{s_Y \in S_Y} \varphi\left(\tau_{X,n}^+, s_Y\right) > (1-\lambda) \sup_{h \in \mathcal{H}_Y \setminus \{\varnothing\}} \Theta\left(h\right) - \frac{1}{2\lambda}\varepsilon, \tag{44}$$

so $\min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) > 0$ in the limit, since $\varepsilon > 0$ was arbitrary.

(iii) If neither (i) nor (ii) applies, then we deduce that $\tau_X^+ \in \Phi_X^+$ and $\tau_X^- \in \Phi_X^-$ by taking the limits of **Eq. 42** and **Eq. 44** as $n \to \infty$ and noting that $\varepsilon > 0$ was arbitrary.

By Corollary 1, we obtain a two-point reactive learning strategy that is $(\varphi, \lambda)$-autocratic. $\qquad \square$

Theorem 1 shows that if $\sigma_X$ is $(\varphi, \lambda)$-autocratic, then there exist $\tau_X^\pm \in \Delta\left(S_X\right)$ satisfying **Eq. 26** and $\max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right) \leqslant 0 \leqslant \min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right)$. If the latter condition is an equality, then **Eq. 26** implies that $\max_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) \leqslant \min_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right)$, and thus $\varphi\left(\tau_X^+, s_Y\right) = \varphi\left(\tau_X^-, s_Y\right) = 0$ for all $s_Y \in S_Y$. Therefore, in this case, **Eq. 26** holds with $\lambda = 0$, and the simple mixed action $\tau_X^+$, played in every round, is $(\varphi, \lambda)$-autocratic for all $\lambda \geqslant 0$. (This statement is also true for $\tau_X^-$, if it is distinct from $\tau_X^+$.) We refer to this unconditional play as a "trivial" autocratic strategy.

On the other hand, if $\max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right) < \min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right)$, then **Eq. 26** gives

$$\lambda \geqslant \max \left\{ \frac{\max_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) - \min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right)}{\max_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) - \max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right)}, \frac{\max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right) - \min_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right)}{\min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) - \min_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right)} \right\}$$

$$= 1 - \frac{\min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) - \max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right)}{\max \left\{ \begin{array}{l} \max_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) - \max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right), \\ \min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) - \min_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right) \end{array} \right\}}. \tag{45}$$

By Lemma 2 and Theorem 1, we immediately have the following result:

**Proposition 4.** There exists a $(\varphi, \lambda)$-autocratic strategy if and only if $\Phi_X^+, \Phi_X^- \neq \{\}$ and either (i) $\Phi_X^+ \cap \Phi_X^- \neq \{\}$ or (ii) $\Phi_X^+ \cap \Phi_X^- = \{\}$ and $\lambda \geqslant \lambda_{\min}$, where

$$\lambda_{\min} := 1 - \sup_{\left(\tau_X^+, \tau_X^-\right) \in \Phi_X^+ \times \Phi_X^-} \frac{\min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) - \max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right)}{\max \left\{ \begin{array}{l} \max_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) - \max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right), \\ \min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) - \min_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right) \end{array} \right\}}. \tag{46}$$

(Note that we use the notation $\{\}$ to denote the empty set to avoid confusion with the null history, $\varnothing$.)

**Remark 2.** The minimum discount factor $\lambda_{\min}$ has a natural interpretation: when $\lambda < \lambda_{\min}$, the game is too short on average (with expected duration $1/(1-\lambda)$ rounds) for the player to enforce the constraint. Enforcement requires the player to "correct" deviations from the target payoff relationship over time. If the game terminates too quickly, there is insufficient opportunity for these corrections to bring the expected payoff to the target. Note that $\lambda_{\min}$ depends on the "distance" between the extreme values of $\varphi$ at $\tau_X^+$ and $\tau_X^-$: larger swings between $\varphi\left(\tau_X^+, \cdot\right)$ and $\varphi\left(\tau_X^-, \cdot\right)$ require more patience (larger $\lambda$) to balance out through the discounting mechanism.

**Remark 3.** With some slight modifications in its proof, Theorem 1 can be extended to the setting where $S_X$ and $S_Y$ are compact sets over $\mathbb{R}^n$ and $\mathbb{R}^m$, respectively, and $\varphi$ is continuous, by invoking the Banach–Alaoglu theorem. More generally, with minimal changes, it can be shown that for $S_X$ and $S_Y$ of any structure, and for any bounded function $\varphi$, reduction to short-memory autocratic strategies is still plausible under arbitrarily small extensions of the game length. Rigorously speaking, for any $\varepsilon > 0$ such that $\lambda + \varepsilon < 1$, every $(\varphi, \lambda)$-autocratic strategy can be replaced by a two-point reactive learning $(\varphi, \lambda + \varepsilon)$-autocratic strategy.

## 3.2 Additive objective functions and reactive strategies

**Definition 5.** A function $\varphi : S_X \times S_Y \to \mathbb{R}$ is additive if there exist functions $\phi_X : S_X \to \mathbb{R}$ and $\phi_Y : S_Y \to \mathbb{R}$ with $\varphi\left(s_X, s_Y\right) = \phi_X\left(s_X\right) + \phi_Y\left(s_Y\right)$ for all $s_X \in S_X$ and $s_Y \in S_Y$.

For additive objective functions, we can strengthen our results considerably. Theorem 1 guarantees that enforcement can be achieved with two-point reactive learning strategies. However, the additive structure, $\varphi\left(s_X, s_Y\right) = \phi_X\left(s_X\right) + \phi_Y\left(s_Y\right)$, allows us to take $\psi = \phi_X$ in the next-round correction condition. Here, we show that two-point reactive learning strategies can be reduced to even simpler reactive strategies that condition solely on the opponent's last action, without tracking $X$'s own actions.

**Theorem 2.** If $\varphi \equiv 0$ is $\lambda$-enforceable by $X$ and $\varphi$ is additive, then $\varphi \equiv 0$ is $\lambda$-enforceable by $X$ using a reactive strategy $\sigma_X : S_Y \to \Delta\left(S_X\right)$.

*Proof.* Let $s_Y \in S_Y$ be fixed and suppose that the opponent plays $s_Y$ unconditionally in every round. Let $\left\{\sigma_X^t\left[s_Y\right]\right\}_{t=1}^{\infty} \subseteq \Delta\left(S_X\right)$ be a sequence for which $\sum_{t=0}^{\infty} \lambda^t \varphi\left(\sigma_X^t\left[s_Y\right], s_Y\right) = 0$ (making the dependence on $s_Y$ explicit). Consider the distribution, $\sigma_X\left[s_Y\right]$, defined by

$$\sigma_X\left[s_Y\right](\cdot) := (1-\lambda) \sum_{t=1}^{\infty} \lambda^{t-1} \sigma_X^t\left[s_Y\right](\cdot) \tag{47}$$

for each $s_X \in S_X$ and $s_Y \in S_Y$. The sequence of measures, $\left\{\widetilde{\sigma}_X^k\left[s_Y\right]\right\}_{k \geqslant 1}$, defined by

$$\widetilde{\sigma}_X^k\left[s_Y\right](\cdot) := (1-\lambda) \sum_{t=1}^{k} \lambda^{t-1} \sigma_X^t\left[s_Y\right](\cdot), \tag{48}$$

then converges in total variation to $\sigma_X\left[s_Y\right]$ because

$$\sigma_X\left[s_Y\right](\cdot) - \widetilde{\sigma}_X^k\left[s_Y\right](\cdot) = (1-\lambda) \sum_{t=k+1}^{\infty} \lambda^{t-1} \sigma_X^t\left[s_Y\right](\cdot)$$

18

$$\leqslant (1 - \lambda) \sum_{t=k+1}^{\infty} \lambda^{t-1}$$

$$= \lambda^k, \tag{49}$$

and thus $\left\| \sigma_X\left[s_Y\right] - \widetilde{\sigma}_X^k\left[s_Y\right] \right\|_{\mathrm{TV}} \leqslant \lambda^k \to 0$ as $k \to \infty$. It follows, then, that for every $s_Y \in S_Y$,

$$
\begin{aligned}
0 &= (1 - \lambda) \sum_{t=0}^{\infty} \lambda^t \varphi\left(\sigma_X^t\left[s_Y\right], s_Y\right) \\
&= (1 - \lambda) \, \varphi\left(\sigma_X^0, s_Y\right) + (1 - \lambda) \sum_{t=1}^{\infty} \lambda^t \varphi\left(\sigma_X^t\left[s_Y\right], s_Y\right) \\
&= (1 - \lambda) \, \varphi\left(\sigma_X^0, s_Y\right) + \lambda \lim_{t \to \infty} \varphi\left(\widetilde{\sigma}_X^t\left[s_Y\right], s_Y\right) \\
&= (1 - \lambda) \, \varphi\left(\sigma_X^0, s_Y\right) + \lambda \varphi\left(\sigma_X\left[s_Y\right], s_Y\right),
\end{aligned}
\tag{50}
$$

where the last equation is due to the boundedness of $\varphi$ and the fact that the sequence $\left\{\widetilde{\sigma}_X^k\left[s_Y\right]\right\}_{k \geqslant 1}$ converges to $\sigma_X\left[s_Y\right]$ in total variation. Since $\varphi\left(s_X, s_Y\right) = \phi_X\left(s_X\right) + \phi_Y\left(s_Y\right)$, we have

$$\phi_Y\left(s_Y\right) = -\lambda \mathbb{E}_{s_X \sim \sigma_X[s_Y]}\left[\phi_X\left(s_X\right)\right] - (1 - \lambda) \, \mathbb{E}_{s_X \sim \sigma_X^0}\left[\phi_X\left(s_X\right)\right], \tag{51}$$

which is the pointwise next-round correction condition for the reactive strategy $\left(\sigma_X^0, \sigma_X\left[s_Y\right]\right)$, giving the desired result. $\qquad \square$

**Remark 4.** Although our focus is on finite action spaces, this theorem readily extends to measurable spaces. Only minor additional justification is needed, such as the fact that for each $s_Y \in S_Y$, $\sigma_X\left[s_Y\right]$ is a probability measure (which follows from the Vitali-Hahn-Saks Theorem [see 30]).

**Remark 5.** An important feature of the reactive strategy constructed in Theorem 2 is that it uses the same initial mixed action, $\sigma_X^0$, as the original (possibly longer-memory) autocratic strategy. This is crucial for the proof, which relies on taking a weighted average of the mixed actions along the trajectory induced by the original strategy when the opponent plays the same action, $s_Y$, repeatedly. The initial action anchors this averaging process, ensuring that the mean converges to the desired enforcement property.

In fact, we can say slightly more about reactive strategies in this context:

**Lemma 3.** Suppose that $\lambda \in [0, 1)$ and that $\tau_X^{\pm} \in \Delta\left(S_X\right)$ satisfy **Eq. 26**. If $\varphi$ is additive, then there exists a function $p^* : S_Y \to [0, 1]$ such that, with the two-point reactive strategy

$$\sigma_X^*\left[s_Y\right] = p^*\left[s_Y\right] \tau_X^+ + \left(1 - p^*\left[s_Y\right]\right) \tau_X^-, \tag{52}$$

$X$ can enforce $\varphi \equiv K$ for any $K \in \left[\phi_X\left(\tau_X^-\right) + \max_{s_Y \in S_Y} \phi_Y\left(s_Y\right), \phi_X\left(\tau_X^+\right) + \min_{s_Y \in S_Y} \phi_Y\left(s_Y\right)\right]$.

*Proof.* Suppose that $\varphi\left(s_X, s_Y\right) = \phi_X\left(s_X\right) + \phi_Y\left(s_Y\right)$ for all $s_X \in S_X$ and $s_Y \in S_Y$, and let $\psi = \phi_X$. The generalized next-round correction condition is equivalent to

$$p^*\left[p, s_Y\right] = \frac{K - p\varphi\left(\tau_X^+, s_Y\right) - (1 - p)\,\varphi\left(\tau_X^-, s_Y\right)}{\lambda\left(\psi\left(\tau_X^+\right) - \psi\left(\tau_X^-\right)\right)} + \frac{p - (1 - \lambda)\,p_0}{\lambda}$$

$$= \frac{K - p\left(\phi_X\left(\tau_X^+\right) - \phi_X\left(\tau_X^-\right)\right) - \phi_X\left(\tau_X^-\right) - \phi_Y\left(s_Y\right)}{\lambda\left(\psi\left(\tau_X^+\right) - \psi\left(\tau_X^-\right)\right)} + \frac{p - (1-\lambda)\,p_0}{\lambda}$$

$$= \frac{K - \phi_X\left(\tau_X^-\right) - \phi_Y\left(s_Y\right)}{\lambda\left(\phi_X\left(\tau_X^+\right) - \phi_X\left(\tau_X^-\right)\right)} - \frac{(1-\lambda)\,p_0}{\lambda}, \tag{53}$$

which is independent of $p$. To ensure $p^*\left[s_Y\right] \in [0,1]$ for all $s_Y \in S_Y$ (dropping $p$), we require

$$\frac{K - \phi_X\left(\tau_X^-\right) - \min_{s_Y \in S_Y}\phi_Y\left(s_Y\right)}{(1-\lambda)\left(\phi_X\left(\tau_X^+\right) - \phi_X\left(\tau_X^-\right)\right)} - \frac{\lambda}{1-\lambda} \leqslant p_0 \leqslant \frac{K - \phi_X\left(\tau_X^-\right) - \max_{s_Y \in S_Y}\phi_Y\left(s_Y\right)}{(1-\lambda)\left(\phi_X\left(\tau_X^+\right) - \phi_X\left(\tau_X^-\right)\right)}. \tag{54}$$

Since $K \in \left[\phi_X\left(\tau_X^-\right) + \max_{s_Y \in S_Y}\phi_Y\left(s_Y\right), \phi_X\left(\tau_X^+\right) + \min_{s_Y \in S_Y}\phi_Y\left(s_Y\right)\right]$, we have the inequalities

$$\frac{K - \phi_X\left(\tau_X^-\right) - \min_{s_Y \in S_Y}\phi_Y\left(s_Y\right)}{(1-\lambda)\left(\phi_X\left(\tau_X^+\right) - \phi_X\left(\tau_X^-\right)\right)} - \frac{\lambda}{1-\lambda} \leqslant 1; \tag{55a}$$

$$\frac{K - \phi_X\left(\tau_X^-\right) - \max_{s_Y \in S_Y}\phi_Y\left(s_Y\right)}{(1-\lambda)\left(\phi_X\left(\tau_X^+\right) - \phi_X\left(\tau_X^-\right)\right)} \geqslant 0. \tag{55b}$$

By **Eq. 26**, we have $\max_{s_Y \in S_Y}\phi_Y\left(s_Y\right) - \min_{s_Y \in S_Y}\phi_Y\left(s_Y\right) \leqslant \lambda\left(\phi_X\left(\tau_X^+\right) - \phi_X\left(\tau_X^-\right)\right)$, which gives

$$\frac{K - \phi_X\left(\tau_X^-\right) - \min_{s_Y \in S_Y}\phi_Y\left(s_Y\right)}{(1-\lambda)\left(\phi_X\left(\tau_X^+\right) - \phi_X\left(\tau_X^-\right)\right)} - \frac{\lambda}{1-\lambda} \leqslant \frac{K - \phi_X\left(\tau_X^-\right) - \max_{s_Y \in S_Y}\phi_Y\left(s_Y\right)}{(1-\lambda)\left(\phi_X\left(\tau_X^+\right) - \phi_X\left(\tau_X^-\right)\right)}. \tag{56}$$

Therefore, the range of acceptable values of $p_0$ is the interval

$$\left[\max\left\{\frac{K - \phi_X\left(\tau_X^-\right) - \min_{s_Y \in S_Y}\phi_Y\left(s_Y\right)}{(1-\lambda)\left(\phi_X\left(\tau_X^+\right) - \phi_X\left(\tau_X^-\right)\right)} - \frac{\lambda}{1-\lambda}, 0\right\},\right.$$

$$\left.\min\left\{\frac{K - \phi_X\left(\tau_X^-\right) - \max_{s_Y \in S_Y}\phi_Y\left(s_Y\right)}{(1-\lambda)\left(\phi_X\left(\tau_X^+\right) - \phi_X\left(\tau_X^-\right)\right)}, 1\right\}\right]. \tag{57}$$

Any $p_0$ in this range translates into a valid response function, $p^*$. $\qquad\square$

**Corollary 2.** If $\varphi \equiv 0$ is $\lambda$-enforceable by $X$ and $\varphi$ is additive, then $\varphi \equiv 0$ is $\lambda$-enforceable by $X$ using a two-point reactive strategy.

We further obtain a simplification of Proposition 4:

**Proposition 5.** Suppose $\varphi$ is additive. Then, there exists a $(\varphi, \lambda)$-autocratic strategy if and only if $\Phi_X^+, \Phi_X^- \neq \{\}$ and either *(i)* $\phi_X$ is constant or *(ii)* $\phi_X$ is non-constant and $\lambda \geqslant \lambda_{\min}$, where

$$\lambda_{\min} = \frac{\max_{s_Y \in S_Y}\phi_Y\left(s_Y\right) - \min_{s_Y \in S_Y}\phi_Y\left(s_Y\right)}{\max_{s_X \in S_X}\phi_X\left(s_X\right) - \min_{s_X \in S_X}\phi_X\left(s_X\right)}. \tag{58}$$

## 3.3 Infinite-horizon (undiscounted) games

Up until this point, our focus has been on discounted games, which terminate in finitely many rounds with probability one ($\lambda < 1$). Although this case is the most realistic from a modeling

perspective, we know from classical results that some linear payoff relationships are enforceable in undiscounted, infinite-horizon games. An example is tit-for-tat in the repeated prisoner's dilemma, which ensures that $\pi_X = \pi_Y$ as $\lambda \to 1$ (and only in the infinite-horizon limit). In this section, we consider this limit more generally.

We assume that the undiscounted expectation of $\varphi$ is Cesàro summable, meaning the limit

$$\lim_{T \to \infty} \frac{1}{T+1} \sum_{t=0}^{T} \mathbb{E}_{\sigma_X, \sigma_Y} \left[ \varphi \left( s_X^t, s_Y^t \right) \right] \tag{59}$$

exists. Since both action spaces are finite, this limit exists whenever both behavioral strategies $\sigma_X$ and $\sigma_Y$ are finite-memory. We begin with the undiscounted analog of Theorem 1:

**Theorem 3.** Suppose that $\sigma_X : \mathcal{H} \to \Delta(S_X)$ is a $(\varphi, 1)$-autocratic strategy of arbitrary memory. Then, there exists a two-point reactive learning strategy that is also $(\varphi, 1)$-autocratic.

The proof requires two technical lemmas that establish the existence of appropriate two-point support and the undiscounted analog of the generalized next-round correction condition.

**Lemma 4.** If $\varphi \equiv 0$ is enforceable with $\lambda_{\min} = 1$, then there exist $\tau_X^{\pm} \in \Delta(S_X)$ such that

$$\max_{s_Y \in S_Y} \varphi \left( \tau_X^+, s_Y \right) > \min_{s_Y \in S_Y} \varphi \left( \tau_X^+, s_Y \right) = 0 = \max_{s_Y \in S_Y} \varphi \left( \tau_X^-, s_Y \right) > \min_{s_Y \in S_Y} \varphi \left( \tau_X^-, s_Y \right). \tag{60}$$

*Proof.* Suppose $\varphi \equiv 0$ is enforceable in expectation by some behavioral strategy $\sigma_X$ with initial action $\sigma_X^0$ and discount factor $\lambda = 1$. We first show that the sets $\Phi_X^+$ and $\Phi_X^-$ defined in **Eq. 32** are non-empty. Suppose by contradiction that this fails. Then either $\min_{s_Y \in S_Y} \varphi(\tau_X, s_Y) < 0$ for all $\tau_X \in \Delta(S_X)$, or $\max_{s_Y \in S_Y} \varphi(\tau_X, s_Y) > 0$ for all $\tau_X \in \Delta(S_X)$. Assume the first case holds (the second is analogous). Let $a := \max_{\tau_X \in \Delta(S_X)} \min_{s_Y \in S_Y} \varphi(\tau_X, s_Y) < 0$ and suppose that $Y$ plays $s_Y^t := \operatorname{argmin}_{s_Y \in S_Y} \varphi(\tau_X^t, s_Y)$ at each round $t$, where $\tau_X^t$ is the mixed action generated by $(\sigma_X^0, \sigma_X)$ at time $t$. This implies $(T+1)^{-1} \sum_{t=0}^{T} \varphi(\tau_X^t, s_Y^t) \leqslant a$ for every $T \geqslant 0$. Taking limits, we obtain $\lim_{T \to \infty} (T+1)^{-1} \sum_{t=0}^{T} \varphi(\tau_X^t, s_Y^t) \leqslant a < 0$, contradicting the assumption that $\sigma_X$ is $(\varphi, 1)$-autocratic. As a result, $\Phi_X^+, \Phi_X^- \neq \{\}$. Next, suppose there exist $\tau_X^+ \in \Phi_X^+$ and $\tau_X^- \in \Phi_X^-$ such that $\min_{s_Y \in S_Y} \varphi(\tau_X^+, s_Y) \geqslant 0 \geqslant \max_{s_Y \in S_Y} \varphi(\tau_X^-, s_Y)$, at least one of which is strict. Since no trivial autocratic strategies exist, by Proposition 4 there exists $\lambda^* < 1$ and a $(\varphi, \lambda^*)$-autocratic strategy, but this finding is a contradiction because $\lambda_{\min} = 1$. $\qquad \square$

**Lemma 5.** Suppose that $\Phi_X^+, \Phi_X^- \neq \{\}$. Then, there exist a set $\mathcal{M}_X \subseteq \Delta(S_X)$, a (two-point) response function $\sigma_X^* : \mathcal{M}_X \times S_Y \to \mathcal{M}_X$, and an enforcement potential $\Psi : \mathcal{M}_X \to \mathbb{R}$ such that

$$\varphi \left( \tau_X, s_Y \right) = \Psi \left( \tau_X \right) - \Psi \left( \sigma_X^* \left[ \tau_X, s_Y \right] \right) \tag{61}$$

for all $\tau_X \in \mathcal{M}_X$ and $s_Y \in S_Y$.

*Proof.* The proof parallels that of Lemma 2. Consider the response function $p^*$ from **Eq. 29** with $K = 0$ and $\lambda = 1$, let $\psi(\tau_X^-) = 0$, and choose $\psi(\tau_X^+)$ satisfying

$$\psi \left( \tau_X^+ \right) \geqslant \max \left\{ \max_{s_Y \in S_Y} \varphi \left( \tau_X^+, s_Y \right), - \min_{s_Y \in S_Y} \varphi \left( \tau_X^-, s_Y \right) \right\}. \tag{62}$$

Then, assuming no unconditional strategies exist (otherwise **Eq. 61** holds trivially), we deduce that $\psi(\tau_X^+) - \psi(\tau_X^-) > 0$, which ensures $\operatorname{Im}(p^*) \subseteq [0, 1]$. Defining $\mathcal{M}_X$ and $\Psi$ as in Lemma 2 (with $p_0$ now free in $[0, 1]$), we obtain **Eq. 61**. $\qquad \square$

**Eq. 61** serves as the undiscounted generalized next-round correction condition, representing the limiting case of **Eq. 12** as $\lambda \to 1$. We now prove Theorem 3.

*Proof of Theorem 3.* Assume, without loss of generality, that no trivial strategies exist. We consider two cases. In the first case, if there exists a $(\varphi, \lambda)$-autocratic strategy for some $\lambda \in [0,1)$, then by the proof of Theorem 1, the separation sets $\Phi_X^+$ and $\Phi_X^-$ are non-empty. In the second case, if no $(\varphi, \lambda)$-autocratic strategy exists for any $\lambda \in [0,1)$, then by Lemma 4 there exist $\tau_X^\pm \in \Delta(S_X)$ satisfying **Eq. 60**, which also implies $\Phi_X^+, \Phi_X^- \neq \{\}$. In either case, Lemma 5 ensures that **Eq. 61** holds. Following the proof of Proposition 1, for any behavioral strategy $\sigma_Y$, we derive a sequence of marginalized distributions $\{v_X^t\}_{t=0}^\infty$ over $\Delta(S_X) \times S_Y$. Telescoping the undiscounted generalized next-round correction condition yields

$$\frac{1}{T+1} \sum_{t=0}^T \mathbb{E}_{(\tau_X, s_Y) \sim v_X^t} [\varphi(\tau_X, s_Y)] = \frac{1}{T+1} \left( \mathbb{E}_{(\tau_X, s_Y) \sim v_X^0} [\Psi(\tau_X)] - \mathbb{E}_{(\tau_X, s_Y) \sim v_X^{T+1}} [\Psi(\tau_X)] \right) \quad (63)$$

for every $T \geqslant 0$. By boundedness of $\Psi$ and the dominated convergence theorem, the right-hand side of **Eq. 63** converges to 0 as $T \to \infty$. $\qquad\square$

**Remark 6.** It is important to highlight that, unlike the discounted case, the initial action $\sigma_X^0$ plays no part in the enforceability of $\varphi$; any "suboptimal" initial choice will eventually be corrected as the stage game repeats itself infinitely many times. It plays a major role, however, in the conditioning structure of the response function, $\sigma_X^*$. If $p^*[p, s_Y] \in \{0, 1\}$, then $p \in \{0, 1\}$ and $s_Y \in \left\{ \text{argmax}_{s_Y \in S_Y} \varphi(\tau_X^-, s_Y), \text{argmin}_{s_Y \in S_Y} \varphi(\tau_X^+, s_Y) \right\}$. Inductively, we deduce that if $p_0 \in (0, 1)$, then $\text{Im}(p^*) \subseteq (0, 1)$. This essentially means that the Markov chain generated by $(\sigma_X^0, \sigma_X^*)$ (and any behavioral strategy of $Y$) is ergodic. However, if $X$ desires to enforce $\varphi \equiv 0$ via a simple reactive learning strategy, they can do so by choosing $p_0 = 0$ or $p_0 = 1$.

By Lemmas 4, 5 and Theorem 3, we obtain a version of Propositions 1-2 and 4 in the case where a payoff constraint is only enforceable in the undiscounted setting.

**Proposition 6.** A function $\varphi : S_X \times S_Y \to \mathbb{R}$ is enforceable with $\lambda_{\min} = 1$ if and only if $\Phi_X^+, \Phi_X^- \neq \{\}$ and

$$\Phi_X^+ = \left\{ \tau_X \in \Delta(S_X) \mid \max_{s_Y \in S_Y} \varphi(\tau_X, s_Y) > \min_{s_Y \in S_Y} \varphi(\tau_X, s_Y) = 0 \right\}; \quad (64)$$

$$\Phi_X^- = \left\{ \tau_X \in \Delta(S_X) \mid \min_{s_Y \in S_Y} \varphi(\tau_X, s_Y) < \max_{s_Y \in S_Y} \varphi(\tau_X, s_Y) = 0 \right\}. \quad (65)$$

**Proposition 7.** A function $\varphi : S_X \times S_Y \to \mathbb{R}$ is enforceable with $\lambda_{\min} = 1$ if and only if there exist a set $\mathcal{M}_X \subseteq \Delta(S_X)$, an initial action $\sigma_X^0 \in \mathcal{M}_X$, a response function $\sigma_X^* : \mathcal{M}_X \times S_Y \to \mathcal{M}_X$, and an enforcement potential $\Psi : \mathcal{M}_X \to \mathbb{R}$ such that **Eq. 61** holds for all $\tau_X \in \mathcal{M}_X$ and $s_Y \in S_Y$, and the generalized next-round correction condition (**Eq. 12**) never holds for any $\lambda < 1$.

The following result extends [11, Proposition 6]. In simple terms, a player's control over the long-run payoff outcome is not altered if the game is arbitrarily extended.

**Proposition 8.** Fix $\varphi : S_X \times S_Y \to \mathbb{R}$ and $\lambda \in [0, 1]$. If there exists a $(\varphi, \lambda)$-autocratic strategy, then there exists a (two-point reactive learning) $(\varphi, \lambda^*)$-autocratic strategy for any $\lambda^* \in [\lambda, 1]$.

*Proof.* If $\lambda = 0$ or $1$, then we have nothing to show. We can thus assume that $\lambda \in (0,1)$. From Proposition 4, for any $1 > \lambda^* \geqslant \lambda$ there exists a two-point reactive learning $(\varphi, \lambda^*)$-autocratic strategy. Suppose $\lambda^* = 1$. As in the first case in the proof of Theorem 3, there exist mixed actions $\tau_X^+ \in \Phi_X^+, \tau_X^- \in \Phi_X^-$ such that $\max_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) > \min_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right)$. Therefore, by the same result, there exists a two-point reactive learning strategy that is $(\varphi, 1)$-autocratic. $\qquad\square$

# 4 Properties of payoff relationships and autocratic strategies

Having established that any enforceable payoff relationship can be implemented using two-point reactive learning strategies (Theorem 1), we now investigate the computational and structural properties of autocratic strategies. We show that verifying enforceability and computing optimal strategies can be accomplished in polynomial time using linear programming. We also establish convexity properties that reveal favorable geometric structure in the space of enforceable relations.

## 4.1 Enforceable relationships and the interval of enforceability

While the generalized next-round correction condition provides a characterization of autocratic strategies, verifying it directly can be challenging in practice, as it requires finding an appropriate map $\Psi$ on the reachable set, $\mathcal{M}_X$. In this section, we provide a simple, computationally tractable criterion for determining enforceability. The key observation is that enforceability depends only on whether an "interval of enforceability" contains zero, which can be checked by solving two saddle point problems: $\min_{\tau_X \in \Delta(S_X)} \max_{s_Y \in S_Y} \varphi\left(\tau_X, s_Y\right)$ and $\max_{\tau_X \in \Delta(S_X)} \min_{s_Y \in S_Y} \varphi\left(\tau_X, s_Y\right)$.

**Definition 6.** For a function $\varphi : S_X \times S_Y \to \mathbb{R}$, the interval of enforceability is

$$J\left(\varphi\right) := \left[\min_{\tau_X \in \Delta(S_X)} \max_{s_Y \in S_Y} \varphi\left(\tau_X, s_Y\right), \max_{\tau_X \in \Delta(S_X)} \min_{s_Y \in S_Y} \varphi\left(\tau_X, s_Y\right)\right]. \tag{66}$$

Of course, $J\left(\varphi\right) \neq \{\}$ if $\min_{\tau_X \in \Delta(S_X)} \max_{s_Y \in S_Y} \varphi\left(\tau_X, s_Y\right) \leqslant \max_{\tau_X \in \Delta(S_X)} \min_{s_Y \in S_Y} \varphi\left(\tau_X, s_Y\right)$.

The interval $J\left(\varphi\right)$ encodes geometric information about the separability of payoff values. Since the functions $\max_{s_Y \in S_Y} \varphi\left(\cdot, s_Y\right)$ and $\min_{s_Y \in S_Y} \varphi\left(\cdot, s_Y\right)$ are continuous over the compact set $\Delta\left(S_X\right)$, the extrema defining $J\left(\varphi\right)$ are always attained. The following characterization connects enforceability to the condition $0 \in J\left(\varphi\right)$, which has a natural interpretation: player $X$ can enforce $\varphi \equiv 0$ if and only if she can find mixed actions that "sandwich" zero between the worst and best values of $\varphi$ across opponent responses.

**Corollary 3.** Consider a function $\varphi : S_X \times S_Y \to \mathbb{R}$. The following are equivalent:

(a) $\varphi \equiv 0$ is enforceable;

(b) $\Phi_X^+, \Phi_X^- \neq \{\}$;

(c) $0 \in J\left(\varphi\right)$.

## 4.2 Computational tractability of enforceable relationships

Suppose that $S_X = \{U, D\}$, $S_Y = \{L, R\}$, and that

$$\varphi = \begin{array}{c} \\ U \\ D \end{array}\begin{pmatrix} \overset{L}{4} & \overset{R}{1} \\ -1 & 0 \end{pmatrix}. \tag{67}$$

Let $\tau_X^+$ be the mixed action that plays $U$ and $D$ uniformly at random, and let $\tau_X^-$ be the pure action $D$. Then, $\varphi(\tau_X^+, L) = 3/2$, $\varphi(\tau_X^+, R) = 1/2$, $\varphi(\tau_X^-, L) = -1$, and $\varphi(\tau_X^-, R) = 0$, from which we see that $\tau_X^+ \in \Phi_X^+$ and $\tau_X^- \in \Phi_X^-$, and the two inequalities of **Eq. 26** hold with $\lambda = 2/3$ (in fact, they are both equalities). However, if we set $\tau_X^+$ to be $U$ and $\tau_X^-$ to be $D$, then the first of these two inequalities fails to hold. In fact, the minimum $\lambda$ for which they hold when restricted to pure actions is $\lambda = 3/4$. Therefore, randomizing between two mixed actions, $\tau_X^+$ and $\tau_X^-$, might be preferable to randomizing between two pure actions, for impatient players.

However, under somewhat restrictive conditions, we can guarantee that $\lambda_{\min}$ is attained for pure actions, $\tau_X^\pm \in \Delta(S_X)$.

**Lemma 6.** Suppose $S_X = \{U, D\}$, $S_Y = \{L, R\}$, and let $\varphi : S_X \times S_Y \to \mathbb{R}$ satisfy $\varphi(U, L) \geqslant \varphi(U, R)$ and $\varphi(D, L) \geqslant \varphi(D, R)$. If $\varphi \equiv 0$ is enforceable, then $\lambda_{\min}$ is attained at pure actions.

*Proof.* Suppose that $\tau_X^+ \in \Phi_X^+$ and $\tau_X^- \in \Phi_X^-$ are any two mixed actions satisfying the inequalities of **Eq. 26** for some $\lambda \in [0, 1)$. Without a loss of generality, we may assume that $\delta_U \in \Phi_X^+$ and $\delta_D \in \Phi_X^-$ because $\varphi(\tau_X^+, L) \geqslant 0$ means that there must be a pure action $s_X^+ \in \{U, D\}$ in the support of $\tau_X^+$ such that $\varphi(s_X^+, L) \geqslant 0$, and similarly for $\tau_X^-$. Since $S_X$ has only two options, $\tau_X^+$ and $\tau_X^-$ can be represented by the probabilities of playing $U$, denoted $p$ and $q$, respectively. Since $\max_{s_Y \in S_Y} \varphi(\tau_X^-, s_Y)$ is a non-decreasing function of $q$, we see that **Eq. 26**a holds when $\tau_X^-$ is replaced by $\delta_D$. This first inequality, **Eq. 26**a, then says

$$\varphi(\tau_X^+, R) \geqslant (1 - \lambda)\varphi(\tau_X^+, L) + \lambda\varphi(D, L). \tag{68}$$

This inequality is linear in $p$, so it must also hold when $p = 0$ or when $p = 1$. If it holds when $p = 0$, then $\varphi(D, R) \geqslant \varphi(D, L)$. Since it must be true that $\varphi(D, L) = \varphi(D, R)$ or $\varphi(D, L) > \varphi(D, R)$, we see at once that the inequality holds for $p = 1$ and thus for the pair $(U, R)$. That **Eq. 26**b holds for the pair $(U, R)$ is analogous, this time increasing $p$ to 1 first and then $q$ to 0. $\square$

To see that Lemma 6 does not extend to larger action spaces for $X$, consider the function

$$\varphi = \begin{array}{c} \\ U \\ M \\ D \end{array}\begin{pmatrix} \overset{L}{-2} & \overset{R}{-\frac{2}{5}} \\ 2 & 5 \\ 1 & 2 \end{pmatrix}. \tag{69}$$

For pure action candidates, we must have $s_X^+ \in \{M, D\}$ and $s_X^- = U$. Suppose that $\lambda = 1/2$, and let $\tau_X^+ = (3/10)\delta_M + (7/10)\delta_D$ and $\tau_X^- = U$. We then calculate that

$$\frac{13}{10} = \min_{s_Y \in S_Y} \varphi(\tau_X^+, s_Y) \geqslant (1 - \lambda)\max_{s_Y \in S_Y} \varphi(\tau_X^+, s_Y) + \lambda\max_{s_Y \in S_Y} \varphi(\tau_X^-, s_Y) = \frac{25}{20}; \tag{70a}$$

$$-\frac{2}{5} = \max_{s_Y \in S_Y} \varphi \left( \tau_X^-, s_Y \right) \leqslant \lambda \min_{s_Y \in S_Y} \varphi \left( \tau_X^+, s_Y \right) + (1 - \lambda) \min_{s_Y \in S_Y} \varphi \left( \tau_X^-, s_Y \right) = -\frac{7}{20}. \tag{70b}$$

However, if we consider the pure-action candidates for $s_X^+$ and $s_X^-$ for $\lambda = 1/2$, then, with $s_X^+ = M$ and $s_X^- = U$, the first inequality is

$$2 = \min_{s_Y \in S_Y} \varphi \left( s_X^+, s_Y \right) \geqslant (1 - \lambda) \max_{s_Y \in S_Y} \varphi \left( s_X^+, s_Y \right) + \lambda \max_{s_Y \in S_Y} \varphi \left( s_X^-, s_Y \right) = \frac{23}{10}, \tag{71}$$

which does not hold. If $s_X^+ = D$ and $s_X^- = U$, then the second inequality is

$$-\frac{2}{5} = \max_{s_Y \in S_Y} \varphi \left( \tau_X^-, s_Y \right) \leqslant \lambda \min_{s_Y \in S_Y} \varphi \left( \tau_X^+, s_Y \right) + (1 - \lambda) \min_{s_Y \in S_Y} \varphi \left( \tau_X^-, s_Y \right) = -\frac{1}{2}, \tag{72}$$

which also does not hold.

These examples demonstrate that mixed actions can achieve better (smaller) values of $\lambda_{\min}$ than pure actions, even in small games. This motivates the need for efficient algorithms that search over the full space $\Delta(S_X)$ rather than just pure actions. We now show that despite this complexity, both the verification of enforceability and the computation of optimal autocratic strategies are tractable via linear programming.

To specify a two-point reactive learning strategy explicitly, we need only provide a small number of parameters. For a strategy $(\sigma_X^0, \sigma_X^*)$ defined by transition probabilities $(p_0, p^*)$ (as in Lemmas 2 and 5), we define its base to be the finite set of probabilities that fully characterize the strategy.

**Definition 7.** Consider a two-point reactive learning strategy $(\sigma_X^0, \sigma_X^*)$ defined by $(p_0, p^*)$ of either Lemma 2 or Lemma 5. The set $\{p_0\} \cup \{p^* [0, s_Y], p^* [1, s_Y]\}_{s_Y \in S_Y}$ is called a "base" of $(\sigma_X^0, \sigma_X^*)$.

The notion of a base arises from the fact that

$$p^* [q, s_Y] = q p^* [1, s_Y] + (1 - q) p^* [0, s_Y] \tag{73}$$

for all $q \in [0, 1]$ and $s_Y \in S_Y$.

**Proposition 9.** Consider a function $\varphi : S_X \times S_Y \to \mathbb{R}$. The problem of identifying whether $\varphi \equiv 0$ is enforceable can be solved in polynomial time. In addition, if $\varphi \equiv 0$ is enforceable, then the minimum discount factor $\lambda_{\min}$ and a base of some two-point reactive learning $(\varphi, \lambda_{\min})$-autocratic strategy can both be computed in polynomial time.

*Proof.* Let $S_X := \{s_{X,1}, \ldots, s_{X,m}\}$ and $S_Y := \{s_{Y,1}, \ldots, s_{Y,n}\}$. We can think of $\varphi$ as an $m \times n$ matrix with entries $\phi_{ij} = \varphi(s_{X,i}, s_{Y,j})$. Recall that $\varphi \equiv 0$ is enforceable if and only if $\Phi_X^+, \Phi_X^- \neq \{\}$. $\Phi_X^+ \neq \{\}$ is equivalent to $\phi^\mathsf{T} x \geqslant 0$ for $x \in \mathbb{R}_+^m$ and $\sum_{i=1}^m x_i = 1$, while $\Phi_X^- \neq \{\}$ is equivalent to $\phi^\mathsf{T} x \leqslant 0$ for $x \in \mathbb{R}_+^m$ and $\sum_{i=1}^m x_i = 1$, both of which can be solved in polynomial time. In addition, $\varphi \equiv 0$ is 0-enforceable if and only if the linear system $\phi^\mathsf{T} x = 0$ with $x \in \mathbb{R}_+^m$ and $\sum_{i=1}^m x_i = 1$ is feasible, which is also a computationally tractable task.

Suppose now that $\varphi \equiv 0$ is enforceable. If $\varphi$ is 0-enforceable, then we can compute the trivial mixed action in polynomial time as above. Therefore, we may assume that $\varphi \equiv 0$ is not 0-enforceable. As $\Phi_X^+, \Phi_X^- \neq \{\}$, the supremum problem in **Eq. 46** is well-defined and its value is attained. In fact, by standard max-min techniques and the Charnes-Cooper transformation, the supremum problem can be equivalently split into the following linear programming forms:

$$\text{maximize} \qquad z^+ - w^-$$

$$
\begin{aligned}
\text{subject to} \quad & z^+ \leqslant \textstyle\sum_{i=1}^m x_i^+ \phi_{ij} && j = 1, \ldots, n \\
& w^+ \geqslant \textstyle\sum_{i=1}^m x_i^+ \phi_{ij} && j = 1, \ldots, n \\
& z^- \leqslant \textstyle\sum_{i=1}^m x_i^- \phi_{ij} && j = 1, \ldots, n \\
& w^- \geqslant \textstyle\sum_{i=1}^m x_i^- \phi_{ij} && j = 1, \ldots, n \\
& x_i^\pm \geqslant 0 && i = 1, \ldots, m \\
& \textstyle\sum_{i=1}^m x_i^\pm = t \\
& w^+ - w^- = 1 \\
& z^+ - z^- \leqslant 1 \\
& z^+ \geqslant 0 \\
& w^- \leqslant 0 \\
& t \geqslant 0
\end{aligned}
$$

$\left(P_{\lambda_{\min}}^1\right)$

$$\text{maximize} \qquad z^+ - w^-$$

$$
\begin{aligned}
\text{subject to} \quad & z^+ \leqslant \textstyle\sum_{i=1}^m x_i^+ \phi_{ij} && j = 1, \ldots, n \\
& w^+ \geqslant \textstyle\sum_{i=1}^m x_i^+ \phi_{ij} && j = 1, \ldots, n \\
& z^- \leqslant \textstyle\sum_{i=1}^m x_i^- \phi_{ij} && j = 1, \ldots, n \\
& w^- \geqslant \textstyle\sum_{i=1}^m x_i^- \phi_{ij} && j = 1, \ldots, n \\
& x_i^\pm \geqslant 0 && i = 1, \ldots, m \\
& \textstyle\sum_{i=1}^m x_i^\pm = t \\
& w^+ - w^- \leqslant 1 \\
& z^+ - z^- = 1 \\
& z^+ \geqslant 0 \\
& w^- \leqslant 0 \\
& t \geqslant 0
\end{aligned}
$$

$\left(P_{\lambda_{\min}}^2\right)$

Note that as $\lambda_{\min} \in (0,1]$, we can always find an optimal solution for $\left(P_{\lambda_{\min}}^1\right)$ and $\left(P_{\lambda_{\min}}^1\right)$ with $t_1, t_2 > 0$. We can then set $\tau_X^{1,\pm} := x^{1,\pm}/t_1$ and $\tau_X^{2,\pm} := x^{2,\pm}/t_2$ to retrieve the desired mixed actions in $\Phi_X^+$ and $\Phi_X^-$, and we set $\lambda_{\min} = 1 - \max\left\{z^{1,+} - w^{1,-}, z^{2,+} - w^{2,-}\right\}$. The pair $\tau_X^{i,\pm}$ of actions that corresponds to the largest optimal value of the two linear programs will then lead to a two-point reactive learning $(\varphi, \lambda_{\min})$-autocratic strategy $\left(\sigma_X^0, \sigma_X^*\right)$, as seen from Theorem 1.

Finally, in order to identify a base of $\left(\sigma_X^0, \sigma_X^*\right)$, we need only compute at most $2n + 1$ values, as indicated by **Eq. 31**, **Eq. 57** and **Eq. 73**. $\qquad\square$

## 4.3 Convexity of autocratic strategies and payoff relationships

The space of enforceable payoff relationships exhibits several convexity properties that facilitate the construction of new autocratic strategies from known ones. These properties have both theoretical and practical significance since they reveal geometric structure in the space of enforceable constraints and provide methods for "interpolating" between different autocratic strategies. Throughout this subsection, we focus on the discounted case ($\lambda \in [0,1)$) for simplicity, though analogous results hold in the limit $\lambda \to 1$.

Consider two autocratic strategies, $(\sigma_X^0, \sigma_X^*)$ and $(\widetilde{\sigma}_X^0, \widetilde{\sigma}_X^*)$ that enforce $\varphi \equiv 0$ and $\widetilde{\varphi} \equiv 0$, respectively. Assume these are generated by mixed actions $\tau_X^\pm, \widetilde{\tau}_X^\pm \in \Delta(S_X)$ and transition probabilities $(p_0, p^*)$ and $(\widetilde{p}_0, \widetilde{p}^*)$, as in Lemma 2, and that they share a common discount factor $\lambda \in [0,1)$ (which can always be arranged by Proposition 8). For $q \in [0,1]$, we investigate when $\varphi_q := (1-q)\varphi + q\widetilde{\varphi} \equiv 0$ is enforceable and how to construct strategies that enforce it.

**Proposition 10.** Consider some $q \in [0,1]$ and suppose $\tau_X^\pm = \widetilde{\tau}_X^\pm$. Then, there exists a two-point reactive learning strategy that is $(\varphi_q, \lambda_q)$-autocratic for some $\lambda_q \in [0,1)$.

*Proof.* Consider the sets $\Phi_q^\pm$ for the function $\varphi_q$ as in **Eq. 32**. Then $\tau_X^\pm \in \Phi_q^\pm$, due to $\tau_X^\pm \in \Phi_i^\pm$ for $i = 1, 2$, where $\Phi_i^\pm$ are also defined as in **Eq. 32**. The result follows from Proposition 4. $\square$

Proposition 10 does not guarantee that we can enforce a convex combination of $\varphi$ and $\widetilde{\varphi}$ by taking a convex combination of the probabilistic mechanisms $p_1^*$ and $p_2^*$, unless the latter are identical (see **Fig. 1**). However, we can give a sufficient condition for this property:

**Proposition 11.** Consider some $q \in [0,1]$ and suppose $\tau_X^\pm = \widetilde{\tau}_X^\pm$. If

$$\min_{s_Y \in S_Y} \varphi\left(\tau_X^+, s_Y\right) - \max_{s_Y \in S_Y} \varphi\left(\tau_X^-, s_Y\right) = \min_{s_Y \in S_Y} \widetilde{\varphi}\left(\tau_X^+, s_Y\right) - \max_{s_Y \in S_Y} \widetilde{\varphi}\left(\tau_X^-, s_Y\right), \tag{74}$$

then the reactive learning strategy generated by the initial action $qp_0 + (1-q)\widetilde{p}_0$, with the response function $qp^* + (1-q)\widetilde{p}^*$ randomizing between $\tau_X^\pm$, is $(\varphi_q, \lambda)$-autocratic.

*Proof.* The assumption implies that $\psi\left(\tau_X^\pm\right) = \widetilde{\psi}\left(\tau_X^\pm\right)$, as these are defined in the denominators of $p^*, \widetilde{p}^*$ in the proof of Lemma 2. Moreover, for $q \in [0,1]$, the response function $qp^* + (1-q)\widetilde{p}^*$ yields valid probabilities with initial action $qp_0 + (1-q)\widetilde{p}_0$. The resulting reactive learning strategy, $(qp_0 + (1-q)\widetilde{p}_0, qp^* + (1-q)\widetilde{p}^*)$, and the function $\Psi := \psi_1\left(\tau_X^+\right) - \psi_1\left(\tau_X^-\right)$ satisfy the generalized next-round correction condition with discount factor $\lambda$, and the result then follows from Proposition 1. $\square$

In the donation game, ALLD enforces the line $cu_X = -bu_Y$, while ALLC enforces the line $cu_X = -bu_Y + b^2 - c^2$. It is easily verified that the condition of Proposition 11 holds. As a result, an agent can enforce all lines of slope $-b/c$ intersecting the feasible region. In the donation game in particular, the convex hull of all such feasible lines is equal to the entire payoff region (see **Fig. 2**).

## 4.4 Equalizing and self-equalizing strategies for generic payoff functions

$X$ can set her own payoff to some constant, $K$, if $u_X - K \equiv 0$ is enforceable. Assuming that the aim of player $X$ is setting her own score and that no trivial strategies exist, we deduce from
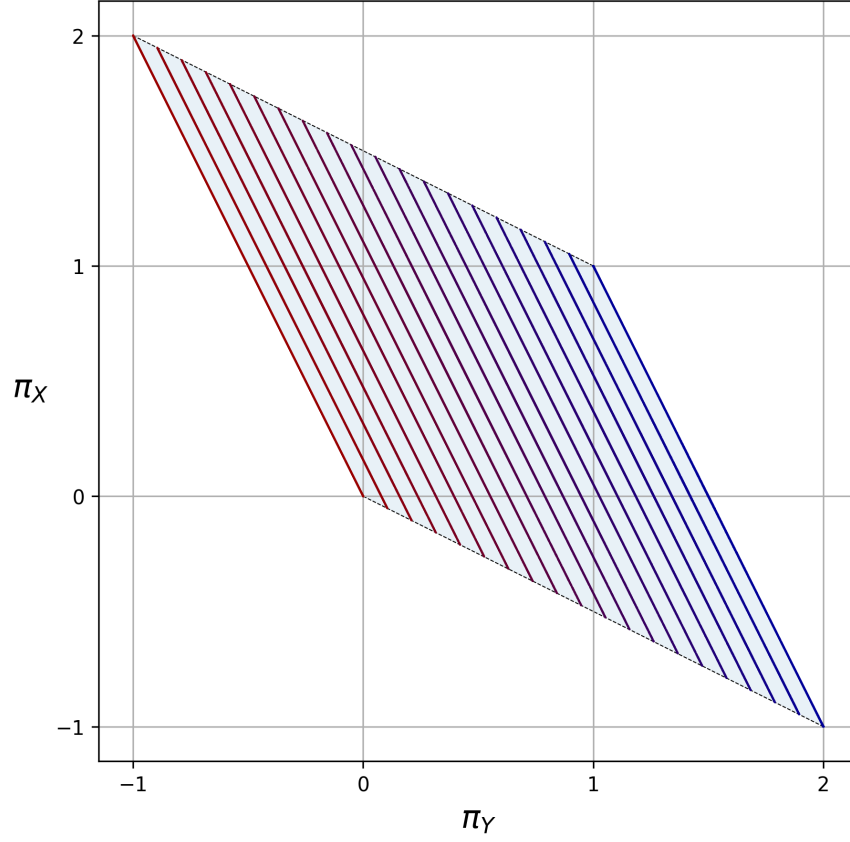
**Figure 2:** Pencil of enforceable lines in the donation game. The boundary lines $cu_X = -bu_Y$ (enforced by ALLD) and $cu_X = -bu_Y + b^2 - c^2$ (enforced by ALLC) are shown, along with intermediate parallel lines that can be enforced by convex combinations of these strategies. By Proposition 11, an agent can enforce any line in this family by appropriately mixing between punishment and forgiveness. The convex hull of all such enforceable lines equals the entire payoff region, demonstrating complete unilateral control over expected payoff outcomes in the repeated donation game.

Corollary 3 that she can set her payoff equal to $K$ if and only if

$$K \in \left[ \min_{\tau_X \in \Delta(S_X)} \max_{s_Y \in S_Y} u_X(\tau_X, s_Y), \max_{\tau_X \in \Delta(S_X)} \min_{s_Y \in S_Y} u_X(\tau_X, s_Y) \right]. \tag{75}$$

Note that the maximum payoff value that can be set by $X$ is $\max_{\tau_X \in \Delta(S_X)} \min_{s_Y \in S_Y} u_X(\tau_X, s_Y)$. This results in a simple relationship between an agent's control of their payoff in a repeated game and that of their optimal return in a one-shot game: an agent can achieve their (average) security level when the enforceability interval **Eq. 75** is well-defined.

Accordingly, suppose player $X$ has an incentive to set her opponent's expected payoff. This case corresponds to enforcing $u_Y - K \equiv 0$. In similar fashion to the above, it can be shown that if she desires to "punish" her opponent in the worst possible way, the best she can hope for is setting $Y$'s payoff to $\min_{\tau_X \in \Delta(S_X)} \max_{s_Y \in S_Y} u_Y(\tau_X, s_Y)$, which by the minimax theorem is exactly equal to $\max_{s_Y \in \Delta(S_Y)} \min_{\tau_X \in \Delta(S_X)} u_Y(\tau_X, s_Y)$. This is precisely the security level of the co-player.

## 4.5  Symmetric objective functions

Next, we prove that an enforceable symmetric relationship is either trivially enforceable or enforceable with $\lambda_{\min} = 1$. In other words, if $X$'s aim is to control a symmetric payoff region, then she either needs a memory-less plan or an "infinite" amount of patience.

**Proposition 12.** Suppose that $S_X = S_Y$ and that $\varphi : S_X \times S_Y \to \mathbb{R}$ is symmetric. Then, $\varphi \equiv 0$ is enforceable if and only if there exists a trivial strategy or $\max_{s_X \in \Delta(S_X)} \min_{s_Y \in \Delta(S_Y)} \varphi(s_X, s_Y) = 0$.

*Proof.* Suppose that $\varphi$ is not 0-enforceable. We know $\max_{s_X \in \Delta(S_X)} \min_{s_Y \in \Delta(S_Y)} \varphi(s_X, s_Y) = 0$ implies that $J(\varphi) = \{0\}$ by the minimax theorem and the symmetry of $\varphi$, thus $\varphi \equiv 0$ is enforceable with $\lambda_{\min} = 1$ by Proposition 6. Conversely, let $\varphi$ be enforceable but not 0-enforceable. Then, it is $\lambda$-enforceable with $\lambda \in (0, 1]$ and $J(\varphi)$ is well-defined and it contains 0 due to Corollary 3. However, the symmetry of $\varphi$ and the minimax theorem imply that $J(\varphi)$ has to be trivial, i.e., the assumption of Proposition 6 holds. $\square$

It is important to note that Proposition 12 also holds for skew-symmetric relations. This follows from the fact that $\varphi \equiv 0$ if and only if $-\varphi \equiv 0$. We therefore deduce that TFT can be an autocratic strategy and enforce fair payoff relationships for symmetric games like IPD, only in the undiscounted setting.

## 4.6  "Zero-sum" strategies for generic payoff functions

Here, we call "zero-sum" autocratic strategies all autocratic plans of $X$ that can enforce $u_X = -u_Y$ in expectation. Corollary 3 directly yields the following:

**Proposition 13.** There exists a zero-sum autocratic strategy for player $X$ if and only if there exist $\tau_X^\pm \in \Delta(S_X)$ such that $u_X(\tau_X^+, s_Y) \geqslant -u_Y(\tau_X^+, s_Y)$ and $u_X(\tau_X^-, s_Y) \leqslant -u_Y(\tau_X^-, s_Y)$ for all $s_Y \in S_Y$.

# 5  Applications and examples

We now apply our theoretical framework to four variants of a classical social dilemma: the prisoner's dilemma, a nonlinear donation game, an asymmetric donation game, and the hawk-dove game. Throughout, we focus on linear payoff relationships relative to a reference line, such

that the desired relationship is either exactly the reference line or else intersects it in a unique point. For example, the line $\kappa - \pi_X = \chi (\kappa - \pi_Y)$ intersects the reference line $\pi_X = \pi_Y$ at $(\kappa, \kappa)$ provided $\chi \neq 1$.

## 5.1  Prisoner's dilemma

In the prisoner's dilemma, the payoff matrix of **Eq. 1** satisfies $T > R > P > S$. We impose the standard constraint of $2P < S + T < 2R$ as well, which means that the the mean payoff for alternating cooperation and defection lies between the payoffs for mutual defection and mutual cooperation. This condition is not strictly necessary, but we use it to simplify the exposition. With the objective function $\varphi (s_X, s_Y) = \kappa - u_X (s_X, s_Y) - \chi (\kappa - u_Y (s_X, s_Y))$, we have

$$\varphi (C, C) = (\chi - 1) (R - \kappa) ; \tag{76a}$$
$$\varphi (C, D) = (\kappa - S) + \chi (T - \kappa) ; \tag{76b}$$
$$\varphi (D, C) = - (T - \kappa) - \chi (\kappa - S) ; \tag{76c}$$
$$\varphi (D, D) = - (\chi - 1) (\kappa - P) . \tag{76d}$$

We first make a few remarks on viable values of $\kappa$ and $\chi$. If $\chi = 1$, then $\kappa$ is irrelevant since it cancels out in the definition of $\varphi$. If $\chi \neq 1$, then we must have $\kappa \in [P, R]$ for an autocratic strategy to exist; otherwise, $\varphi (C, C)$ and $\varphi (D, D)$ have the same sign, which means that $\Phi_X^+$ and $\Phi_X^-$ cannot simultaneously be non-empty (by Proposition 4). If $\chi < 1$ and $\kappa \in [P, R]$, then for $\Phi_X^+$ and $\Phi_X^-$ to both be non-empty, we must have [11]

$$\chi \leqslant \min \left\{ - \frac{T - \kappa}{\kappa - S}, - \frac{\kappa - S}{T - \kappa} \right\} . \tag{77}$$

### 5.1.1  Trivial autocratic strategies

In this game, a trivial autocratic strategy is a value $p \in [0, 1]$ such that $\varphi (p, C) = \varphi (p, D) = 0$. For every $p \in [0, 1]$, we can simply solve for $\kappa$ and $\chi$ in these equations to obtain

$$\kappa = P + p (S + T - 2P) + p^2 (R - S - T + P) ; \tag{78a}$$
$$\chi = \frac{T - P + p (R - S - T + P)}{S - P + p (R - S - T + P)} . \tag{78b}$$

For this pair $(\kappa, \chi)$, playing $C$ with probability $p$ in every round is $(\varphi, 0)$-autocratic.

### 5.1.2  Non-trivial autocratic strategies

We now assume that we have already classified trivial autocratic strategies and we are in the situation in which $\varphi \equiv 0$ is not 0-enforceable. Assuming $(\kappa, \chi)$ is a viable pair, we see that $\delta_C \in \Phi_X^+$ and $\delta_D \in \Phi_X^-$ when $\chi \geqslant 1$, and $\delta_D \in \Phi_X^+$ and $\delta_C \in \Phi_X^-$ when $\chi < 1$. In the non-additive prisoner's dilemma, where $R - T \neq S - P$, we cannot guarantee that the minimizer $(\tau_X^+, \tau_X^-)$ for $\lambda_{\min}$ is attained in pure actions. By Lemma 6, we can restrict the search to pairs of pure actions provided that the differences

$$\varphi (C, C) - \varphi (C, D) = -\chi (T - R) - (R - S) ; \tag{79a}$$

$$\varphi\left(D,C\right)-\varphi\left(D,D\right)=-\chi\left(P-S\right)-\left(T-P\right) \tag{79b}$$

both have the same sign. For $\chi < 1$, this condition requires

$$\chi \leqslant \min\left\{-\frac{R-S}{T-R},-\frac{T-P}{P-S}\right\}=\min_{\kappa\in[P,R]}\min\left\{-\frac{T-\kappa}{\kappa-S},-\frac{\kappa-S}{T-\kappa}\right\} \tag{80}$$

since $2P < S+T < 2R$. In particular, for any $(\kappa,\chi)$ with $\kappa \in [P,R]$ and

$$\min\left\{-\frac{R-S}{T-R},-\frac{T-P}{P-S}\right\} < \chi \leqslant \min\left\{-\frac{T-\kappa}{\kappa-S},-\frac{\kappa-S}{T-\kappa}\right\}, \tag{81}$$

we resort to Proposition 9 to calculate $\lambda_{\min}$. For all other viable values of $(\kappa,\chi)$, we can use Proposition 4, which gives a $(\varphi,\lambda)$-autocratic strategy if and only if $\lambda \geqslant \lambda_{\min}$, where

$$\lambda_{\min} = \begin{cases} 1-\dfrac{\min_{s_Y\in S_Y}\varphi\left(C,s_Y\right)-\max_{s_Y\in S_Y}\varphi\left(D,s_Y\right)}{\max\left\{\begin{smallmatrix}\max_{s_Y\in S_Y}\varphi(C,s_Y)-\max_{s_Y\in S_Y}\varphi(D,s_Y),\\ \min_{s_Y\in S_Y}\varphi(C,s_Y)-\min_{s_Y\in S_Y}\varphi(D,s_Y)\end{smallmatrix}\right\}} & \chi \geqslant 1, \\[20pt] 1-\dfrac{\min_{s_Y\in S_Y}\varphi\left(D,s_Y\right)-\max_{s_Y\in S_Y}\varphi\left(C,s_Y\right)}{\max\left\{\begin{smallmatrix}\max_{s_Y\in S_Y}\varphi(D,s_Y)-\max_{s_Y\in S_Y}\varphi(C,s_Y),\\ \min_{s_Y\in S_Y}\varphi(D,s_Y)-\min_{s_Y\in S_Y}\varphi(C,s_Y)\end{smallmatrix}\right\}} & \chi \leqslant \min\left\{-\dfrac{R-S}{T-R},-\dfrac{T-P}{P-S}\right\}. \end{cases} \tag{82}$$

If $\chi = 1$, then $\kappa$ is irrelevant and $\lambda_{\min} = 1$. If $\chi > 1$ and $\kappa \in [P,R]$, then $\varphi\left(p,D\right) > \varphi\left(p,C\right)$ for all $p \in [0,1]$, so **Eq. 82** reduces to

$$\begin{aligned} \lambda_{\min} &= 1-\frac{\varphi\left(C,C\right)-\varphi\left(D,D\right)}{\max\left\{\varphi\left(C,D\right)-\varphi\left(D,D\right),\varphi\left(C,C\right)-\varphi\left(D,C\right)\right\}} \\ &= 1-\frac{\left(\chi-1\right)\left(R-P\right)}{\max\left\{-S+\chi T-\left(\chi-1\right)P,\left(\chi-1\right)R+T-\chi S\right\}}. \end{aligned} \tag{83}$$

Since $\chi > 1$, we see that $\varphi\left(C,D\right)-\varphi\left(D,D\right) \geqslant \varphi\left(C,C\right)-\varphi\left(D,C\right)$ if and only if $S+T \geqslant R+P$. Therefore, we have

$$\lambda_{\min} = \begin{cases} 1-\dfrac{\left(\chi-1\right)\left(R-P\right)}{-S+\chi T-\left(\chi-1\right)P} & S+T \geqslant R+P, \\[14pt] 1-\dfrac{\left(\chi-1\right)\left(R-P\right)}{\left(\chi-1\right)R-\chi S+T} & S+T < R+P. \end{cases} \tag{84}$$

Finally, if $\chi \leqslant \min\left\{-\left(R-S\right)/\left(T-R\right),-\left(T-P\right)/\left(P-S\right)\right\}$ and $\kappa \in [P,R]$, then **Eq. 82** gives

$$\begin{aligned} \lambda_{\min} &= 1-\frac{\varphi\left(D,D\right)-\varphi\left(C,C\right)}{\max\left\{\varphi\left(D,C\right)-\varphi\left(C,C\right),\varphi\left(D,D\right)-\varphi\left(C,D\right)\right\}} \\ &= 1-\frac{\left(1-\chi\right)\left(R-P\right)}{\max\left\{\left(1-\chi\right)R-T+\chi S,S-\chi T-\left(1-\chi\right)P\right\}}. \end{aligned} \tag{85}$$
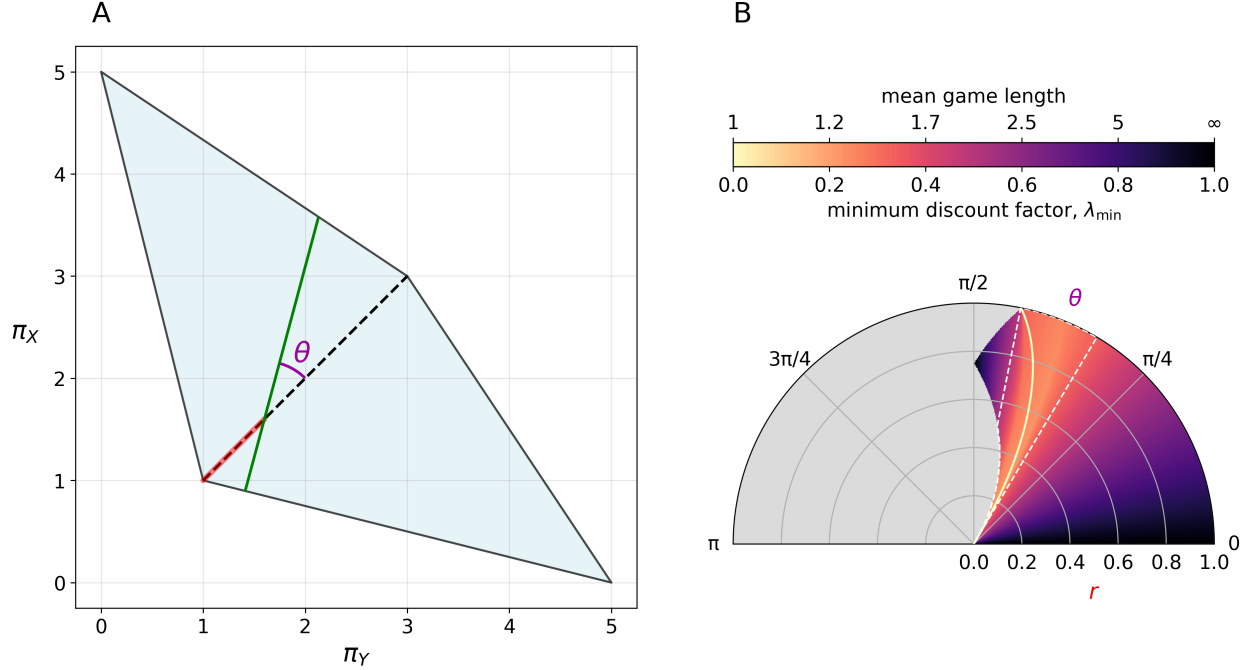
**Figure 3:** Heat map on enforceability of the linear payoff relationship $\varphi = \kappa - u_X(s_X, s_Y) - \chi(\kappa - u_Y(s_X, s_Y)) \equiv 0$ for all values $\kappa, \chi \in \mathbb{R}$, for the repeated prisoner's dilemma with $(R, S, T, P) = (3, 0, 5, 1)$. $\theta$ is the angle between the payoff relationship (green) and the reference line (dashed), and $r$ represents the fraction (red) of the reference line made up by the intersection point. The region enclosed by the white dashed line is the set of $(r, \theta)$ for which at least one of $\tau_X^+$ and $\tau_X^-$ is non-pure in the optimizer for $\lambda_{\min}$ (**Eq. 46**). For this particular game, we have $\kappa = P + r(R - P)$ and $\chi = \tan(\theta + \pi/4)$.

Given the upper bound on $\chi$, we then see that

$$
\lambda_{\min} = \begin{cases} 1 - \dfrac{(1 - \chi)(R - P)}{S - \chi T - (1 - \chi) P} & S + T < R + P, \\[3mm] 1 - \dfrac{(1 - \chi)(R - P)}{(1 - \chi) R - T + \chi S} & S + T \geqslant R + P. \end{cases} \tag{86}
$$

Of course, in all cases, if $\lambda_{\min} > 1$, then there exists no strategy enforcing $\varphi \equiv 0$ for a repeated game with discounting, since the discount factor represents a probability. For the standard (non-additive) payoff parameters of a prisoner's dilemma ($R = 3$, $S = 0$, $T = 5$, and $P = 1$), **Fig. 3** summarizes the enforceable lines.

## 5.2 Nonlinear donation game

The donation game has a relatively simple structure and has been studied extensively in the context of ZD strategies [6, 10, 31]. A slightly more realistic social dilemma occurs when the agents have multiple choices when it comes to paying up a cost and gaining a benefit. For

**Figure 4:** Enforceability of linear payoff relationships in the three-action nonlinear donation game with parameters $b_1 = 3$, $c_1 = 1$, $b_2 = 4$, $c_2 = 2.5$ (satisfying $b_2 - c_2 < b_1 - c_1$). The heatmap shows the minimum discount factor $\lambda_{\min}$ required to enforce $\varphi = \kappa - u_X - \chi (\kappa - u_Y) \equiv 0$ across different values of $\kappa$ and $\chi$. Due to the game's additive structure, any enforceable relationship can be implemented using a two-point reactive strategy, significantly simplifying the strategy space compared to general memory-one approaches.

simplicity, consider the three-action extension of the donation game, with payoff matrix

$$
\begin{array}{c}
\phantom{C_1} \\
\begin{array}{c} C_1 \\ C_2 \\ D \end{array}
\end{array}
\begin{array}{c}
\begin{array}{ccc} C_1 & C_2 & D \end{array} \\
\left(
\begin{array}{ccc}
b_1 - c_1,\ b_1 - c_1 & b_2 - c_1,\ b_1 - c_2 & -c_1,\ b_1 \\
b_1 - c_2,\ b_2 - c_1 & b_2 - c_2,\ b_2 - c_2 & -c_2,\ b_2 \\
b_1,\ 0 & b_2,\ 0 & 0,\ 0
\end{array}
\right).
\end{array}
\tag{87}
$$

Here, we assume that $0 < c_1 < c_2$, $0 < b_1 < b_2$, $c_1 < b_1$, and $c_2 < b_2$, but that $b_2 - c_2 < b_1 - c_1$. In other words, playing $C_2$ is more costly, but more beneficial than playing $C_1$. Yet, when both players choose the same level, mutual "$C_1$" strictly Pareto-dominates mutual "$C_2$."

The game's payoff structure is additive: for each action profile, we can write $u_X (s_X, s_Y) = \phi_X (s_X) + \phi_Y (s_Y)$. By symmetry, $u_Y$ is also additive. Consequently, for linear payoff relationships $\varphi = \kappa - u_X - \chi (\kappa - u_Y)$, the function $\varphi$ inherits this additive structure.

Since this game is additive and $\varphi$ is linear, Theorem 2 tells us that any enforceable payoff relationship can be implemented using a reactive strategy that conditions solely on the opponent's most recent action. Moreover, by Lemma 3, enforcement can be achieved using a two-point reactive strategy that mixes between two fixed distributions based on the opponent's last action.

The presence of three actions introduces additional strategic complexity compared to the standard two-action donation game. While mutual cooperation at level $C_1$ yields the Pareto-efficient outcome $(b_1 - c_1, b_1 - c_1)$, players face a tension between contributing at the higher level

$C_2$ (which benefits the opponent more) and defecting entirely. This structure creates richer possibilities for autocratic strategies, as player $X$ can condition their response not just on whether $Y$ cooperated, but also on the level of cooperation chosen. **Fig. 4** illustrates that the enforceability landscape extends naturally from the two-action case, with the minimum discount factor varying smoothly across the parameter space $(\kappa, \chi)$.

## 5.3 Fairness and equality in asymmetric games

In symmetric games like the standard prisoner's dilemma, "fair" strategies enforcing $\varphi = u_X - u_Y$ are well-studied [6], with TFT being the canonical example. And we have seen that payoff equality can be enforced in expectation if and only if $\lambda \to 1$ (Proposition 12). However, the distinction between equality and "fairness" becomes crucial in asymmetric social dilemmas, where two agents might have different abilities or resources. Consider the asymmetric donation game with payoff matrix

$$
\begin{array}{cc}
 & \begin{array}{cc} C & \qquad\qquad D \end{array} \\
\begin{array}{c} C \\ D \end{array} &
\left( \begin{array}{cc}
b_Y - c_X,\ b_X - c_Y & -c_X,\ b_X \\
b_Y,\ -c_Y & 0,\ 0
\end{array} \right),
\end{array}
\tag{88}
$$

where $b_X > c_X > 0$ and $b_Y > c_Y > 0$. In other words, as cooperators, $X$ pays $c_X$ to donate $b_X$ and $Y$ pays $c_Y$ to donate $b_Y$. As defectors, they both pay nothing and donate nothing.

In this game, the natural reference line is not $\pi_X = \pi_Y$ but rather the line through $(0,0)$ and $(b_Y - c_X, b_X - c_Y)$ since this represents the line between payoffs for mutual defection and payoffs for mutual cooperation. A strategy is "fair" if it enforces a payoff relationship along this reference line, reflecting proportional sharing that accounts for the asymmetric costs and benefits. In contrast, a strategy emphasizing equality enforces $\pi_X = \pi_Y$ regardless of the players' differing contributions.

Like the standard donation game, this asymmetric variant is additive. By Theorem 2, any enforceable payoff relationship can be implemented using a reactive strategy. The key distinction from symmetric games emerges when comparing fairness and equality: while TFT enforces equality in the symmetric donation game (where $b_X = b_Y$ and $c_X = c_Y$), in the asymmetric case these two objectives diverge. Fair strategies that respect the natural reference line require $\lambda \to 1$, making this constraint practically infeasible in finite-horizon interactions. By contrast, enforcing equality can be achieved with smaller discount factors, as shown in **Fig. 5**. This illustrates a fundamental principle: in asymmetric settings, unilateral enforcement of fair outcomes (proportional to players' contributions), is significantly more demanding than enforcement of equal outcomes.

## 5.4 Hawk-dove game: a "weak" social dilemma

Finally, we consider a weak social dilemma [32] known as the "hawk-dove" [33] (or "snowdrift" [34, 35]) game. The toy narrative behind this game is that there is an aggressive type (hawk) and a peaceful type (dove), and there is a common resource of value $V$ for which two individuals compete. When two doves meet, they share the resource, each receiving $V/2$. When a hawk meets a dove, the hawk takes the resource and leaves nothing for the dove. However, when two
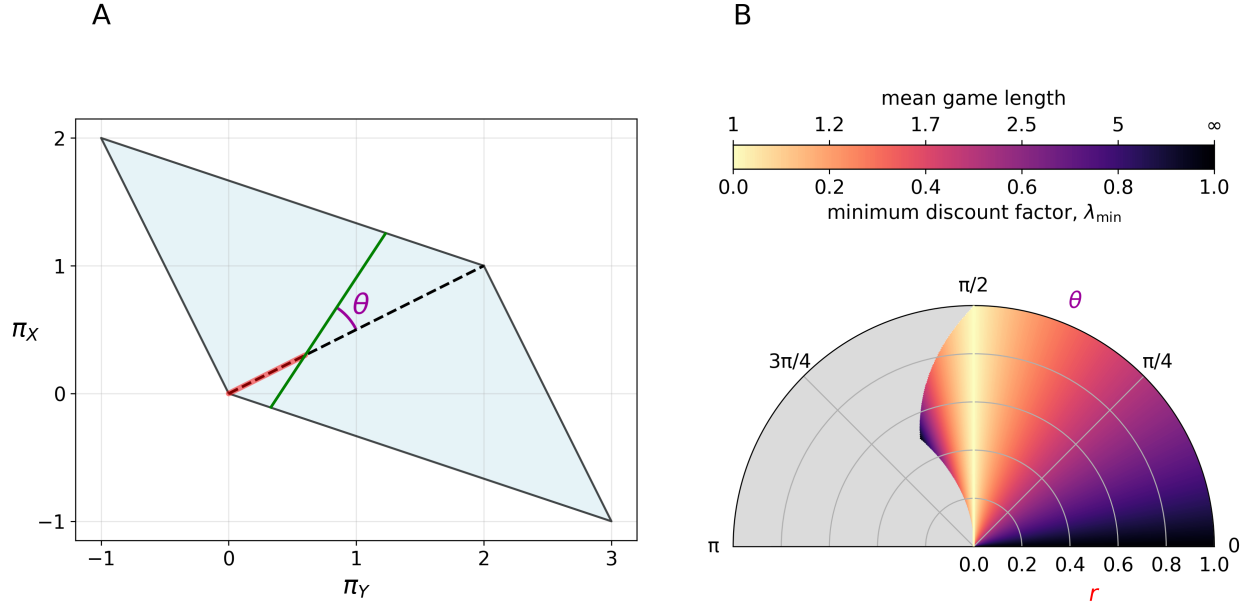
**Figure 5:** Enforceability in the asymmetric donation game with $b_X = 3$, $c_X = 1$, $b_Y = 2$, and $c_Y = 1$. The heatmap displays the minimum discount factor $\lambda_{\min}$ for enforcing $\varphi = \kappa - u_X - \chi (\kappa - u_Y) \equiv 0$. The natural reference line (indicated by the dashed line in the payoff space) connects mutual defection $(0,0)$ to mutual cooperation $(b_Y - c_X, b_X - c_Y)$, reflecting the asymmetric costs and benefits. Strategies enforcing equality ($\pi_X = \pi_Y$, corresponding to $\chi = 1$) require $\lambda \to 1$, while "fair" strategies that enforce proportional sharing along the reference line are more readily achievable.

**Figure 6:** Heat map on enforceability of the linear payoff relationship $\varphi = \kappa - u_X(s_X, s_Y) - \chi(\kappa - u_Y(s_X, s_Y)) \equiv 0$ for all values $\kappa, \chi \in \mathbb{R}$, for the repeated hawk-dove game with $(R, S, T, P) = (-1, 2, 0, 1)$. $\theta$ is the angle between the payoff relationship (green) and the reference line (dashed), and $r$ represents the fraction (red) of the reference line made up by the intersection point. The region enclosed by the white dashed line is the set of $(r, \theta)$ for which at least one of $\tau_X^+$ and $\tau_X^-$ is non-pure in the optimizer for $\lambda_{\min}$ (**Eq. 46**). For this game, we have $\kappa = P + r(R - P)$ and $\chi = \tan(\theta + \pi/4)$.

hawks meet, they fight and incur a cost that effectively lowers the value of the resource by $C > V$, resulting in $(V - C)/2$ to each.

**Fig. 6** shows which linear relationships are enforceable in this game, using parameters $V = 2$ and $C = 4$ as examples. Here, like in the non-additive prisoner's dilemma, we observe regions for which $\lambda_{\min}$ can be attained only when at least one of the two mixed actions in a two-point reactive learning strategy is non-pure.

## 6   Discussion

Our results establish a complete characterization of enforceable payoff relationships in discounted repeated games. The main theoretical contribution answers a question about the role of strategic complexity: extending memory beyond a simple reactive learning structure provides no additional power for enforcing payoff constraints. Any payoff relationship that can be enforced by a strategy of arbitrary memory can be implemented using a two-point reactive learning strategy. This universality result demonstrates that the space of reactive learning strategies captures all possible enforceable constraints, regardless of the opponent's strategic sophistication.

One consequence is that a search for autocratic strategies need only consider the computationally tractable space of reactive learning strategies, rather than the vast and intractable space of behavioral strategies. Our constructive characterization provides explicit formulas for both

verifying enforceability and computing the minimum discount factor required. For arbitrary functions, $\varphi$, these computations can be performed in polynomial time using linear programming. The computational tractability of our characterization stands in stark contrast to the general intractability of analyzing finite-memory strategies. For memory-$m$ strategies in games with $k$ actions per player, the strategy space has dimension $k^{m+1}$, making exhaustive analysis infeasible even for modest values of $m$ and $k$ [27]. The reduction to two-point reactive learning strategies eliminates this exponential dependence on memory length, reducing the problem to a fixed-dimensional space determined only by the number of actions in the stage game.

These findings resolve several open questions in the literature on zero-determinant strategies. First, we provide a definitive answer to the question raised by Hilbe et al. [14] regarding the existence of autocratic strategies in discounted games. Our necessary and sufficient conditions establish precisely when a payoff relationship is enforceable, and our formulas for the minimum discount factor extend the existence results of Hilbe et al. [11] by providing exact thresholds. Second, we demonstrate that the initial action plays a crucial role in determining enforceability in discounted games, contrasting with undiscounted settings where initial conditions are often irrelevant. This distinction highlights fundamental differences between discounted and undiscounted repeated games that have not been fully appreciated in prior work. The structural properties we establish (e.g., convexity of enforceable relationships, the dichotomy for symmetric relationships, and the polynomial-time computability) reveal that the space of autocratic strategies has favorable geometric and algorithmic properties. The convexity result implies that if two payoff relationships are enforceable, any convex combination is also enforceable (under appropriate conditions on the underlying strategies), providing a way to construct new autocratic strategies from known ones. The dichotomy for symmetric relationships shows that fairness constraints are fundamentally incompatible with discounting: symmetric relationships are either trivially enforceable or require the limiting case of an infinite horizon. This explains why strategies like tit-for-tat, which enforce equal payoffs in the prisoner's dilemma, lose this property in any discounted setting.

Although we have framed our results primarily for two-player games, the framework extends naturally to multiplayer settings where coalitions of players coordinate to enforce payoff constraints on the larger group. A coalition $I \subseteq \{1, \ldots, N\}$ can use a correlated strategy to enforce a linear relationship $\sum_{i=1}^{N} \alpha_i \pi_i + \alpha_0 = 0$ on the expected payoffs of all players. Our characterization carries over to this setting: any enforceable coalitional constraint can be implemented using a reactive learning strategy for the coalition that conditions on the previous actions of players outside the coalition and the coalition's own previous mixed action. In fact, one need only replace $X$ by $I$ (the coalition) and $Y$ by $-I$ (the anti-coalition) in Theorem 1. The only subtlety is that a strategy for $I$ allows for correlations, in the sense that it is a map $\sigma_I : \mathcal{H} \to \Delta(S_I)$ rather than from $\mathcal{H}$ to $\prod_{i \in I} \Delta(S_i)$. This, we believe, is reasonable, when a subset of a larger group strategically coordinate to control outcomes for others. There has been limited work in the area of multiplayer payoff enforcement [9, 25, 36, 37], but this area is not well-understood.

This coalitional perspective is especially relevant for applications in multi-agent reinforcement learning and algorithmic game theory [38–40]. Zero-determinant strategies allow a coalition to effectively reshape the incentive landscape faced by learning agents outside the coalition. If external agents are optimizing their policies using gradient-based methods or other adaptive algorithms, the coalition can enforce constraints that guide the learning dynamics toward desirable equilibria or prevent convergence to undesirable outcomes. Doing so is relevant in, for

example, specific domains such as climate agreements [41] and algorithmic collusion [42]. Our results provide concrete tools for designing such coalitional strategies and understanding their limitations.

Autocratic strategies are closely linked to evolutionary game theory, where populations of agents adapt their strategies over time through selection, mutation, or learning [15, 43–45]. In fact, this is the setting in which linear payoff constraints were first recognized [6, 46]. An autocratic strategy that enforces a favorable payoff relationship can maintain a consistent advantage over a wide range of opponents, potentially allowing it to persist in evolutionary competition despite not being a Nash equilibrium of the stage game. The robustness of reactive learning strategies (in the sense of their ability to enforce constraints against arbitrary opponents) suggests they may be especially resilient to invasion attempts and environmental perturbations. In learning dynamics, the presence of an autocratic player fundamentally alters the optimization landscape for other agents. If player $Y$ is adapting through reinforcement learning or evolutionary search, player $X$'s autocratic strategy constrains the payoffs $Y$ can achieve along any learning trajectory. Understanding these constraints is crucial for predicting the long-run outcomes of multi-agent learning systems and for designing interventions that guide learning toward socially beneficial equilibria.

Several important questions remain open. First, while we characterize enforceable relationships for arbitrary payoff functions, our explicit formulas for minimum discount factors apply most directly to finite action spaces. Extensions to continuous action spaces or games with state-dependent payoffs would require additional technical machinery, though we expect the fundamental principles to carry over. Second, our framework assumes that players observe actions perfectly and condition their strategies on these observations. In settings with imperfect monitoring or private information, enforcement becomes more subtle, and the relationship between memory and enforceability may change [26]. Future work could explore the robustness of autocratic strategies to noise and misperception, characterize the set of enforceable relationships when both players simultaneously attempt to enforce constraints, and investigate the evolutionary stability of populations containing autocratic strategists. While much of the existing literature focuses on linear payoff relationships, our framework naturally extends to arbitrary nonlinear constraints on expected payoffs. This generality opens new avenues for studying strategic behavior in environments where traditional zero-determinant techniques fail. For instance, payoff relationships involving products, maxima, or other nonlinear combinations of player payoffs can be analyzed using our framework, potentially revealing new classes of enforceable constraints in economic and biological applications. The extension to nonlinear relationships is particularly relevant for settings where players care about relative performance, inequity aversion, or other behavioral considerations that induce nonlinear preferences over payoff profiles [47]. Our results suggest that even in these complex preference structures, the fundamental limits on enforceability are determined by the same geometric separation conditions that govern linear relations.

The universality of reactive learning strategies within the class of autocratic strategies is perhaps surprising given the apparent richness of the space of behavioral strategies with arbitrary memory. Our results show that this richness is illusory for the purpose of enforcing payoff constraints: the geometric separation conditions that determine enforceability depend only on the stage game payoffs and the discount factor, not on the complexity of the enforcing strategy. On the other hand, reactive learning strategies can be thought of as longer-memory behavioral strategies with a "right-invariance" property. This property allows longer histories of simple

information to be "rolled up" into shorter memories of richer information. Thus, we believe that this finding provides both theoretical closure on the role of memory in zero-determinant strategies and practical guidance for finding autocratic strategies in strategic environments.

## Acknowledgments

## References

[1] D. Fudenberg and E. Maskin. The Folk Theorem in Repeated Games with Discounting or with Incomplete Information. *Econometrica*, 54(3):533–554, 1986. doi: 10.2307/1911307.

[2] J. Foerster, R. Y. Chen, M. Al-Shedivat, S. Whiteson, P. Abbeel, and I. Mordatch. Learning with Opponent-Learning Awareness. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '18, pages 122–130. International Foundation for Autonomous Agents and Multiagent Systems, 2018.

[3] A. McAvoy, J. Kates-Harbeck, K. Chatterjee, and C. Hilbe. Evolutionary instability of selfish learning in repeated games. *PNAS Nexus*, 1(4), 2022. doi: 10.1093/pnasnexus/pgac141.

[4] M. A. Nowak. Five rules for the evolution of cooperation. *Science*, 314(5805):1560–1563, 2006. doi: 10.1126/science.1133755.

[5] R. L. Trivers. The Evolution of Reciprocal Altruism. *The Quarterly Review of Biology*, 46(1): 35–57, 1971. doi: 10.1086/406755.

[6] W. H. Press and F. J. Dyson. Iterated prisoner's dilemma contains strategies that dominate any evolutionary opponent. *Proceedings of the National Academy of Sciences*, 109(26):10409–10413, 2012. doi: 10.1073/pnas.1206569109.

[7] R. Axelrod and W. Hamilton. The evolution of cooperation. *Science*, 211(4489):1390–1396, 1981. doi: 10.1126/science.7466396.

[8] R. Axelrod. *The Evolution of Cooperation*. Basic Books, 1984.

[9] L. Pan, D. Hao, Z. Rong, and T. Zhou. Zero-Determinant Strategies in Iterated Public Goods Game. *Scientific Reports*, 5:13096, 2015. doi: 10.1038/srep13096.

[10] A. McAvoy and C. Hauert. Autocratic strategies for iterated games with arbitrary action spaces. *Proceedings of the National Academy of Sciences*, 113(13):3573–3578, 2016. doi: 10.1073/pnas.1520163113.

[11] C. Hilbe, A. Traulsen, and K. Sigmund. Partners or rivals? Strategies for the iterated prisoner's dilemma. *Games and Economic Behavior*, 92:41–52, 2015. doi: 10.1016/j.geb.2015.05.005.

[12] A. McAvoy and C. Hauert. Autocratic strategies for alternating games. *Theoretical Population Biology*, 113:13–22, 2017. doi: 10.1016/j.tpb.2016.09.004.

[13] C. Adami and A. Hintze. Evolutionary instability of zero-determinant strategies demonstrates that winning is not everything. *Nature Communications*, 4, 2013. doi: 10.1038/ncomms3193.

[14] C. Hilbe, M. A. Nowak, and A. Traulsen. Adaptive Dynamics of Extortion and Compliance. *PLoS ONE*, 8(11):e77886, 2013. doi: 10.1371/journal.pone.0077886.

[15] A. J. Stewart and J. B. Plotkin. From extortion to generosity, evolution in the Iterated Prisoner's Dilemma. *Proceedings of the National Academy of Sciences*, 110(38):15348–15353, 2013. doi: 10.1073/pnas.1306246110.

[16] C. Hilbe, B. Wu, A. Traulsen, and M. A. Nowak. Cooperation and control in multiplayer social dilemmas. *Proceedings of the National Academy of Sciences*, 111(46):16425–16430, 2014. doi: 10.1073/pnas.1407887111.

[17] Z. Wang, Y. Zhou, J. W. Lien, J. Zheng, and B. Xu. Extortion can outperform generosity in the iterated prisoner's dilemma. *Nature Communications*, 7:11125, 2016. doi: 10.1038/ncomms11125.

[18] M. Milinski, C. Hilbe, D. Semmann, R. Sommerfeld, and J. Marotzke. Humans choose representatives who enforce cooperation in social dilemmas through extortion. *Nature Communications*, 7:10915, 2016. doi: 10.1038/ncomms10915.

[19] A. McAvoy, U. M. Sehwag, C. Hilbe, K. Chatterjee, W. Barfuss, Q. Su, N. E. Leonard, and J. B. Plotkin. Unilateral incentive alignment in two-agent stochastic games. *Proceedings of the National Academy of Sciences*, 122(25), 2025. doi: 10.1073/pnas.2319927121.

[20] M. Nowak and K. Sigmund. A strategy of win-stay, lose-shift that outperforms tit-for-tat in the Prisoner's Dilemma game. *Nature*, 364(6432):56–58, 1993. doi: 10.1038/364056a0.

[21] C. Hilbe, L. A. Martinez-Vaquero, K. Chatterjee, and M. A. Nowak. Memory-n strategies of direct reciprocity. *Proceedings of the National Academy of Sciences*, 114(18):4715–4720, 2017.

[22] F. Lesigang, C. Hilbe, and N. E. Glynatsi. Can i afford to remember less than you? best responses in repeated additive games. *Economics Letters*, page 112300, 2025.

[23] M. Barlo, G. Carmona, and H. Sabourian. Repeated games with one-memory. *Journal of Economic Theory*, 144(1):312–336, 2009.

[24] M. Ueda. Memory-two zero-determinant strategies in repeated games. *Royal Society open science*, 8(5):202186, 2021.

[25] A. Govaert and M. Cao. Zero-determinant strategies in finitely repeated n-player games. *IFAC-PapersOnLine*, 52(3):150–155, 2019.

[26] A. Mamiya and G. Ichinose. Zero-determinant strategies under observation errors in repeated games. *Physical Review E*, 102(3):032115, 2020.

[27] C. Hauert and H. G. Schuster. Effects of increasing the number of players and memory size in the iterated Prisoner's Dilemma: a numerical approach. *Proceedings of the Royal Society of London. Series B: Biological Sciences*, 264(1381):513–519, 1997.

[28] R. M. Dawes. Social dilemmas. *Annual Review of Psychology*, 31(1):169–193, 1980. doi: 10.1146/annurev.ps.31.020180.001125.

[29] A. McAvoy and M. A. Nowak. Reactive learning strategies for iterated games. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 475(2223):20180819, 2019. doi: 10.1098/rspa.2018.0819.

[30] J. L. Doob. *Measure Theory*. Springer New York, 1994. doi: 10.1007/978-1-4612-0877-8.

[31] C. Hilbe, M. A. Nowak, and K. Sigmund. Evolution of extortion in Iterated Prisoner's Dilemma games. *Proceedings of the National Academy of Sciences*, 110:6913–6918, 2013. doi: 10.1073/pnas.1214834110.

[32] C. Hauert, F. Michor, M. A. Nowak, and M. Doebeli. Synergy and discounting of co-operation in social dilemmas. *Journal of Theoretical Biology*, 239(2):195–202, 2006. doi: 10.1016/j.jtbi.2005.08.040.

[33] J. Maynard Smith and G. R. Price. The Logic of Animal Conflict. *Nature*, 246(5427):15–18, 1973. doi: 10.1038/246015a0.

[34] R. Sugden. *The Economics of Rights, Co-operation, and Welfare*. B. Blackwell, 1986.

[35] C. Hauert and M. Doebeli. Spatial structure often inhibits the evolution of cooperation in the snowdrift game. *Nature*, 428(6983):643–646, 2004. doi: 10.1038/nature02360.

[36] C. Hilbe, B. Wu, A. Traulsen, and M. A. Nowak. Evolutionary performance of zero-determinant strategies in multiplayer games. *Journal of Theoretical Biology*, 374:115–124, 2015. doi: 10.1016/j.jtbi.2015.03.032.

[37] F. Chen, T. Wu, and L. Wang. Evolutionary dynamics of zero-determinant strategies in repeated multiplayer games. *Journal of Theoretical Biology*, 549:111209, 2022. doi: 10.1016/j.jtbi.2022.111209.

[38] D. Su, H. Peng, G. Zeng, P. Li, A. Li, and Y. Pan. Emergence of cooperation in multi-agent reinforcement learning via coalition labeling and structural entropy. In *Proceedings of the 2025 SIAM International Conference on Data Mining (SDM)*, pages 507–515. SIAM, 2025.

[39] B. D. Bernheim, B. Peleg, and M. D. Whinston. Coalition-proof nash equilibria i. concepts. *Journal of economic theory*, 42(1):1–12, 1987.

[40] T. W. Sandholm and R. H Crites. Multiagent reinforcement learning in the iterated prisoner's dilemma. *Biosystems*, 37(1-2):147–166, 1996.

[41] W. Nordhaus. Dynamic climate clubs: On the effectiveness of incentives in global climate agreements. *Proceedings of the National Academy of Sciences*, 118(45):e2109988118, 2021.

[42] E. Calvano, G. Calzolari, V. Denicolo, and S. Pastorello. Artificial intelligence, algorithmic pricing, and collusion. *American Economic Review*, 110(10):3267–3297, 2020.

[43] J. W. Weibull. *Evolutionary game theory*. MIT Press, 1995.

[44] K. Sigmund and M. A. Nowak. Evolutionary game theory. *Current Biology*, 9(14):R503–R505, 1999.

[45] J. Chen and A. Zinger. The robustness of zero-determinant strategies in Iterated Prisoner's Dilemma games. *Journal of Theoretical Biology*, 357:46–54, 2014. doi: 10.1016/j.jtbi.2014.05.004.

[46] A. J. Stewart and J. B. Plotkin. Extortion and cooperation in the Prisoner's Dilemma. *Proceedings of the National Academy of Sciences*, 109(26):10134–10135, 2012. doi: 10.1073/pnas.1208087109.

[47] E. Fehr and K. M. Schmidt. A theory of fairness, competition, and cooperation. *The Quarterly Journal of Economics*, 114(3):817–868, 1999.