
Strategic stability under regularized learning in games

Abstract

1 In this paper, we examine the long-run behavior of regularized, no-regret learning in
2 finite games. A well-known result in the field states that the empirical frequencies
3 of no-regret play converge to the game’s set of coarse correlated equilibria; however,
4 our understanding of how the players’ *actual* strategies evolve over time is much
5 more limited – and, in many cases, non-existent. This issue is exacerbated by
6 a series of recent results showing that *only* strict Nash equilibria are stable and
7 attracting under regularized learning, thus making the relation between learning
8 and pointwise solution concepts particularly elusive. In lieu of this, we take a more
9 general approach and instead seek to characterize the *setwise* rationality properties
10 of the players’ day-to-day play. To that end, we focus on one of the most stringent
11 criteria of setwise strategic stability, namely that any unilateral deviation from the
12 set in question incurs a cost to the deviator – a property known as *closedness under*
13 *better replies* (club). In so doing, we obtain a remarkable equivalence between
14 strategic and dynamic stability: *a product of pure strategies is closed under better*
15 *replies if and only if its span is stable and attracting under regularized learning.*
16 In addition, we estimate the rate of convergence to such sets, and we show that
17 methods based on entropic regularization (like the exponential weights algorithm)
18 converge at a geometric rate, while projection-based methods converge within a
19 *finite* number of iterations, even with bandit, payoff-based feedback.

20 1 Introduction

21 **Background.** The question of whether players can learn to emulate rational behavior through
22 repeated interactions has been one of the mainstays of non-cooperative game theory – and it has
23 recently gained increased momentum owing to a surge of breakthrough applications to machine
24 learning and data science (from online advertising to auctions and multi-agent reinforcement learning).
25 Informally, this question can be stated as follows:

26 *If each player follows an iterative procedure aiming to increase their individual payoff,*
27 *does the players’ long-run behavior converge to a rationally admissible state?*

28 A natural setting for studying this question is to assume that each player is following a no-regret
29 algorithm, i.e., a policy which is asymptotically as good against a given sequence of payoff functions
30 as the best fixed strategy in hindsight. In this framework, the link between learning and rationality
31 is provided by a folk result which states that, under no-regret learning, the empirical frequency of
32 play converges to the game’s set of *coarse correlated equilibria* (CCE) – also known as the game’s
33 *Hannan set* [22]. This result has been of seminal importance to the field because no-regret play can
34 be achieved via a wide class of “regularized learning” policies, as exemplified by the “*follow-the-*
35 *regularized-leader*” (FTRL) family of algorithms [41, 42] and its variants – optimistic mirror descent
36 [13, 36, 37, 43], HEDGE / EXP3 [4, 5, 9, 10], implicitly normalized forecasters [1, 3], etc.

37 All these policies have (at least) one thing in common: they seek to provide the tightest possible
38 guarantees for each player’s individual regret, thus accelerating convergence to the game’s Hannan set.
39 As such, in games where the marginalization of CCE coincides with the game’s Nash equilibria (like
40 two-player zero-sum games), we obtain a positive equilibrium convergence guarantee: the long-run
41 average frequency of play evolves “as if” the players were rational to begin with – i.e., as if they
42 had full knowledge of the game, common knowledge of rationality, the ability to communicate this
43 knowledge, etc.

44 On the other hand, in many concrete applications – and, in particular, in the context of regularized
45 learning – players learn *independently* from one another, with no common correlating device. By
46 comparison, the Hannan set consists of *correlated* strategies which, when marginalized, may fail
47 even the weakest axioms of rational behavior and rationalizability (such as the elimination of strictly
48 dominated strategies). In particular, a well-known example of Viossat & Zapechelnyuk [45] (which
49 we discuss in detail in Section 4) shows that it is possible to have negative regret for all time, but still
50 employ only strictly dominated strategies throughout the entire horizon of play.

51 The reason for this disconnect is that no-regret play has significant predictive power for the empirical
52 frequency of play – i.e., the empirical distribution of pure strategy *profiles* – but much less so for the
53 players’ day-to-day sequence of play – that is, the evolution of the players’ *actual* mixed strategies
54 over time. In particular, even when the marginalization of the Hannan set is Nash, the actual trajectory
55 of play may – and, in fact, often *does* – diverge away from the game’s set of equilibria [14, 20, 29–31]
56 or exhibits chaotic, unpredictable oscillations [11, 33].

57 Motivated by the above, our paper seeks to understand the rationality properties of the players’ *actual*
58 sequence of play under regularized learning, as encoded by the following question:

59 *Which sets of mixed strategies are stable and attracting under regularized learning?*
60 *Are these sets robust to strategic deviations? And, if so, is the converse also true?*

61 **Our contributions in the context of related work.** This question has attracted significant interest
62 in the literature, especially in its pointwise version, namely: Which mixed strategy profiles are
63 stable and attracting under regularized learning? Are the dynamics’ stable states robust to unilateral
64 deviations? And, if so, are these the only stable states of regularized learning?

65 In the related setting of population games, the answer to this question is sometimes referred to as the
66 “folk theorem of evolutionary game theory” [12, 23, 40, 47]. Somewhat informally, this theorem states
67 that, under the replicator dynamics (the continuous-time analogue of the exponential / multiplicative
68 weights algorithm, itself an archetypal regularized learning method), the following is true for *all*
69 games: only Nash equilibria are (Lyapunov) stable, and a state is stable and attracting under the
70 replicator dynamics if and only if it is a strict Nash equilibrium of the underlying game [23, 47].

71 In the context of regularized learning, [17, 28] recently showed that a similar equivalence holds for
72 the dynamics of FTRL in *continuous* time: a state is stable and attracting under the FTRL dynamics
73 if and only if it is a strict Nash equilibrium. Subsequently, [19] extended this equivalence to an
74 entire class of regularized learning schemes, with different types of feedback and/or possible update
75 structures – from optimistic methods to algorithms run with bandit, payoff-based information. In
76 all these cases, the same principle emerges: a state is asymptotically stable and attracting under
77 regularized learning if and only if it is a strict Nash equilibrium.

78 This is an important pointwise prediction but it does not cover cases where regularized learning
79 algorithms converge to a *set* – not a point. In this case, the very definition of strategic stability is an
80 intricate affair, and there are several definitions that come into play [6, 15, 18, 38]. The first such
81 notion that we consider is that of “resilience to strategic deviations”, namely that every unilateral
82 deviation from said set is deterred by some other element thereof. Our first contribution in this
83 direction is a universal guarantee to the effect that, with probability 1, in any game, and from any
84 initial condition, *the long-run limit of any regularized learning algorithm is a resilient set.*

85 This result is significant in its universality, but the notion of resilience is not sufficiently strong to
86 disallow irrational behavior – and, in fact, it is subject to similar shortcomings as Hannan consistency.
87 On that account, we turn to a much more stringent criterion of *setwise strategic stability*, that of
88 *minimal closedness under better replies* (m-club). This notion, originally due to Ritzberger & Weibull
89 [38], states that any deviation from a product of pure strategies is costly, and it is one of the strictest
90 setwise refinements in game theory; in particular, it refines the notion of closedness under rational
91 behavior (curb) [6], and it satisfies all the seminal *strategic stability* requirements of Kohlberg &
92 Mertens [25], including robustness to strategic payoff perturbations.

93 In this general context, we show that regularized learning enjoys a striking relation with club sets: *A*
94 *product of pure strategies is closed under better replies if and only if its span is stable and attracting*
95 *under regularized learning.* More to the point, we also estimate the rate of convergence to club sets,
96 and we show that convergence occurs at a geometric rate for entropically regularized methods – like
97 HEDGE and EXP3 – and in a *finite number* of iterations under projection-based methods.

98 In light of the above, our results can be seen both as a far-reaching setwise generalization of the folk
 99 theorem of evolutionary game theory, as well as a bona fide algorithmic analogue of a precursor
 100 result for the replicator dynamics [39]. Importantly, our analysis covers several different update
 101 structures – “vanilla” regularized methods, but also their optimistic variants – as well as a wide range
 102 of information models – from full payoff information to bandit, payoff-based feedback.

103 2 Preliminaries

104 We start by recalling some basics from game theory, roughly following the classical treatise of
 105 Fudenberg & Tirole [18]. First, a *finite game in normal form* consists of (i) a finite set of *players*
 106 $i \in \mathcal{N} \equiv \{1, \dots, N\}$; (ii) a finite set of *actions* – or *pure strategies* – \mathcal{A}_i per player $i \in \mathcal{N}$; and (iii) an
 107 ensemble of *payoff functions* $u_i: \prod_j \mathcal{A}_j \rightarrow \mathbb{R}$, each determining the reward $u_i(\alpha)$ of player $i \in \mathcal{N}$
 108 in a given *action profile* $\alpha = (\alpha_1, \dots, \alpha_N)$. Collectively, we will write $\mathcal{A} = \prod_j \mathcal{A}_j$ for the game’s
 109 *action space* and $\Gamma \equiv \Gamma(\mathcal{N}, \mathcal{A}, u)$ for the game with primitives as above.

110 During play, each player $i \in \mathcal{N}$ may randomize their choice of action by playing a *mixed strategy*,
 111 i.e., a probability distribution $x_i \in \mathcal{X}_i := \Delta(\mathcal{A}_i)$ over \mathcal{A}_i that selects $\alpha_i \in \mathcal{A}_i$ with probability $x_{i\alpha_i}$. To
 112 lighten notation, we will identify $\alpha_i \in \mathcal{A}_i$ with the mixed strategy that assigns all weight to α_i (thus
 113 justifying the terminology “pure strategies”). Then, writing $x = (x_i)_{i \in \mathcal{N}}$ for the players’ *strategy*
 114 *profile* and $\mathcal{X} = \prod_i \mathcal{X}_i$ for the game’s *strategy space*, the players’ payoff functions may be extended
 115 to all of \mathcal{X} by setting $u_i(x) := \mathbb{E}_{\alpha \sim x} [u_i(\alpha)] = \sum_{\alpha \in \mathcal{A}} u_i(\alpha) x_\alpha$ where, in a slight abuse of notation,
 116 we write x_α for the joint probability of playing $\alpha \in \mathcal{A}$ under x , i.e., $x_\alpha = \prod_i x_{i\alpha_i}$. This randomized
 117 framework will be referred to as the *mixed extension* of Γ and we will denote it by $\Delta(\Gamma)$.

118 For concision, we will also write $(x_i; x_{-i}) = (x_1, \dots, x_i, \dots, x_N)$ for the strategy profile where player
 119 i plays $x_i \in \mathcal{X}_i$ against the strategy profile $x_{-i} \in \prod_{j \neq i} \mathcal{X}_j$ of all other players (and likewise for pure
 120 strategies). In this notation, we also define each player’s *mixed payoff vector* as

$$v_i(x) = (u_i(\alpha_i; x_{-i}))_{\alpha_i \in \mathcal{A}_i} \quad (1)$$

121 so the payoff to player $i \in \mathcal{N}$ under $x \in \mathcal{X}$ becomes $u_i(x) = \sum_{\alpha_i \in \mathcal{A}_i} u_i(\alpha_i; x_{-i}) x_{i\alpha_i} = \langle v_i(x), x_i \rangle$.
 122 The *best-response correspondence* of player $i \in \mathcal{N}$ is then defined as the set-valued mapping
 123 $\text{br}_i: \mathcal{X} \rightrightarrows \mathcal{X}_i$ given by $\text{br}_i(x) = \arg \max_{x'_i \in \mathcal{X}_i} u_i(x'_i; x_{-i})$ for all $x \in \mathcal{X}$. Extending this over all
 124 players, we will write $\text{br} = \prod_i \text{br}_i$ for the product correspondence $\text{br}(x) = \text{br}_1(x) \times \dots \times \text{br}_N(x)$,
 125 and we will say that $x^* \in \mathcal{X}$ is a *Nash equilibrium* (NE) if $x^* \in \text{br}(x^*)$. Equivalently, given that
 126 $u_i(x'_i; x_{-i})$ is linear in x'_i , we conclude that x^* is a Nash equilibrium if and only if $u_i(x^*) \geq u_i(\alpha_i; x_{-i}^*)$
 127 for all $\alpha_i \in \mathcal{A}_i$ and all $i \in \mathcal{N}$.

128 As a final point of note, if x^* is a Nash equilibrium where each player has a unique best response – that
 129 is, $\text{br}_i(x^*) = \{x_i^*\}$ for all $i \in \mathcal{N}$ – we will say that x^* is *strict* because, in this case, $u_i(x^*) > u_i(x_i; x_{-i}^*)$
 130 for all $x_i \neq x_i^*$, $i \in \mathcal{N}$. Among Nash equilibria, strict equilibria are the only ones that are “structurally
 131 robust” (in the sense that they remain invariant to small perturbations of the underlying game), so
 132 they play a particularly important role in game theory.

133 3 Regularized learning in games

134 Throughout our paper, we will consider iterative decision processes that unfold as follows:

- 135 1. At each stage $t = 1, 2, \dots$, every participating agent selects an action.
- 136 2. Agents receive a reward determined by their chosen actions and their individual payoff functions.
- 137 3. Based on this reward (or other feedback), the agents update their strategies and the process repeats.

138 In this online setting, a crucial requirement is the minimization of the players’ *regret*, i.e., the
 139 difference between a player’s cumulative payoff over time and the player’s best possible strategy in
 140 hindsight. Formally, if the players’ actions at each epoch $t = 1, 2, \dots$ are collectively drawn by the
 141 probability distribution $z_t \in \Delta(\mathcal{A})$, the *regret* of each player $i \in \mathcal{N}$ is defined as

$$\text{Reg}_i(T) = \max_{\alpha_i \in \mathcal{A}_i} \sum_{t=1}^T [u_i(\alpha_i; z_{-i,t}) - u_i(z_t)], \quad (2)$$

142 and we will say that player i has *no regret* if $\text{Reg}_i(T) = o(T)$.

143 One of the most widely used policies to achieve no-regret play is the so-called “*follow-the-regularized-*
144 *leader*” (FTRL) family of algorithms and its variants [41, 42]. For completeness, we will work
145 with a more general *regularized learning* (RL) template which allows us to simultaneously consider
146 different types of feedback, strategy sampling policies, update structures, etc. To lighten notation
147 below, we will drop the player index $i \in \mathcal{N}$ when the meaning can be inferred from the context; also,
148 to stress the distinction between “strategy-like” and “payoff-like” variables, we will write throughout
149 $\mathcal{Y}_i := \mathbb{R}^{\mathcal{A}_i}$ and $\mathcal{Y} := \prod_i \mathcal{Y}_i$ for the game’s “*payoff space*”, in direct analogy to \mathcal{X}_i and $\mathcal{X} = \prod_i \mathcal{X}_i$ for
150 the game’s *strategy space*.

151 **3.1. The regularized learning template.** The general class of *regularized learning* (RL) methods
152 that we will consider proceed in an iterative, two-stage fashion as follows:

$$\begin{aligned} \text{Aggregate payoff information:} \quad & Y_{i,t+1} = Y_{i,t} + \gamma_t \hat{v}_{i,t} \\ \text{Update choice probabilities:} \quad & X_{i,t+1} = Q_i(Y_{i,t+1}) \end{aligned} \tag{RL}$$

153 In the above:

- 154 1. $X_{i,t} \in \mathcal{X}_i$ denotes the mixed strategy of player i at time $t = 1, 2, \dots$
- 155 2. $Y_{i,t} \in \mathcal{Y}_i$ is a “score vector” that measures the performance of the player’s actions over time.
- 156 3. $Q_i: \mathcal{Y}_i \rightarrow \mathcal{X}_i$ is a “regularized best response” that maps score vectors to choice probabilities.
- 157 4. $\hat{v}_{i,t}$ is a surrogate / approximation of the mixed payoff vector $v_i(X_t)$ of player i at time t .
- 158 5. $\gamma_t > 0$ is a step-size / sensitivity parameter of the form $\gamma_t \propto 1/t^{\ell_\gamma}$ for some $\ell_\gamma \in [0, 1]$.

159 In words, at each stage of the process, every player $i \in \mathcal{N}$ observes – or otherwise estimates – a proxy
160 $\hat{v}_{i,t}$ of their individual payoff vector; subsequently, players augment their actions’ scores based on this
161 information, they select a mixed strategy via the regularized choice map Q_i , and the process repeats.
162 To streamline our presentation, we discuss in detail the precise definition of \hat{v} and Q in [Sections 3.2](#)
163 and [3.3](#) below, and we present a series of examples of (RL) in [Section 3.4](#) right after.

164 **3.2. Aggregating payoff information.** As noted above, the main idea of regularized learning is
165 to track the players’ payoff vector $v(X_t)$. Importantly, there are several different modeling choices
166 that can be made here: players may have direct access to their payoff vectors (in the full information
167 setting), or some noisy approximation obtained by an inner randomization of the algorithm (e.g.,
168 when they receive information on their pure actions); they may have to recreate their payoff vectors
169 altogether (as in the bandit setting), or their estimates may be based on a strategy other than the one
170 they actually played (as in the case of optimistic algorithms). In all these cases, the surrogate vector
171 \hat{v}_t can be written concisely as

$$\hat{v}_t = v(X_t) + U_t + b_t \tag{3}$$

172 where $b_t = \mathbb{E}[\hat{v}_t | \mathcal{F}_t] - v(X_t)$ and $U_t = \hat{v}_t - \mathbb{E}[\hat{v}_t | \mathcal{F}_t]$ respectively denote the offset and the random
173 error of \hat{v}_t relative to $v(X_t)$. To streamline our presentation, we will also assume that $\|b_t\| = \mathcal{O}(1/t^{\ell_b})$
174 and $\|U_t\| = \mathcal{O}(t^{\ell_\sigma})$ for some $\ell_b, \ell_\sigma \geq 0$; we discuss the specifics of these bounds later in the paper.

175 **3.3. From scores to strategies.** Regarding the “scores-to-strategies” step of (RL), we will follow
176 the classical approach of Shalev-Shwartz [41] and assume that each player is employing a *choice*
177 *map* – or *regularized best response* – of the general form

$$Q_i(y_i) = \arg \max_{x_i \in \mathcal{X}_i} \{ \langle y_i, x_i \rangle - h_i(x_i) \} \quad \text{for all } y_i \in \mathcal{Y}_i. \tag{4}$$

178 In the above, the *regularizer* $h_i: \mathcal{X}_i \rightarrow \mathbb{R}$ acts as a penalty that smooths out the “hard” argmax
179 correspondence $y_i \mapsto \arg \max_{x_i \in \mathcal{X}_i} \langle y_i, x_i \rangle$. Accordingly, instead of following the “leader” (i.e.,
180 playing the strategy with the highest propensity score), players follow the “regularized leader” – that
181 is, they allow for a certain degree of uncertainty in their choice of strategy [9, 28, 41, 42].

182 To ease notation, we will work with kernelized regularizers of the form $h_i(x_i) = \sum_{\alpha_i \in \mathcal{A}_i} \theta(x_i \alpha_i)$ for
183 some continuous function $\theta: [0, 1] \rightarrow \mathbb{R}$ with $\inf_{z \in (0,1)} \theta''(z) > 0$. We will also say that the players’
184 regularizers are *steep* if $\lim_{z \rightarrow 0^+} \theta'(z) = -\infty$, and non-steep otherwise.

185 **Example 3.1.** A standard family of kernelized regularizers is given by $\theta(z) = z^\rho / [\rho(\rho - 1)]$ for
186 $\rho \in (0, 1) \cup (1, 2]$ and $\theta(z) = z \log z$ for $\rho = 1$ [9, 26, 28, 49]. This family includes:

- 187 • For $\rho = 2$, the quadratic regularizer $\theta(z) = z^2/2$, which yields the Euclidean projection map

$$Q_i(y_i) = \Pi_{\mathcal{X}_i}(y_i) \equiv \arg \min_{x_i \in \mathcal{X}_i} \|y_i - x_i\|_2. \tag{5}$$

188 • For $\rho = 1$, the *entropic regularizer* $\theta(z) = z \log z$, which induces the *logit choice map*

$$Q_i(y_i) = \Lambda_i(y_i) \equiv \frac{(\exp(y_i \alpha_i))_{\alpha_i \in \mathcal{A}_i}}{\sum_{\alpha_i \in \mathcal{A}_i} \exp(y_i \alpha_i)} \quad (6)$$

189 • For $\rho = 1/2$, the *fractional power regularizer* $\theta(z) = -4\sqrt{z}$ that underlies the TSALLIS-INF
190 algorithm of [1, 49] (see also Section 3.4 below). ♦

191 **3.4. Specific algorithms.** We now proceed to discuss some archetypal examples of (RL).

192 **Algorithm 1** (Follow the regularized leader). The standard “*follow-the-regularized-leader*” (FTRL)
193 method of Shalev-Shwartz & Singer [42] is obtained when players observe their full payoff vectors,
194 that is, $\hat{v}_{i,t} = v_i(X_t)$. In this case, (RL) boils down to the deterministic update rule

$$Y_{i,t+1} = Y_{i,t} + \gamma_t v_i(X_t) \quad X_{i,t+1} = Q_i(Y_{i,t+1})$$

195 or, more explicitly

$$X_{i,t+1} = \arg \max_{x_i \in \mathcal{X}_i} \left\{ \sum_{s=1}^t \gamma_s u_i(x_i; X_{-i,s}) - h_i(x_i) \right\} \quad (\text{FTRL})$$

196 For a detailed discussion of (FTRL), see [9, 26, 41]. We only note here that, as a special case,
197 when (FTRL) is run with the logit choice setup of Eq. (6), a standard calculation yields the seminal
198 *exponential/multiplicative weights* algorithm – or HEDGE [4, 27, 46] – namely

$$X_{i\alpha_i,t+1} = \frac{X_{i\alpha_i,t} \exp(\gamma_t u_i(\alpha_i; X_{-i,t}))}{\sum_{\alpha'_i \in \mathcal{A}_i} X_{i\alpha'_i,t} \exp(\gamma_t u_i(\alpha'_i; X_{-i,t}))} \quad (\text{HEDGE})$$

199 For an appetizer to the literature on (HEDGE), see [2, 9, 10, 26, 41] and references therein. ♦

200 **Algorithm 2** (Optimistic FTRL). A notable variant of FTRL – originally due to Popov [35] and
201 subsequently popularized by Rakhlin & Sridharan [36, 37] – is the so-called *optimistic FTRL* method.
202 This scheme employs an “optimistic” correction intended to anticipate future steps, and it updates as

$$Y_{i,t+1} = Y_{i,t} + \gamma_t [2v_i(X_t) - v_i(X_{t-1})] \quad (\text{Opt-FTRL})$$

203 with $X_{i,t} = Q_i(Y_{i,t})$. As a special case, if (Opt-FTRL) is run with the logit choice map (6), we obtain
204 the familiar update rule known as *optimistic multiplicative weights* (OMW) [13, 36, 37, 43].

205 Compared to (FTRL), the gain vector $\hat{v}_t = 2v(X_t) - v(X_{t-1})$ of (Opt-FTRL) has offset $b_t = v(X_t) -$
206 $v(X_{t-1})$ relative to $v(X_t)$. Thus, even though (Opt-FTRL) assumes full access to the players’ mixed
207 payoff vectors, it uses this information differently than (FTRL): in particular, the offset of (Opt-FTRL)
208 is non-zero *by design*, not because of some systematic error in the payoff measurement process. ♦

209 Now, up to this point, we have not detailed how players might observe their full, mixed payoff
210 vectors. This assumption simplifies the analysis immensely, but it is not realistic in applications to
211 e.g., online advertising and network science, where players may only be able to observe their realized
212 payoffs, and have no information about the strategies of other players or actions they did not play.
213 On that account, we describe below a range of *payoff-based* policies where players estimate their
214 counterfactual, “what-if” payoffs *indirectly*.

215 The most common way to achieve this is via the *importance-weighted estimator*

$$\text{IWE}_{i\alpha_i}(x) = \frac{\mathbb{1}\{\hat{\alpha}_i = \alpha_i\}}{x_{i\alpha_i}} u_i(\hat{\alpha}) \quad \text{for all } \alpha_i \in \mathcal{A}_i, i \in \mathcal{N}, \quad (\text{IWE})$$

216 where $x \in \mathcal{X}$ is the players’ strategy profile, and $\hat{\alpha} \in \mathcal{A}$ is drawn according to x . This estimator is at
217 the heart of the online learning literature [9, 10, 26, 41] and it leads to the following methods:

218 **Algorithm 3** (Bandit FTRL). Plugging (IWE) directly into (RL) yields the *bandit FTRL* policy

$$Y_{i,t+1} = Y_{i,t} + \gamma_t \text{IWE}_i(\hat{X}_t) \quad X_{i,t+1} = Q_i(Y_{i,t+1}) \quad (\text{B-FTRL})$$

219 where (IWE) is sampled at the mixed strategy profile

$$\hat{X}_{i,t} = (1 - \delta_t) X_{i,t} + \delta_t \text{unif}_{\mathcal{A}_i} \quad (7)$$

220 for some “explicit exploration” parameter $\delta_t \propto 1/t^{\ell_\delta}$, $\ell_\delta > 0$, which specifies the mix between $X_{i,t}$
221 and the uniform distribution $\text{unif}_{\mathcal{A}_i}$ on \mathcal{A}_i . As we discuss in the sequel, this combination of (IWE)
222 with the explicit exploration mechanism (7) means that the surrogate payoff vector $\hat{v}_t = \text{IWE}(\hat{X}_t)$
223 used to update (B-FTRL) has offset and noise bounded respectively as $b_t = \mathcal{O}(\delta_t)$ and $U_t = \mathcal{O}(1/\delta_t)$.

224 Two special cases of (B-FTRL) that have attracted significant attention in the literature are:

- 225 1. The *exponential weights algorithm for exploration and exploitation* (EXP3) [5, 10, 26], obtained
 226 by running (B-FTRL) with the logit choice map (6).
 227 2. The *Tsallis implicitly normalized forecaster* (TSALLIS-INF) [1, 3, 48, 49] that was proposed as a
 228 more efficient alternative to EXP3, and which updates as

$$X_{i,t} = \arg \max_{x_i \in \mathcal{X}_i} \left\{ \langle Y_{i,t}, x_i \rangle + 4 \sum_{\alpha_i \in \mathcal{A}_i} \sqrt{x_i \alpha_i} \right\} \quad (\text{TSALLIS-INF})$$

229 i.e., as (B-FTRL) with the fractional power regularizer $\theta(z) = -4\sqrt{z}$ of Example 3.1. \blacklozenge

230 For illustration purposes, we provide some more examples of (RL) in Appendix B.

231 4 First results: resilience to strategic deviations

232 We are now in a position to begin our analysis of the rationality properties of the players' long-
 233 run behavior under (RL). To that end, we should first note that no-regret play may *still* lead to
 234 counterintuitive and highly non-rationalizable outcomes, e.g., with all players selecting dominated
 235 strategies for all time. The example below is adapted from Viossat & Zapechelnyuk [45]:

236 **Example 4.1.** Consider the 4×4 symmetric 2-player game with payoff bimatrix

	A	B	C	D
A	(1, 1)	(1, 2/3)	(0, 0)	(0, -1/3)
B	(2/3, 1)	(2/3, 2/3)	(-1/3, 0)	(-1/3, -1/3)
C	(0, 0)	(0, -1/3)	(1, 1)	(1, 2/3)
D	(-1/3, 0)	(-1/3, -1/3)	(2/3, 1)	(2/3, 2/3)

237 In this game, B and D are strictly dominated for both players by their stronger “twins” (A and C
 238 respectively). However, it is easy to check that if both players choose between (B, B) and (D, D) with
 239 probability 1/2 each, the resulting distribution of play $z \in \Delta(\mathcal{A})$ satisfies $u_i(\alpha_i; z_{-i}) - u_i(z) \leq -1/6$
 240 for all $\alpha_i \in \{A, B, C, D\}$, $i = 1, 2$. As a result, the players' regret under $z_t \equiv z$ is *negative*, even
 241 though both players play strictly dominated strategies at all times. \blacklozenge

242 The example above shows that the no-regret property does not suffice to exclude non-rationalizable
 243 outcomes by itself. In addition, it also shows that predictions based on correlated play are not always
 244 appropriate for describing the players' behavior under (RL): the end-state of any regularized learning
 245 algorithm will be a closed connected set of mixed strategies, so it is not possible to play *only* (B, B)
 246 or (D, D) in the long run. We are thus led to the following natural question: *What are the rationality*
 247 *properties of long-run play under (RL)? Is the players' behavior robust to strategic deviations?*

248 To study this question formally, we will focus on the *limit set* $\mathcal{L}(X)$ of X_t under (RL), viz.

$$\mathcal{L}(X) := \bigcap_t \text{cl}\{X_s : s \geq t\} \equiv \{\hat{x} \in \mathcal{X} : X_{t_k} \rightarrow \hat{x} \text{ for some subsequence } X_{t_k} \text{ of } X_t\}. \quad (8)$$

249 In words, $\mathcal{L}(X)$ is the set of limit points of X_t or, equivalently, the *smallest* subset of \mathcal{X} to which
 250 X_t converges. Clearly, the simplest instance of a limit set is when $\mathcal{L}(X)$ is a singleton, i.e., when
 251 X_t converges to a point. This case has attracted significant interest in the literature: for example, if
 252 $\mathcal{L}(X) = \{x^*\}$ then, for certain special cases of (RL), it is known that x^* is a Nash equilibrium of Γ
 253 [29]. However, beyond this relatively simple regime, the structure of the limit sets of (RL) could be
 254 arbitrarily complicated and their rationality properties are not well-understood.

255 With this in mind, as a first attempt to study whether the long-run behavior of (RL) is “robust to
 256 strategic deviations”, we will consider the notion of *resilience*, as defined below:

257 **Definition 1.** A closed subset \mathcal{S} of \mathcal{X} is *resilient to strategic deviations* – or simply *resilient* – if, for
 258 every deviation $x_i \in \mathcal{X}_i$ of every player $i \in \mathcal{N}$, we have $u_i(x^*) \geq u_i(x_i, x_{-i}^*)$ for some $x^* \in \mathcal{S}$.

259 Informally, \mathcal{S} is resilient if every unilateral deviation from \mathcal{S} is deterred by some (possibly different)
 260 element thereof. In particular, if \mathcal{S} is a singleton, we immediately recover the definition of a Nash
 261 equilibrium; beyond this base case, other examples include the set of undominated strategies of a
 262 game, the support face of the equilibria of two-player zero-sum games, etc.

263 Importantly, as we show below, the limit sets of (RL) are almost surely resilient *in all games*:

264 **Theorem 1.** *Let X_t , $t = 1, 2, \dots$, be the sequence of play generated by (RL) with step-size/gain*
 265 *parameters $\ell_\gamma > 2\ell_\sigma$ and $\ell_b > 0$. Then, with probability 1, the limit set $\mathcal{L}(X)$ of X_t is resilient.*

266 *Proof sketch.* The proof of [Theorem 1](#) boils down to two interleaved arguments that we detail in
267 [??](#). The first hinges on showing that, if $\mathbb{P}(\mathcal{L}(X) = \mathcal{S}) > 0$ for some *non-random* $\mathcal{S} \subseteq \mathcal{X}$, \mathcal{S}
268 must be resilient. This is argued by contradiction: if $p_i \in \mathcal{X}_i$ is a unilateral deviation violating
269 [Definition 1](#), we must also have $\liminf_{t \rightarrow \infty} [u_i(p_i; X_{-i,t}) - u_i(X_t)] > 0$ with positive probability.
270 However, the existence of a strategy that consistently outperforms X_t runs contrary to the fact that
271 strategies that [\(RL\)](#) selects against underperforming strategies. We make this intuition precise via
272 an energy argument that leverages a series of results from martingale limit theory (which is where
273 the requirements for γ_t , b_t and U_t come in). Then, to get the stronger statement that the *random* set
274 $\mathcal{L}(X)$ is resilient w.p.1, we show that the above remains true if p_i is replaced by a deviation q_i which
275 is close enough to p_i and has *rational* entries. Since there is a countable number of such profiles,
276 we can use a union bound on an enumeration of the rationals to isolate a deviation witnessing the
277 negation of [Definition 1](#) and apply our argument for non-random sets to conclude our proof. ■

278 [Theorem 1](#) is our first universal guarantee for [\(RL\)](#), so some remarks are in order. First, we should
279 point out that the requirements $\ell_b > 0$ and $2\ell_\sigma < \ell_\gamma$ are a priori *implicit* because they depend on the
280 offset and magnitude statistics of the feedback sequence \hat{v}_t . However, in most learning algorithms,
281 these quantities are under the *explicit* control of the players: for example, as we show in [Appendix B](#),
282 [Algorithm 2](#) has $\ell_b = \ell_\gamma$ while, for [Algorithm 3](#), we have $\ell_b = \ell_\sigma = \ell_\delta$. In this way, when instantiated
283 to [Algorithms 1–3](#) (and special cases thereof), [Theorem 1](#) yields the following corollary:

284 **Corollary 1.** *Suppose that [Algorithms 1–3](#) are run with $\ell_\gamma \in (0, 1]$ and, for [Algorithm 3](#), $\ell_\delta \in$
285 $(0, \ell_\gamma/2)$. Then, with probability 1, the limit set $\mathcal{L}(X)$ of X_t is resilient.*

286 Now, since [Theorem 1](#) applies to all games, it would seem to provide a universally positive answer
287 to whether [\(RL\)](#) is robust to strategic deviations. However, this is not so: a direct calculation shows
288 that the face of \mathcal{X} that is spanned by the dominated strategies (B, B) and (D, D) of [Example 4.1](#)
289 is resilient, so [Theorem 1](#) cannot exclude convergence to a set where dominated strategies survive.
290 Thus, just like no-regret play, the notion of resilience does not suffice by itself to capture the idea
291 of rational behavior. This is because, albeit natural, resilience is too lax to provide a meaningful
292 link between robustness to unilateral deviations – a *game-theoretic* requirement – and stability under
293 regularized learning – a *dynamic* requirement. We address this question in detail in the next section.

294 5 A characterization of strategic stability under regularized learning

295 Similar to the set of pure strategies that arise from no-regret play, the main limitation of resilience
296 is that a payoff-improving deviation may be countered by an action profile where the deviator also
297 switched to a *different* strategy; in other words, resilience is not a *self-enforcing* barrier to deviations.
298 In view of this, we will focus below on a much more stringent criterion of strategic stability, namely
299 that *any* deviation from the set in question incur a cost to the deviating agent.

300 **Club sets.** The above idea can be made precise as follows: First, define the *better-reply correspon-*
301 *dence* of player $i \in \mathcal{N}$ as $\text{btr}_i(x) = \{x'_i \in \mathcal{X}_i : u_i(x'_i; x_{-i}) \geq u_i(x)\}$, and write $\text{btr} = \prod_i \text{btr}_i$ for
302 the product correspondence $\text{btr}(x) = \text{btr}_1(x) \times \cdots \times \text{btr}_N(x)$. [In words, btr_i assigns to each
303 $x \in \mathcal{X}$ those strategies of player i that are (weakly) better against x than x_i .] In addition, given a
304 product of pure strategies $\mathcal{C} = \prod_{i \in \mathcal{N}} \mathcal{C}_i$ with $\mathcal{C}_i \subseteq \mathcal{A}_i$ for all $i \in \mathcal{N}$, let $\mathcal{S} = \Delta(\mathcal{C})$ denote the span of
305 \mathcal{C} , and let $\mathcal{P}(\mathcal{X})$ denote the collection of all such sets. We then say that $\mathcal{S} \in \mathcal{P}(\mathcal{X})$ is *closed under*
306 *better replies* – a *club set* for short – if it is closed under btr , i.e., $\text{btr}(\mathcal{S}) \subseteq \mathcal{S}$; finally, \mathcal{S} is said to
307 be *minimally club* (m-club) if it does not admit a proper club subset.

308 Of course, the entire strategy space \mathcal{X} is closed under better replies so, a priori, club sets could also
309 contain dominated strategies and/or other non-rationalizable outcomes. By contrast, *minimal* club
310 sets are much more rigid in their relation to rational behavior because any unilateral deviation from
311 an m-club set is *costly*, and m-club sets are *minimal* in this regard. On that account, m-club sets can
312 be seen as *the closest setwise analogue to strict Nash equilibria*.

313 This analogy is accentuated further by the following properties of m-club sets (all due to Ritzberger
314 & Weibull [[38](#)], who introduced the concept):

- 315 1. Every game admits an m-club set; and if this set is a singleton, then it is a *strict* Nash equilibrium.

- 316 2. Any m -club set \mathcal{S} is *fixed* under better replies, that is, $\text{btr}(\mathcal{S}) = \mathcal{S}$ (implying in turn that \mathcal{S} cannot
317 contain any dominated strategies, including iteratively dominated ones).
318 3. Any m -club set \mathcal{S} contains an *essential equilibrium component*, i.e., a component of Nash equilibria
319 such that every small perturbation of the game admits a nearby equilibrium; in addition, this
320 component has *full support* on \mathcal{S} , i.e., it employs all pure strategy profiles that lie in \mathcal{S} .¹

321 Going back to our online learning setting, the above leads to the following natural set of questions:

322 *Are club sets (minimal or not) stable under the dynamics of regularized learning?*
323 *Are they attracting? And, if so, are they the only such sets?*

324 Any answer to these questions – positive or negative – would be an important step in delineating the
325 relation between *strategic stability* (in the above sense) and *dynamic stability* under (RL). To that
326 end, we start by formalizing some notions of dynamic stability that will be central in the sequel:

327 **Definition 2.** Fix some subset \mathcal{S} of \mathcal{X} and a tolerance level $\epsilon > 0$. We then say that \mathcal{S} is:

328 1. *Stochastically stable* if, for every neighborhood \mathcal{U} of \mathcal{S} in \mathcal{X} , there exists a neighborhood \mathcal{U}_1 of \mathcal{S}
329 such that

$$\mathbb{P}(X_t \in \mathcal{U} \text{ for all } t = 1, 2, \dots) \geq 1 - \epsilon \quad \text{whenever } X_1 \in \mathcal{U}_1. \quad (9)$$

330 2. *Stochastically attracting* if there exists a neighborhood \mathcal{U}_1 of \mathcal{S} such that

$$\mathbb{P}(\lim_{t \rightarrow \infty} \text{dist}(X_t, \mathcal{S}) = 0) \geq 1 - \epsilon \quad \text{whenever } X_1 \in \mathcal{U}_1. \quad (10)$$

331 3. *Stochastically asymptotically stable* if it is stochastically stable and attracting.

332 4. *Irreducibly stable* if \mathcal{S} is stochastically asymptotically stable and it does not admit a strictly smaller
333 stochastically asymptotically subset \mathcal{S}' with $\text{supp}(\mathcal{S}') \subsetneq \text{supp}(\mathcal{S})$.

334 With all this in hand, our main result below provides a sharp characterization of strategic stability in
335 the context of regularized learning:

336 **Theorem 2.** Fix some set $\mathcal{S} \in \mathcal{P}(\mathcal{X})$ and suppose that (RL) is run with a steep regularizer and
337 step-size/gain parameters $\ell_\gamma \in [0, 1]$, $\ell_b > 0$, and $\ell_\sigma < 1/2$. Then:

- 338 1. \mathcal{S} is stochastically asymptotically stable under (RL) if and only if it is a club set.
339 2. \mathcal{S} is irreducibly stable under (RL) if and only if it is an m -club set.

340 In addition, we also get the following convergence rate estimates for club sets:

341 **Theorem 3.** Let $\mathcal{S} \in \mathcal{P}(\mathcal{X})$ be a club set, and let X_t , $t = 1, 2, \dots$, be the sequence of play generated
342 by (RL) with parameters $\ell_\gamma \in [0, 1]$, $\ell_b > 0$, and $\ell_\sigma < 1/2$. Then, for all $\epsilon > 0$, there exists an (open,
343 unbounded) initialization domain $\mathcal{D} \subseteq \mathcal{Y}$ such that, with probability at least $1 - \epsilon$, we have

$$\text{dist}(X_t, \mathcal{S}) \leq C\varphi\left(c_1 - c_2 \sum_{s=1}^t \gamma_s\right) \quad \text{whenever } Y_1 \in \mathcal{D} \quad (11)$$

344 where C, c_1, c_2 are constants ($C, c_2 > 0$), and the rate function φ is given by $\varphi(z) = (\theta')^{-1}(z)$ if
345 $z > \lim_{z \rightarrow 0^+} \theta'(z)$, and $\varphi(z) = 0$ otherwise.

346 Specifically, if we instantiate Theorem 3 to Algorithms 1–3, we get the explicit estimates:

347 **Corollary 2.** Suppose that Algorithms 1–3 are run with $\ell_\gamma \in [0, 1]$ and, for Algorithm 3, $\ell_\delta \in$
348 $(0, 1/2)$. Then, with notation as in Theorem 3, X_t converges to \mathcal{S} at a rate of

$$\text{dist}(X_t, \mathcal{S}) \leq C \cdot \begin{cases} [1 - c \sum_{s=1}^t \gamma_s]_+ & \text{if } \theta(z) = z^2/2 & \# \text{quadratic regularization} \\ \exp(-c \sum_{s=1}^t \gamma_s) & \text{if } \theta(z) = z \log z & \# \text{entropic regularization} \\ 1 / (c + \sum_{s=1}^t \gamma_s)^2 & \text{if } \theta(z) = -4\sqrt{z} & \# \text{fractional regularization} \end{cases} \quad (12)$$

349 for positive constants $C, c > 0$. In particular, the projection-based variants of Algorithms 1–3
350 converge to m -club sets in a *finite* number of steps.

¹Formally, a component \mathcal{X}^* of Nash equilibria of Γ is *essential* if, for all $\epsilon > 0$, there exists $\delta > 0$ such that any perturbation of the payoffs of Γ by at most δ produces a Nash equilibrium that is ϵ -close to \mathcal{X}^* [44]. This property – known as “*essentiality*” – has a long history as one of the strictest setwise solution refinements in game theory; in particular, it satisfies all the seminal *strategic stability* requirements of Kohlberg & Mertens [25], including robustness to strategic payoff perturbations. For an in-depth discussion, see van Damme [44].

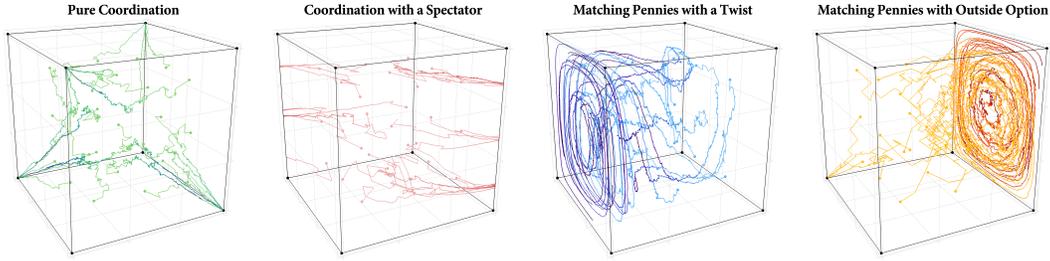


Figure 1: The long-run behavior of EXP3 (Algorithm 3) in four representative $2 \times 2 \times 2$ games. In all cases, the dynamics converge to m-club sets, either *strict equilibria* themselves, or spanning an *essential component* of Nash equilibria. The details of the numerics and the games being played are provided in the appendix.

351 *Proof sketch.* The proof of Theorems 2 and 3 is quite involved so we defer it to Appendix D. At
 352 a high level, it hinges on constructing a family of “primal-dual” energy functions, one per pure
 353 deviation from the set \mathcal{S} under study. If unilateral deviations from \mathcal{S} incur a cost to the deviator (that
 354 is, if \mathcal{S} is club), these energy functions can be “bundled together” to produce a suitable Lyapunov-like
 355 function for \mathcal{S} . In more detail, the minimization of each individual energy function implies that the
 356 score variable Y_t of (RL) diverges along an “astral direction” in the payoff space \mathcal{Y} – i.e., it escapes
 357 to infinity along the interior of a certain convex cone of \mathcal{Y} [16]. Because this minimization occurs at
 358 infinity, the aggregation of offsets and random errors in (RL) affords some extra “wobble room” in
 359 our martingale analysis, so we are able to show that $X_t = Q(Y_t)$ remains close to \mathcal{S} under a much
 360 wider range of parameters compared to Theorem 1. Then, a series of convex analysis arguments in
 361 the spirit of [28] coupled with the definition of Q allows us to show that the escape of Y_t along the
 362 intersection of all these cones implies convergence to \mathcal{S} at the specified rate.

363 On the converse side, if an asymptotically stable set is not club, we can find a non-costly (and possibly
 364 profitable) deviation z from \mathcal{S} which is selected against by (RL). However, this extinction runs
 365 contrary to the reinforcement of better replies under (RL), an argument which can be made precise
 366 by applying the martingale law of large numbers to $\langle Y_t, z \rangle$ [21]. The irreducible stability of m-club
 367 sets then follows by invoking this criterion reductively for any potentially stable subset \mathcal{S}' of \mathcal{S} . ■

368 **Discussion and remarks.** Theorems 2 and 3 are our main results linking dynamic and strategic
 369 stability, so we conclude with a series of remarks. First, we should note that Theorem 2 can be
 370 summed up as follows: *a product of pure strategies is (minimally) closed under better replies if and*
 371 *only if its span is (irreducibly) stable under regularized learning.* Importantly, this equivalence is
 372 based solely on the game’s payoff data: it does not depend on the specific choices underlying (RL),
 373 including the choice map employed by each player, whether some players are using an optimistic
 374 adjustment or not, if they have access to their full payoff vectors, etc. As such, this equivalence
 375 provides a crisp operational criterion for identifying which pure strategy combinations ultimately
 376 persist under regularized learning – and, via Theorem 3, *how fast* this identification takes place.

377 In this light, Theorem 2 essentially states that the only robust prediction that can be made for
 378 the outcome of a regularized learning process is (minimal) closedness under better replies. This
 379 interpretation has significant cutting power for the emergence of rational behavior. To begin, in terms
 380 of equilibrium play, it effortlessly implies that a pure strategy profile is stochastically asymptotically
 381 stable under (RL) if and only if it is a strict Nash equilibrium. A version of this equivalence was
 382 only recently proved in [17] and [19] (in continuous and discrete time respectively), so Theorem 2
 383 can be seen as a far-reaching generalization of these recent results. More to the point, since every
 384 m-club set \mathcal{S} contains an essential equilibrium component that is fully supported in \mathcal{S} , Theorem 2
 385 also provides an important link between dynamic and structural stability: if an equilibrium – or a
 386 component of equilibria – is not robust to perturbations of the underlying game, *it cannot be robustly*
 387 *identified by a regularized learning process* (and vice versa). This remark is of particular importance
 388 for extensive-form games as such games often have non-generic equilibrium components that cannot
 389 be treated otherwise by the existing theory.

390 Finally, we should stress that Theorems 2 and 3 guarantee convergence even with a constant step-size.
 391 Together with the finite-time convergence guarantees of Corollary 2 for projection-based methods,
 392 this feature is a testament to the robustness of club sets as, in the presence of uncertainty, convergence
 393 almost always requires a vanishing step-size which can slow convergence down to a crawl. We find
 394 this robust convergence landscape particularly intriguing for future research on the topic.

- 396 [1] Abernethy, J., Lee, C., and Tewari, A. Fighting bandits with a new kind of smoothness. In *NIPS '15: Proceedings of the 29th International Conference on Neural Information Processing Systems*, 2015.
- 397
- 398 [2] Arora, S., Hazan, E., and Kale, S. The multiplicative weights update method: A meta-algorithm and applications. *Theory of Computing*, 8(1):121–164, 2012.
- 399
- 400 [3] Audibert, J.-Y. and Bubeck, S. Regret bounds and minimax policies under partial monitoring. *Journal of Machine Learning Research*, 11:2635–2686, 2010.
- 401
- 402 [4] Auer, P., Cesa-Bianchi, N., Freund, Y., and Schapire, R. E. Gambling in a rigged casino: The adversarial multi-armed bandit problem. In *Proceedings of the 36th Annual Symposium on Foundations of Computer Science*, 1995.
- 403
- 404
- 405 [5] Auer, P., Cesa-Bianchi, N., Freund, Y., and Schapire, R. E. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.
- 406
- 407 [6] Basu, K. and Weibull, J. W. Strategy subsets closed under rational behavior. *Economics Letters*, 36: 141–146, 1991.
- 408
- 409 [7] Bauschke, H. H., Borwein, J. M., and Combettes, P. L. Bregman monotone optimization algorithms. *SIAM Journal on Control and Optimization*, 42(2):596–636, 2003.
- 410
- 411 [8] Benaïm, M. Dynamics of stochastic approximation algorithms. In Azéma, J., Émery, M., Ledoux, M., and Yor, M. (eds.), *Séminaire de Probabilités XXXIII*, volume 1709 of *Lecture Notes in Mathematics*, pp. 1–68. Springer Berlin Heidelberg, 1999.
- 412
- 413
- 414 [9] Bubeck, S. and Cesa-Bianchi, N. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012.
- 415
- 416 [10] Cesa-Bianchi, N. and Lugosi, G. *Prediction, Learning, and Games*. Cambridge University Press, 2006.
- 417 [11] Chotibut, T., Falniowski, F., Misiurewicz, M., and Piliouras, G. Family of chaotic maps from game theory. *Dynamical Systems*, 2020.
- 418
- 419 [12] Cressman, R. *Evolutionary Dynamics and Extensive Form Games*. The MIT Press, 2003.
- 420 [13] Daskalakis, C. and Panageas, I. Last-iterate convergence: Zero-sum games and constrained min-max optimization. In *ITCS '19: Proceedings of the 10th Conference on Innovations in Theoretical Computer Science*, 2019.
- 421
- 422
- 423 [14] Daskalakis, C., Ilyas, A., Syrgkanis, V., and Zeng, H. Training GANs with optimism. In *ICLR '18: Proceedings of the 2018 International Conference on Learning Representations*, 2018.
- 424
- 425 [15] Demichelis, S. and Ritzberger, K. From evolutionary to strategic stability. *Journal of Economic Theory*, 113:51–75, 2003.
- 426
- 427 [16] Dudík, M., Schapire, R. E., and Telgarsky, M. Convex analysis at infinity: An introduction to astral space. <https://arxiv.org/abs/2205.03260>, 2022.
- 428
- 429 [17] Flokas, L., Vlatakis-Gkaragkounis, E. V., Lianas, T., Mertikopoulos, P., and Piliouras, G. No-regret learning and mixed Nash equilibria: They do not mix. In *NeurIPS '20: Proceedings of the 34th International Conference on Neural Information Processing Systems*, 2020.
- 430
- 431
- 432 [18] Fudenberg, D. and Tirole, J. *Game Theory*. The MIT Press, 1991.
- 433 [19] Giannou, A., Vlatakis-Gkaragkounis, E. V., and Mertikopoulos, P. Survival of the strictest: Stable and unstable equilibria under regularized learning with partial information. In *COLT '21: Proceedings of the 34th Annual Conference on Learning Theory*, 2021.
- 434
- 435
- 436 [20] Gidel, G., Berard, H., Vignoud, G., Vincent, P., and Lacoste-Julien, S. A variational inequality perspective on generative adversarial networks. In *ICLR '19: Proceedings of the 2019 International Conference on Learning Representations*, 2019.
- 437
- 438
- 439 [21] Hall, P. and Heyde, C. C. *Martingale Limit Theory and Its Application*. Probability and Mathematical Statistics. Academic Press, New York, 1980.
- 440
- 441 [22] Hannan, J. Approximation to Bayes risk in repeated play. In Dresher, M., Tucker, A. W., and Wolfe, P. (eds.), *Contributions to the Theory of Games, Volume III*, volume 39 of *Annals of Mathematics Studies*, pp. 97–139. Princeton University Press, Princeton, NJ, 1957.
- 442
- 443
- 444 [23] Hofbauer, J. and Sigmund, K. Evolutionary game dynamics. *Bulletin of the American Mathematical Society*, 40(4):479–519, July 2003.
- 445
- 446 [24] Juditsky, A., Nemirovski, A. S., and Tauvel, C. Solving variational inequalities with stochastic mirror-prox algorithm. *Stochastic Systems*, 1(1):17–58, 2011.
- 447
- 448 [25] Kohlberg, E. and Mertens, J.-F. On the strategic stability of equilibria. *Econometrica*, 54(5):1003–1037, September 1986.
- 449
- 450 [26] Lattimore, T. and Szepesvári, C. *Bandit Algorithms*. Cambridge University Press, Cambridge, UK, 2020.

- 451 [27] Littlestone, N. and Warmuth, M. K. The weighted majority algorithm. *Information and Computation*, 108
452 (2):212–261, 1994.
- 453 [28] Mertikopoulos, P. and Sandholm, W. H. Learning in games via reinforcement and regularization. *Mathe-*
454 *matics of Operations Research*, 41(4):1297–1324, November 2016.
- 455 [29] Mertikopoulos, P. and Zhou, Z. Learning in games with continuous action sets and unknown payoff
456 functions. *Mathematical Programming*, 173(1-2):465–507, January 2019.
- 457 [30] Mertikopoulos, P., Papadimitriou, C. H., and Piliouras, G. Cycles in adversarial regularized learning. In
458 *SODA '18: Proceedings of the 29th annual ACM-SIAM Symposium on Discrete Algorithms*, 2018.
- 459 [31] Mertikopoulos, P., Lecouat, B., Zenati, H., Foo, C.-S., Chandrasekhar, V., and Piliouras, G. Optimistic
460 mirror descent in saddle-point problems: Going the extra (gradient) mile. In *ICLR '19: Proceedings of the*
461 *2019 International Conference on Learning Representations*, 2019.
- 462 [32] Nemirovski, A. S. Prox-method with rate of convergence $O(1/t)$ for variational inequalities with Lipschitz
463 continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on*
464 *Optimization*, 15(1):229–251, 2004.
- 465 [33] Palaiopoulos, G., Panageas, I., and Piliouras, G. Multiplicative weights update with constant step-size in
466 congestion games: Convergence, limit cycles and chaos. In *NIPS '17: Proceedings of the 31st International*
467 *Conference on Neural Information Processing Systems*, 2017.
- 468 [34] Piliouras, G., Sim, R., and Skoulakis, S. Optimal no-regret learning in general games: Bounded regret with
469 unbounded step-sizes via clairvoyant mwu. <https://arxiv.org/abs/2111.14737>, 2021.
- 470 [35] Popov, L. D. A modification of the Arrow–Hurwicz method for search of saddle points. *Mathematical*
471 *Notes of the Academy of Sciences of the USSR*, 28(5):845–848, 1980.
- 472 [36] Rakhlin, A. and Sridharan, K. Online learning with predictable sequences. In *COLT '13: Proceedings of*
473 *the 26th Annual Conference on Learning Theory*, 2013.
- 474 [37] Rakhlin, A. and Sridharan, K. Optimization, learning, and games with predictable sequences. In *NIPS '13:*
475 *Proceedings of the 27th International Conference on Neural Information Processing Systems*, 2013.
- 476 [38] Ritzberger, K. and Weibull, J. W. Evolutionary selection in normal-form games. *Econometrica*, 63(6):
477 1371–99, November 1995.
- 478 [39] Rockafellar, R. T. and Wets, R. J. B. *Variational Analysis*, volume 317 of *A Series of Comprehensive*
479 *Studies in Mathematics*. Springer-Verlag, Berlin, 1998.
- 480 [40] Sandholm, W. H. *Population Games and Evolutionary Dynamics*. MIT Press, Cambridge, MA, 2010.
- 481 [41] Shalev-Shwartz, S. Online learning and online convex optimization. *Foundations and Trends in Machine*
482 *Learning*, 4(2):107–194, 2011.
- 483 [42] Shalev-Shwartz, S. and Singer, Y. Convex repeated games and Fenchel duality. In *NIPS' 06: Proceedings*
484 *of the 19th Annual Conference on Neural Information Processing Systems*, pp. 1265–1272. MIT Press,
485 2006.
- 486 [43] Syrgkanis, V., Agarwal, A., Luo, H., and Schapire, R. E. Fast convergence of regularized learning in games.
487 In *NIPS '15: Proceedings of the 29th International Conference on Neural Information Processing Systems*,
488 pp. 2989–2997, 2015.
- 489 [44] van Damme, E. *Stability and perfection of Nash equilibria*. Springer-Verlag, Berlin, 1987.
- 490 [45] Viosat, Y. and Zapechelnuyk, A. No-regret dynamics and fictitious play. *Journal of Economic Theory*,
491 148(2):825–842, March 2013.
- 492 [46] Vovk, V. G. Aggregating strategies. In *COLT '90: Proceedings of the 3rd Workshop on Computational*
493 *Learning Theory*, pp. 371–383, 1990.
- 494 [47] Weibull, J. W. *Evolutionary Game Theory*. MIT Press, Cambridge, MA, 1995.
- 495 [48] Zimmert, J. and Seldin, Y. An optimal algorithm for stochastic and adversarial bandits. In *AISTATS '19:*
496 *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics*, 2019.
- 497 [49] Zimmert, J. and Seldin, Y. Tsallis-INF: An optimal algorithm for stochastic and adversarial bandits.
498 *Journal of Machine Learning Research*, 22(28):1–49, 2021.

499 **A Auxiliary results**

500 In this appendix we collect some basic properties of the regularized choice maps and some results
501 from probability theory that will be useful in the sequel.

502 **A.1. Regularized choice maps and their properties.** Throughout this appendix, we will suppress
503 the player index $i \in \mathcal{N}$, and we will follow standard conventions in convex analysis [39] that treat h
504 as an extended-real-valued function $h: \mathcal{Y} \rightarrow \mathbb{R} \cup \{\infty\}$ with $h(x) = \infty$ for all $x \in \mathcal{Y} \setminus \mathcal{X}$. With this in
505 mind, the subdifferential of a h at $x \in \mathcal{X}$ is defined as

$$\partial h(x) := \{y \in \mathcal{Y} : h(x') \geq h(x) + \langle y, x' - x \rangle \text{ for all } x' \in \mathcal{X}\}, \quad (\text{A.1})$$

506 where \mathcal{Y} denotes here the algebraic dual \mathcal{V}^* of \mathcal{V} . Accordingly, the *domain of subdifferentiability* of
507 h is $\text{dom } \partial h := \{x \in \text{dom } h : \partial h \neq \emptyset\}$, and the convex conjugate of h is defined as

$$h^*(y) = \max_{x \in \mathcal{X}} \{\langle y, x \rangle - h(x)\} \quad (\text{A.2})$$

508 for all $y \in \mathcal{Y}$. We then have the following basic results.

509 **Lemma A.1.** *Let h be a regularizer on \mathcal{X} , and let $Q: \mathcal{Y} \rightarrow \mathcal{X}$ be the induced choice map. Then:*

- 510 1. Q is single-valued, and, for all $x \in \mathcal{X}$, $y \in \mathcal{Y}$, we have $x = Q(y) \iff y \in \partial h(x)$.
- 511 2. For all $x \in \text{ri } \mathcal{X}$, we have $\partial h(x) = \{(\theta'(x_\alpha) + \mu)_{\alpha \in \mathcal{A}} : \mu \in \mathbb{R}\}$.
- 512 3. For all $y \in \mathcal{Y}$, we have $Q(y) = \nabla h^*(y)$.
- 513 4. Q is $(1/K)$ -Lipschitz continuous with $K := \inf_{(0,1]} \theta''(z)$. In particular, as a special case, the logit
514 choice map Λ is 1-Lipschitz continuous in the (L^1, L^∞) pair of norms on \mathcal{Y} and \mathcal{X} respectively.
- 515 5. If $y_\alpha - y_{\alpha'} \rightarrow -\infty$ for some $\alpha' \neq \alpha$, then $Q_\alpha(y) \rightarrow 0$.

516 *Remark.* Some of the properties presented in Lemma A.1 are well known in the literature on
517 regularized learning methods (see e.g., [28] and references therein), but we provide a proof of the
518 entire lemma for completeness. \blacklozenge

519 *Proof of Lemma A.1.* For the first property of Q , note that the maximum in (4) is attained for all
520 $y \in \mathcal{Y}$ because h is lower-semicontinuous (l.s.c.) and strongly convex. Furthermore, x solves (4) if
521 and only if $y - \partial h(x) \ni 0$, i.e., if and only if $y \in \partial h(x)$.

522 For our second claim, if $x \in \text{ri}(\mathcal{X})$, the first-order stationarity conditions for the convex problem (4)
523 that defines Q become

$$y_\alpha - \theta'(x_\alpha) = \mu \quad \text{for all } \alpha \in \mathcal{A}, \quad (\text{A.3})$$

524 because the inequality constraints $x_\alpha \geq 0$ are all inactive (recall that $x \in \text{ri}(\mathcal{X})$ by assumption). Now,
525 by the first part of the theorem we have $x = Q(y)$ if and only if $y \in \partial h(x)$, so we conclude that
526 $\partial h(x) = \{(\theta'(x_\alpha) + \mu)_{\alpha \in \mathcal{A}} : \mu \in \mathbb{R}\}$, as claimed.

527 For the fourth item, the expression $Q = \nabla h^*$ is an immediate consequence of Danskin's theorem,
528 while the Lipschitz continuity of Q follows from standard results, see e.g., [39, Theorem 12.60(b)].

529 For our last claim, let y_t be a sequence in \mathcal{Y} such that $y_{\alpha,t} - y_{\alpha',t} \rightarrow -\infty$ and let $x_t = Q(y_t)$. Then,
530 by descending to a subsequence if necessary, assume there exists some $\varepsilon > 0$ such that $x_{\alpha,t} \geq \varepsilon > 0$
531 for all t . Then, by the defining relation $Q(y) = \arg \max \{\langle y, x \rangle - h(x)\}$ of Q , we have:

$$\langle y_t, x_t \rangle - h(x_t) \geq \langle y_t, x' \rangle - h(x') \quad (\text{A.4})$$

532 for all $x' \in \mathcal{X}$. Therefore, taking $x'_t = x_t + \varepsilon(e_{\alpha'} - e_\alpha)$, we readily obtain

$$\varepsilon(y_{\alpha,t} - y_{\alpha',t}) \geq h(x_t) - h(x'_t) \geq \min h - \max h \quad (\text{A.5})$$

533 which contradicts our original assumption that $y_{\alpha,t} - y_{\alpha',t} \rightarrow -\infty$. With \mathcal{X} compact, the above shows
534 that $x_\alpha^* = 0$ for any limit point x^* of x_t , i.e. $Q_\alpha(y_t) \rightarrow 0$. \blacksquare

535 The second collection of results concerns the *Fenchel coupling*, an energy function that was first
536 introduced in [28, 29] and is defined as follows:

$$F(p, y) = h(p) + h^*(y) - \langle y, p \rangle \quad \text{for all } p \in \mathcal{X} \text{ and } y \in \mathcal{Y}. \quad (\text{A.6})$$

537 This coupling will play a major role in the proofs of Theorem 1, so we prove two of its most basic
538 properties below.

539 **Lemma A.2.** For all $p \in \mathcal{X}$ and all $y, y' \in \mathcal{Y}$, we have:

$$a) \quad F(p, y) \geq \frac{1}{2}K \|Q(y) - p\|^2. \quad (\text{A.7a})$$

$$b) \quad F(p, y') \leq F(p, y) + \langle y' - y, Q(y) - p \rangle + \frac{1}{2K} \|y' - y\|_\infty^2. \quad (\text{A.7b})$$

540 In particular, if $h(0) = 0$, we have

$$(K/2)\|Q(y)\|^2 \leq h^*(y) \leq -\min h + \langle y, Q(y) \rangle + (2/K)\|y\|_\infty^2 \quad \text{for all } y \in \mathcal{Y}. \quad (\text{A.8})$$

541 *Proof of Lemma A.2.* By the strong convexity of h relative to $\|\cdot\|$ (cf. Lemma A.1), we have

$$\begin{aligned} h(x) + t\langle y, p - x \rangle &\leq h(x + t(p - x)) \\ &\leq th(p) + (1 - t)h(x) - \frac{1}{2}Kt(1 - t)\|x - p\|^2, \end{aligned} \quad (\text{A.9})$$

542 leading to the bound

$$\frac{1}{2}K(1 - t)\|x - p\|^2 \leq h(p) - h(x) - \langle y, p - x \rangle = F(p, y) \quad (\text{A.10})$$

543 for all $t \in (0, 1]$. The bound (A.7a) then follows by letting $t \rightarrow 0^+$ in (A.10).

544 For our second claim, we have

$$\begin{aligned} F(p, y') &= h(p) + h^*(y') - \langle y', p \rangle \\ &\leq h(p) + h^*(y) + \langle y' - y, \nabla h^*(y) \rangle + \frac{1}{2K} \|y' - y\|_\infty^2 - \langle y', p \rangle \\ &= F(p, y) + \langle y' - y, Q(y) - p \rangle + \frac{1}{2K} \|y' - y\|_\infty^2, \end{aligned} \quad (\text{A.11})$$

545 where the inequality in the second line follows from the fact that h^* is $(1/K)$ -strongly smooth [39, Theorem 12.60(e)]. ■

547 **A.2. Basic results from probability theory.** We conclude this appendix with some useful results
548 from probability theory that we will use freely throughout the sequel. For a complete treatment, we
549 refer the reader to Hall & Heyde [21].

550 **Lemma A.3** (Azuma-Hoeffding inequality). Let $M_t \in \mathbb{R}$, $t = 1, 2, \dots$, be a martingale with
551 $\|M_t - M_{t-1}\|_\infty \leq \sigma_t$ (a.s.). Then, for all $\eta > 0$, we have

$$\mathbb{P}\left(|M_t| \leq \left(2 \log(2t^2/\eta) \sum_{s=1}^t \sigma_s^2\right)^{1/2} \text{ for all } t\right) \geq 1 - \eta. \quad (\text{A.12})$$

552 **Lemma A.4** (Kolmogorov's inequality). Let $Z_t \in \mathbb{R}$, $t = 1, 2, \dots$, be a martingale difference
553 sequence that is bounded in L^2 . Then:

$$\mathbb{P}\left(\max_{s \leq t} \sum_{\ell=1}^s Z_\ell \geq \varepsilon\right) \leq \frac{1}{\varepsilon^2} \mathbb{E}\left[\left(\sum_{s=1}^t Z_s\right)^2\right] \quad \text{for all } \varepsilon > 0. \quad (\text{A.13})$$

554 **Lemma A.5** (Doob's maximal inequality). Let $Z_t \in \mathbb{R}$, $t = 1, 2, \dots$, be a martingale difference
555 sequence that is bounded in L^p for some $p \geq 1$. Then

$$\mathbb{P}\left(\max_{s \leq t} |Z_s| > \varepsilon\right) \leq \frac{1}{\varepsilon^p} \mathbb{E}[|Z_t|^p] \quad \text{for all } \varepsilon > 0. \quad (\text{A.14})$$

556 **Lemma A.6** (Burkholder–Davis–Gundy inequality). Let Z_t , $t = 1, 2, \dots$, be a martingale difference
557 sequence in \mathbb{R}^n . Then, for all $p > 1$, there exist constants c_p, C_p that depend only on p and are such
558 that

$$c_p \mathbb{E}\left[\sum_{s=1}^t \|Z_s\|_2^2\right]^{p/2} \leq \mathbb{E}\left[\max_{s \leq t} \left\|\sum_{\ell=1}^s Z_\ell\right\|_2^p\right] \leq C_p \mathbb{E}\left[\sum_{s=1}^t \|Z_s\|_2^2\right]^{p/2}. \quad (\text{A.15})$$

559 **Lemma A.7** (Robbins–Siegmund). Let \mathcal{F}_t , $t = 1, 2, \dots$, be a filtration on a complete probability
560 space $(\Omega, \mathcal{F}, \mathbb{P})$, and suppose that the sequences X_t , L_t and K_t \mathcal{F}_t -measurable, nonnegative, and
561 such that

$$\mathbb{E}[X_{t+1} | \mathcal{F}_t] \leq X_t(1 + L_t) + K_t \quad \text{with probability 1.} \quad (\text{A.16})$$

562 Then, X_t converges to some random variable X_∞ with probability 1 on the event
563 $\{\sum_{t=1}^\infty L_t < \infty \text{ and } \sum_{t=1}^\infty K_t < \infty\}$.

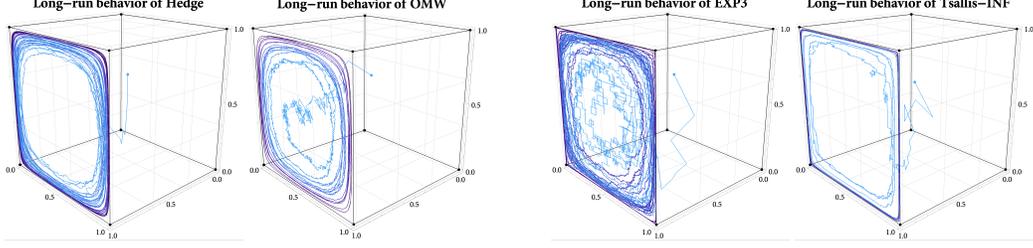


Figure 2: The long-run behavior of Algorithms 1–3 in a $2 \times 2 \times 2$ game. Algorithms 1 and 2 were run with a logit choice map as per (HEDGE); Algorithm 3 was run with both variants, EXP3 and TSALLIS-INF. All algorithms were run for 5×10^5 iterations with $\gamma_t = 1/t^{0.4}$ and $\delta_t = 0.1/t^{0.15}$; color indicates time, with darker hues indicating later iterations. The face to the left is closed under better replies, so X_t converges quickly to said face (as per Theorems 2 and 3).

564 B Specific algorithms and their properties

565 **B.1. Known algorithms as special cases of (RL).** To complement our analysis in the main part of
 566 our paper, we detail below how Algorithms 1–3 can be recast in the general framework of (RL). To
 567 lighten notation, we will assume that b_t , U_t and \hat{v}_t are respectively bounded as

$$\|b_t\|_\infty \leq B_t \quad \|U_t\|_\infty \leq \sigma_t \quad \text{and} \quad \|\hat{v}_t\|_\infty \leq M_t \quad (\text{B.1})$$

568 and we will set

$$G := \max_{i \in \mathcal{N}} \max_{\alpha \in \mathcal{A}} |v_i(\alpha)| \quad (\text{B.2})$$

569 so we can take $M_t = G + B_t + \sigma_t$ in (B.1). We will also make free use of the fact that v is Lipschitz
 570 continuous on \mathcal{X} , and we will write L for its Lipschitz modulus in the (L^1, L^∞) pair of norms on \mathcal{X}
 571 and \mathcal{Y} respectively, viz.

$$\|v(x') - v(x)\|_\infty \leq L \|x' - x\|_1 \quad \text{for all } x, x' \in \mathcal{X}. \quad (\text{B.3})$$

572 We now proceed to establish the required bounds for Algorithms 1–3:

573 **Algorithm 1.** Since $\hat{v}_t = v(X_t)$, we readily get $b_t = U_t = 0$ by definition, so Algorithm 1 fits the
 574 scheme (RL) for free with $\ell_b = \infty$, $\ell_\sigma = 0$. \blacklozenge

575 **Algorithm 2.** For the case of (Opt-FTRL), we have $\hat{v}_t = 2v(X_t) - v(X_{t-1})$ so $b_t = v(X_t) - v(X_{t-1})$,
 576 which is \mathcal{F}_t -measurable. We thus get

$$\begin{aligned} \|b_t\|_\infty &= \|\mathbb{E}[\hat{v}_t | \mathcal{F}_t] - v(X_t)\|_\infty \leq \mathbb{E}[\|v(X_t) - v(X_{t-1})\|_\infty | \mathcal{F}_t] \\ &\leq L \mathbb{E}[\|X_t - X_{t-1}\| | \mathcal{F}_t] && \# \text{ by (B.3)} \\ &= L \mathbb{E}[\|Q(Y_t) - Q(Y_{t-1})\|_\infty | \mathcal{F}_t] && \# \text{ by (Opt-FTRL)} \\ &\leq (L/K) \mathbb{E}[\|Y_t - Y_{t-1}\|_\infty | \mathcal{F}_t] && \# \text{ by Lemma A.1} \\ &\leq \gamma_t (L/K) \mathbb{E}[2v(X_t) - v(X_{t-1}) | \mathcal{F}_t] && \# \text{ by (Opt-FTRL)} \\ &\leq 3LG/K \cdot \gamma_t && \# \text{ by (B.2)} \\ &= \mathcal{O}(\gamma_t) = \mathcal{O}(1/t^{\ell_\gamma}) && (\text{B.4}) \end{aligned}$$

577 Moreover, given that \hat{v} is \mathcal{F}_t -measurable, we readily get $U_t = 0$. \blacklozenge

578 **Algorithm 3.** Since $\hat{\alpha}_t$ is sampled according to $\hat{X}_t = (1 - \delta_t)X_{i,t} + \delta_t \text{unif}_{A_i}$ (cf. Eq. (7) in
 579 Section 3), we readily obtain $\mathbb{E}[\hat{v}_{i,t} | \mathcal{F}_t] = v_i(\hat{X}_t)$, and hence, by (B.3), we get

$$B_t = \mathcal{O}(\|\hat{X}_t - X_t\|) = \mathcal{O}(\delta_t) = \mathcal{O}(1/t^{\ell_\delta}). \quad (\text{B.5})$$

580 Moreover, since $\hat{X}_{i\alpha_i,t} \geq \delta_t/A_i$, it follows that $\|\hat{v}_t\|_\infty = \mathcal{O}(1/\delta_t) = \mathcal{O}(t^{\ell_\delta})$. \blacklozenge

581 For comparison purposes, we illustrate the algorithms' behavior in a simple $2 \times 2 \times 2$ game in Fig. 2.

	Representative	Regularizer (θ)	Feedback	Bias (B_t)	Variance (σ_t)
Algorithm 1	HEDGE	$z \log z$	full info	0	0
Algorithm 2	OMW	$z \log z$	full info	$\mathcal{O}(1/t^{\ell_\gamma})$	0
Algorithm 3	EXP3	$z \log z$	payoff	$\mathcal{O}(1/t^{\ell_\delta})$	$\mathcal{O}(t^{\ell_\delta})$
Algorithm 3	TSALLIS-INF	$-4\sqrt{z}$	payoff	$\mathcal{O}(1/t^{\ell_\delta})$	$\mathcal{O}(t^{\ell_\delta})$
Algorithm 4	MP	general	full info	$\mathcal{O}(1/t^{\ell_\gamma})$	0
Algorithm 5	CMW	$z \log z$	full info	$\mathcal{O}(1/t^{\ell_\delta})$	0

Table 1: A range of algorithms adhering to the general template (RL) and their bias and variance characteristics when run with a step-size sequence of the form $\gamma_t = \gamma/t^{\ell_\gamma}$, $\ell_\gamma \in (0, 1]$, and, where applicable, a sampling parameter $\delta_t = \delta/t^{\ell_\delta}$.

582 **B.2. Further algorithms and illustrations.** To demonstrate the breadth of (RL) as an algorithmic
583 template, we provide below some more examples of algorithms from the game-theoretic literature
584 that can be recast as special cases thereof (see also Table 1 for a recap).

585 **Algorithm 4** (Mirror-prox). A progenitor of (Opt-FTRL) is the so-called *mirror-prox* (MP) algorithm
586 [24, 32], which updates as:

$$\begin{aligned} \tilde{Y}_t &= Y_t + \gamma_t v(X_t) & Y_{t+1} &= Y_t + \gamma_t v(\tilde{X}_t) \\ \tilde{X}_t &= Q(\tilde{Y}_t) & X_{t+1} &= Q(Y_{t+1}). \end{aligned} \quad (\text{MP})$$

587 The main difference between (MP) and (Opt-FTRL) is that the former utilizes two surrogate gain
588 vectors per iteration – meaning in particular that the interim, leading state \tilde{X}_t is generated with payoff
589 information from X_t , not \tilde{X}_{t-1} . This method has been used extensively in the literature for solving
590 variational inequalities and two-player, zero-sum games, cf. Juditsky et al. [24] and references therein.

591 A calculation similar to that for (Opt-FTRL) shows that Algorithm 4 has $B_t = \mathcal{O}(1/t^{\ell_\gamma})$ and $\sigma_t = 0$
592 because the algorithm has no further randomization. \blacklozenge

593 **Algorithm 5** (Clairvoyant multiplicative weights). A recent variant of the HEDGE algorithm is the
594 so-called *clairvoyant multiplicative weights* (CMW) algorithm [34]

$$Y_{i,t+1} = Y_{i,t} + \gamma_t v_i(X_{t+1}) \quad X_{i,t+1} = \Lambda_i(Y_{i,t+1}). \quad (\text{CMW})$$

595 The main difference between (CMW) and (HEDGE) is that the proxy payoff vector \hat{v}_t in (CMW) is
596 based on the *future* state X_{t+1} and *not* the current state X_t . To perform this “clairvoyant” update,
597 the players of the game must coordinate to solve an implicit fixed point problem, so (CMW) is only
598 meaningful when one has access to the payoff function $v(\cdot)$. In this regard, (CMW) can be seen as a
599 Bregman proximal point method in the general spirit of Bauschke et al. [7].

600 To cast (CMW) as an instance of the generalized template (RL), simply note that the sequence of
601 input signals is given by $\hat{v}_t = v(X_{t+1})$, so $U_t = 0$ and $b_t = v(X_{t+1}) - v(X_t) = \mathcal{O}(\gamma_t) = \mathcal{O}(1/t^{\ell_\gamma})$. \blacklozenge

602 C Proof of Theorem 1

603 Our main goal in this appendix will be to prove Theorem 1 on the resilience properties of (RL). For
604 convenience, we restate below the relevant result for ease of reference:

605 **Theorem 1.** *Let X_t , $t = 1, 2, \dots$, be the sequence of play generated by (RL) with step-size/gain*
606 *parameters $\ell_\gamma > 2\ell_\sigma$ and $\ell_b > 0$. Then, with probability 1, the limit set $\mathcal{L}(X)$ of X_t is resilient.*

607 *Proof.* Our proof that $\mathcal{L}(X)$ is resilient hinges on an energy-based technique that we will employ
608 repeatedly in other parts of our analysis. To begin, introduce a player-strategy deviation pair (i, z_i) ,
609 and say that a set is resilient to (i, z_i) if there exists an element of the set, say x^* , which counters
610 said deviation, i.e., such that $u_i(x^*) \geq u_i(z_i; x_{-i}^*)$. In this specific case, our proof proceeds by
611 contradiction, namely by assuming that, with positive probability, $\mathcal{L}(X)$ is *not* resilient to (i, z_i) . The
612 main steps of our proof unfold as follows:

613 **Step 1.** Assume that $\mathcal{L}(X)$ is not resilient to (i, z_i) with positive probability. Then there exists
614 $c, \epsilon, t_0 > 0$ such that

$$\mathbb{P}(u_i(z_i; X_{t,-i}) \geq u_i(X_t) + c \text{ for all } t \geq t_0) \geq \epsilon. \quad (\text{C.1})$$

615 *Proof of Step 1.* The function $f : x \in \mathcal{X} \mapsto u_i(z_i; x_{-i}) - u_i(x)$ is continuous and \mathcal{X} is compact, so
616 there is a definite function $\eta \equiv \eta(\delta)$ such that if $\|x - x'\| \leq \eta(\delta)$, then $|f(x) - f(x')| \leq \delta$. Now, by
617 assumption, $\{\forall x^* \in \mathcal{L}(X), u_i(z_i; x_{-i}^*) > u_i(x^*)\}$ is of positive probability. We thus get

$$\begin{aligned} 0 &< \mathbb{P}\{\forall x^* \in \mathcal{L}(X), u_i(z_i; x_{-i}^*) > u_i(x^*)\} \\ &= \mathbb{P}\left\{\inf_{x^* \in \mathcal{L}(X)} (u_i(z_i; x_{-i}^*) - u_i(x^*)) > 0\right\} \end{aligned} \quad (\text{C.2a})$$

$$= \mathbb{P}\left(\bigcup_{m>0} \left\{\inf_{x^* \in \mathcal{L}(X)} (u_i(z_i; x_{-i}^*) - u_i(x^*)) > 2^{-m}\right\}\right) \quad (\text{C.2b})$$

$$\leq \frac{1}{2} \mathbb{P}\{\forall x^* \in \mathcal{L}(X), u_i(z_i; x_{-i}^*) - u_i(x^*) > 2c\} \quad (\text{C.2c})$$

618 for some $c > 0$ in (C.2c), and where (C.2a) is because $\mathcal{L}(X)$ is closed – hence compact – almost
619 surely. Therefore, by definition of $\eta(\cdot)$,

$$0 < \mathbb{P}\{\forall x^* \in \mathcal{X}, \text{dist}(x^*, \mathcal{L}(X)) \leq \eta(c) \Rightarrow u_i(z_i; x_{-i}^*) - u_i(x^*) > c\} = 2\epsilon \quad (\text{C.2d})$$

620 Now, let t_0 such that $\mathbb{P}\{\forall t \geq t_0, \text{dist}(X_t, \mathcal{L}(X)) \leq \eta(c)\} > 1 - \frac{\epsilon}{2}$. Then by construction, we get

$$\mathbb{P}\{\forall t \geq t_0, u_i(z_i; X_{t,-i}) > u_i(X_t) + c\} > \epsilon. \quad (\text{C.3})$$

621 and our proof is complete. \blacksquare

622 Intuitively, the existence of an action that consistently outperforms X_t runs contrary to the behavior
623 that one would expect from any regularized learning algorithm. We will proceed to make this intuition
624 precise below by means of an energy argument. To that end, consider the Fenchel coupling

$$F_t = h_i(z_i) + h_i^*(Y_{i,t}) - \langle Y_{i,t}, z_i \rangle \quad (\text{C.4})$$

625 Then, by Lemma A.2 in Appendix A, we readily get that

$$F_{t+1} \leq F_t - \gamma_t \langle \hat{v}_{i,t}, z_i - X_{i,t} \rangle + \frac{\gamma_t^2}{2\kappa_h} \|\hat{v}_{i,t}\|_\infty^2. \quad (\text{C.5})$$

626 where, in obvious notation, we are identifying $z_i \in \mathcal{A}_i$ with the corresponding vertex e_{z_i} of $\mathcal{X}_i =$
627 $\Delta(\mathcal{A}_i)$. To proceed, the main idea will be to relate $\gamma_t \langle \hat{v}_{i,t}, z_i - X_{i,t} \rangle$ to its “perfect” counterpart
628 $\gamma_t \langle v_i(X_t), z_i - X_{i,t} \rangle$. We formalize this below.

629 **Step 2.** If $\mathcal{L}(X)$ is not resilient to (i, z_i) , there exists $t_1 \geq t_0$ such that, with probability $\epsilon'/2 > 0$, and
630 for all $t \geq t_1$, we have

$$F_t \leq F_{t_0} - \frac{c}{2} \sum_{s=t_0}^t \gamma_s. \quad (\text{C.6})$$

631 *Proof of Step 2.* With probability ϵ' and for all $t \geq t_0$, we have

$$\gamma_t \langle \hat{v}_{i,t}, z_i - X_{i,t} \rangle = \gamma_t \langle v_i(X_t), z_i - X_{i,t} \rangle + \gamma_t \langle U_{i,t}, z_i - X_{i,t} \rangle + \gamma_t \langle b_{i,t}, z_i - X_{i,t} \rangle \quad (\text{C.7})$$

$$\geq [c + \langle U_{i,t}, z_i - X_{i,t} \rangle + \langle b_{i,t}, z_i - X_{i,t} \rangle] \gamma_t. \quad (\text{C.8})$$

632 The combination of Eqs. (C.5) and (C.8) then provides the following upper bound of F_{t+1} :

$$F_{t+1} \leq F_t - c\gamma_t + \gamma_t \langle U_{i,t}, z_i - X_{i,t} \rangle + \gamma_t \langle b_{i,t}, z_i - X_{i,t} \rangle + \frac{\gamma_t^2}{2\kappa_h} \|\hat{v}_{i,t}\|_\infty^2 \quad (\text{C.9})$$

$$\leq F_{t_0} - c \sum_{s=t_0}^t \gamma_s + \underbrace{\sum_{s=t_0}^t \gamma_s \langle U_{s,i}, z_i - X_{s,i} \rangle}_{E_{U,t}} + \underbrace{\sum_{s=t_0}^t \gamma_s \langle b_{s,i}, z_i - X_{s,i} \rangle}_{E_{b,t}} + \sum_{s=t_0}^t \frac{\|\hat{v}_{s,i}\|_\infty^2}{2\kappa_h} \gamma_s^2. \quad (\text{C.10})$$

633 We are thus left to show is that $c \sum_{s=t_0}^t \gamma_s$ is the dominant term above. To do so, we proceed to
634 examine each term individually:

635 • *Second-order term:* We first deal with the second-order term $\sum_{s=t_0}^t \frac{\|\hat{v}_{s,i}\|_\infty^2}{2\kappa_h} \gamma_s^2$. By expanding the
 636 $\|\hat{v}_{s,i}\|_\infty^2$, we readily get

$$\frac{\sum_{s=t_0}^t \|\hat{v}_{s,i}\|_\infty^2 \gamma_s^2}{\tau_t} = \mathcal{O}\left(\frac{\sum_{s=1}^t \gamma_s^2 (1 + B_s^2 + \sigma_s^2)}{\sum_{s=1}^t \gamma_s}\right). \quad (\text{C.11})$$

637 However, by our assumptions on the parameters of (RL), we readily get

$$\lim_{t \rightarrow \infty} \frac{\gamma_t^2 (1 + B_t^2 + \sigma_t^2)}{\gamma_t} = 0 \quad (\text{C.12})$$

638 so we conclude that

$$\lim_{t \rightarrow \infty} \frac{\sum_{s=1}^t \gamma_s^2 (1 + B_s^2 + \sigma_s^2)}{\sum_{s=1}^t \gamma_s} \quad (\text{C.13})$$

639 by the Stolz-Cesàro theorem.

640 • *Bias term:* By far the most immediate, the bias term $E_{b,t}$ is bounded as

$$E_{b,t} \leq 2 \sum_{s=t_0}^t \|b_{s,i}\|_\infty \gamma_s \leq 2 \sum_{s=t_0}^t B_s \gamma_s = o\left(\sum_{s=t_0}^t \gamma_s\right) \quad \text{as } t \rightarrow \infty. \quad (\text{C.14})$$

641 • *Noise term:* Finally, the noise term $E_{U,t}$ is bounded by means of the Azuma-Hoeffding inequality,
 642 cf. Lemma A.3 in Appendix A. Specifically, with probability at least $1 - \varepsilon'/2$, we have

$$\begin{aligned} E_{U,t} &:= \sum_{s=t_0}^t \gamma_s \langle U_{s,i}, z_i - X_{s,i} \rangle \\ &\leq 2 \left(\sum_{s=t_0}^t \|U_{s,i}\|_\infty^2 \gamma_s^2 \right)^{1/2} \sqrt{2 \log\left(\frac{4t^2}{\varepsilon'}\right)} \\ &\leq 2 \left(\sum_{s=t_0}^t \sigma_s^2 \gamma_s^2 \right)^{1/2} \sqrt{2 \log\left(\frac{4t^2}{\varepsilon'}\right)}. \end{aligned} \quad (\text{C.15})$$

643 for all $t \geq t_0$. To proceed, note that a second application of the Stolz-Cesàro theorem yields
 644 $\sum_{s=t_0}^t \sigma_s^2 \gamma_s^2 = o(\sum_{s=t_0}^t \gamma_s)$ and, moreover, note that $\log(4t^2/\varepsilon') = \mathcal{O}(\sum_{s=t_0}^t \gamma_s)$. Taking square
 645 roots and multiplying then yields that

$$E_{U,t} = o\left(\sum_{s=t_0}^t \gamma_s\right) \quad (\text{C.16})$$

646 with probability at least $1 - \varepsilon'/2$.

647 We are now in a position to establish the bound Eq. (C.6). Indeed, putting Eqs. (C.13), (C.14)
 648 and (C.16) together, we readily infer that there exists $t_1 \geq t_0$ such that, with probability at least
 649 $1 - \varepsilon'/2$, we have

$$\sum_{s=t_0}^t \gamma_s \langle U_{s,i}, z_i - X_{s,i} \rangle + \sum_{s=t_0}^t \gamma_s \langle b_{s,i}, z_i - X_{s,i} \rangle + \sum_{s=t_0}^t \frac{\|\hat{v}_{s,i}\|_\infty^2}{2\kappa_h} \gamma_s^2 \leq \frac{c}{2} \sum_{s=t_0}^t \gamma_s \quad (\text{C.17})$$

650 for all $t \geq t_1$. This proves Eq. (C.6) and concludes our proof. ■

651 Summarizing the above, we have shown that, with probability at least $1 - \varepsilon'/2$, we have

$$F_{t+1} \leq F_{t_0} - \frac{c}{2} \sum_{s=t_0}^t \gamma_s \rightarrow -\infty \quad \text{as } t \rightarrow \infty. \quad (\text{C.18})$$

652 Since F is nonnegative (by Lemma A.2), we have established that the event where $\mathcal{L}(X)$ is not
 653 resilient to (i, z_i) is an event of probability zero. However, since there are uncountably many strategic
 654 deviations, the proof is not yet complete; the last step involves an approximation by deviations with
 655 rational entries.

656 **Step 3.** $\mathcal{L}(X)$ is almost-surely resilient.

657 *Proof of Step 3.* The key point of the proof is the observation that a closed set is resilient if and
 658 only if it is *rationally* resilient, i.e., it nullifies all *rational* deviations $z_i \in \mathcal{X}_i \cap \mathbb{Q}^{\mathcal{A}_i}$ (which are
 659 countably many). Indeed, if $\mathcal{L}(X)$ is not resilient with positive probability, then, likewise, $\mathcal{L}(X)$ will
 660 not be rationally resilient with positive probability either. Because there are countably many rational
 661 deviations, there must be a rational strategic deviation (i, z_i) (with $z_i \in \mathcal{X}_i \cap \mathbb{Q}^{\mathcal{A}_i}$) to which $\mathcal{L}(X)$ is
 662 not resilient. This comes in contradiction with the conclusions of [Step 2](#). ■

663 This concludes the last required step, so the proof of [Theorem 1](#) is now complete. ■

664 **D Proof of Theorems 2 and 3**

665 In this last appendix, our goal is to prove our characterization of club sets, namely:

666 **Theorem 2.** Fix some set $\mathcal{S} \in \mathcal{P}(\mathcal{X})$ and suppose that (RL) is run with a steep regularizer and
 667 step-size/gain parameters $\ell_\gamma \in [0, 1]$, $\ell_b > 0$, and $\ell_\sigma < 1/2$. Then:

- 668 1. \mathcal{S} is stochastically asymptotically stable under (RL) if and only if it is a club set.
 669 2. \mathcal{S} is irreducibly stable under (RL) if and only if it is an m -club set.

670 **Theorem 3.** Let $\mathcal{S} \in \mathcal{P}(\mathcal{X})$ be a club set, and let X_t , $t = 1, 2, \dots$, be the sequence of play generated
 671 by (RL) with parameters $\ell_\gamma \in [0, 1]$, $\ell_b > 0$, and $\ell_\sigma < 1/2$. Then, for all $\epsilon > 0$, there exists an (open,
 672 unbounded) initialization domain $\mathcal{D} \subseteq \mathcal{Y}$ such that, with probability at least $1 - \epsilon$, we have

$$\text{dist}(X_t, \mathcal{S}) \leq C\varphi\left(c_1 - c_2 \sum_{s=1}^t \gamma_s\right) \quad \text{whenever } Y_1 \in \mathcal{D} \quad (11)$$

673 where C, c_1, c_2 are constants ($C, c_2 > 0$), and the rate function φ is given by $\varphi(z) = (\theta')^{-1}(z)$ if
 674 $z > \lim_{z \rightarrow 0^+} \theta'(z)$, and $\varphi(z) = 0$ otherwise.

675 Our proof strategy will be to construct a sheaf of “linearized” energy functions which, when bundled
 676 together, yield a suitable Lyapunov-like function for \mathcal{S} . To do so, let $\mathcal{C} = \prod_i \mathcal{C}_i$ denote the support of
 677 \mathcal{S} (cf. the definition of club sets), and let

$$\mathcal{Z}_i = \{e_{i\alpha'_i} - e_{i\alpha_i} : \alpha_i \in \mathcal{C}_i, \alpha'_i \in \mathcal{A}_i \setminus \mathcal{C}_i\} \quad (\text{D.1})$$

678 and

$$\mathcal{Z} = \bigcup_{i \in \mathcal{N}} \mathcal{Z}_i \quad (\text{D.2})$$

679 denote the set of all pure strategic deviations from \mathcal{S} . Then, our ensemble of candidate energy
 680 functions will be given by

$$E_z(y) = \langle y, z \rangle \quad \text{for } z \in \mathcal{Z}, y \in \mathcal{V}^*. \quad (\text{D.3})$$

681 The motivation for this definition is given by the following lemma.

682 **Lemma D.1.** Suppose that the sequence $y_t \in \mathcal{V}^*$, $t = 1, 2, \dots$, has $E_z(y_t) \rightarrow -\infty$ for all $z \in \mathcal{Z}$ as
 683 $t \rightarrow \infty$. Then the sequence $x_t = Q(y_t)$ converges to \mathcal{S} as $t \rightarrow \infty$.

684 *Proof.* Let $z = e_{i\alpha'_i} - e_{i\alpha_i}$ for some $i \in \mathcal{N}$, $\alpha_i \in \mathcal{C}_i$, and $\alpha'_i \in \mathcal{A}_i \setminus \mathcal{C}_i$. Since $E_z(y_t) \rightarrow -\infty$ by
 685 assumption, we get $y_{i\alpha'_i, t} - y_{i\alpha_i, t} \rightarrow -\infty$ and hence, by [Lemma A.1](#), we conclude that $Q_{i\alpha'_i}(x_t) \rightarrow 0$
 686 as $t \rightarrow \infty$. In turn, given that this holds for all $i \in \mathcal{N}$ and all $\alpha'_i \in \mathcal{A}_i \setminus \mathcal{C}_i$, we conclude that $x_t = Q(y_t)$
 687 converges to \mathcal{S} . ■

688 In view of the above, we will focus on showing that $E_z(Y_t) \rightarrow -\infty$ for all $z \in \mathcal{Z}$. As a first step, we
 689 establish a basic template inequality for the evolution of E_z under (RL).

690 **Lemma D.2.** Fix some $z \in \mathcal{Z}$ and let $E_t := E_z(Y_t)$. Then, for all $t = 1, 2, \dots$, we have

$$E_{t+1} \leq E_t + \gamma_t \langle v(X_t), z \rangle + \gamma_t \xi_t + \gamma_t \psi_t \quad (\text{D.4})$$

691 where the error terms ξ_t and ψ_t are given by

$$\xi_t = \langle U_t, z \rangle \quad \text{and} \quad \psi_t = 2B_t. \quad (\text{D.5})$$

692 *Proof.* Simply set $y \leftarrow Y_{t+1}$ in $E_z(y)$, invoke the definition of the update $Y_t \leftarrow Y_{t+1}$ in (RL), and note
 693 that $|\langle b_t, z \rangle| \leq \|z\| \|b_t\|_\infty \leq 2B_t$ by the definition of \mathcal{Z} . ■

694 The key take-away from (D.4) is that, if X_t is close to \mathcal{S} and $\alpha_i \in \mathcal{C}_i$, $\alpha'_i \in \mathcal{A}_i \setminus \mathcal{C}_i$, we will have

$$\langle v(X_t), z \rangle = v_{i\alpha'_i}(X_t) - v_{i\alpha_i}(X_t) = u_i(\alpha'_i; X_{-i,t}) - u_i(\alpha_i; X_{-i,t}) < 0 \quad (\text{D.6})$$

695 by the continuity of u_i and the assumption that \mathcal{S} is a club set. More concretely, by the definition of
 696 the better-reply correspondence, we have

$$\langle v(x^*), z \rangle < 0 \quad \text{for all } x^* \in \mathcal{S} \text{ and all } z \in \mathcal{Z} \quad (\text{D.7})$$

697 and hence, by continuity, there exists a neighborhood \mathcal{B} of \mathcal{S} such that

$$\langle v(x), z \rangle < 0 \quad \text{for all } x \in \mathcal{B} \text{ and all } z \in \mathcal{Z}. \quad (\text{D.8})$$

698 In other words, as long as X_t is sufficiently close to \mathcal{S} , (D.4) exhibits a consistent negative drift
 699 pushing E_t towards $-\infty$.

700 To exploit this “dynamic consistency” property of \mathcal{S} , it will be convenient to introduce the family of
 701 sets

$$\mathcal{D}(\epsilon) = \{y \in \mathcal{V}^* : \langle y, z \rangle < -\epsilon \text{ for all } z \in \mathcal{Z}\} \quad (\text{D.9})$$

702 As we show below, these sets are mapped under Q to neighborhoods of \mathcal{S} , so they are particularly
 703 well-suited to serve as initialization domains for (RL). This is encoded in the following properties:

704 **Lemma D.3.** *Let $x = Q(y)$ for some $y \in \mathcal{V}^*$. Then, for all $\alpha_i, \alpha'_i, i \in \mathcal{N}$, we have*

$$x_{i\alpha_i} \leq \varphi \left(\theta(1^-) + y_{i\alpha'_i} - y_{i\alpha_i} \right) \quad (\text{D.10})$$

705 with φ defined as per Theorem 3, i.e.,

$$\varphi(z) = \begin{cases} 0 & \text{if } z \leq \theta'(0^+), \\ (\theta')^{-1}(z) & \text{if } \theta'(0^+) < z < \theta'(1^-), \\ 1 & \text{if } z \geq \theta'(1^-). \end{cases} \quad (\text{D.11})$$

706 **Corollary D.1.** *For all $\delta > 0$ there exists some $\epsilon_\delta \in \mathbb{R}$ such that, for all $\epsilon > \epsilon_\delta$ and all $y \in \mathcal{D}_\epsilon$, we
 707 have*

$$Q_{i\alpha'_i}(y_i) < \delta \quad \text{for all } \alpha'_i \in \mathcal{A}_i \setminus \mathcal{C}_i \text{ and all } i \in \mathcal{N}. \quad (\text{D.12})$$

708 *Proof of Lemma D.3.* Suppressing the player index for simplicity, the first-order stationarity condi-
 709 tions for the convex problem (4) readily give

$$y_\alpha - \theta'(x_\alpha) = \mu - \nu_\alpha, \quad (\text{D.13})$$

710 where μ is the Lagrange multiplier for the equality constraint $\sum_\alpha x_\alpha = 1$, and ν_α is the complementary
 711 slackness multiplier of the inequality constraint $x_\alpha \geq 0$ (so $\nu_\alpha = 0$ whenever $x_\alpha > 0$). Thus, rewriting
 712 (D.13) for some $\alpha \in \mathcal{A}$, we get

$$y_{\alpha'} - y_\alpha = \theta'(x_{\alpha'}) - \theta'(x_\alpha) + \nu_\alpha - \nu_{\alpha'} \quad (\text{D.14})$$

713 and hence

$$\theta'(x_{\alpha'}) = \theta'(x_\alpha) + \nu_{\alpha'} - \nu_\alpha + y_{\alpha'} - y_\alpha \leq \theta'(1^-) + \nu_{\alpha'} + y_{\alpha'} - y_\alpha, \quad (\text{D.15})$$

714 where we used the fact that $\nu_\alpha \geq 0$. Now, if $\theta'(1^-) + y_{\alpha'} - y_\alpha < \theta'(0^+)$ and $x_{\alpha'} > 0$ (so $\nu_{\alpha'} = 0$), we
 715 will have $\theta'(x_{\alpha'}) < \theta'(0^+)$, a contradiction. This shows that $x_{\alpha'} = 0$ if $\theta'(1^-) + y_{\alpha'} - y_\alpha < \theta'(0^+)$,
 716 so (D.10) is satisfied in this case. Otherwise, if $x_{\alpha'} > 0$, we must have $\nu_{\alpha'} = 0$ by complementary
 717 slackness, so (D.10) follows by applying the second branch of (D.11) to (D.15). ■

718 The above provides us with a fairly good handle on the local geometric and dynamic properties of
 719 \mathcal{S} . On the flip side however, the various error terms in (D.5) may be positive, so E_t may fail to be
 720 decreasing and X_t may drift away from \mathcal{S} . On that account, it will be convenient to introduce the
 721 aggregate error processes

$$\mathbf{I}_t = \sum_{s=1}^t \gamma_s \xi_s \quad \text{and} \quad \mathbf{\Pi}_t = \sum_{s=1}^t \gamma_s \psi_s. \quad (\text{D.16})$$

722 Intuitively, the aggregates (D.16) measure the total effect of each error term in (D.4), so we will
 723 establish a first series of results under the following general requirements:

724 1. *Subleading error growth:*

$$\lim_{t \rightarrow \infty} I_t / \tau_t = 0 \quad (\text{Sub.I})$$

$$\lim_{t \rightarrow \infty} \Pi_t / \tau_t = 0 \quad (\text{Sub.II})$$

725 where $\tau_t = \sum_{s=1}^t \gamma_s$ and both limits are to be interpreted in the almost sure sense.

726 2. *Drift dominance:*

$$\mathbb{P}(I_t \leq C\tau_t^\alpha / 2 \text{ for all } t) \geq 1 - \eta \quad (\text{Dom.I})$$

$$\mathbb{P}(\Pi_t \leq C\tau_t^\alpha / 2 \text{ for all } t) \geq 1 - \eta \quad (\text{Dom.II})$$

727 for some $C > 0$ and $\alpha \in [0, 1)$.

728 In a nutshell, (Sub) posits that the aggregate error processes I_t and Π_t of (D.16) are subleading
729 relative to the long-run drift of (D.4), while (Dom) goes a step further and asks that said errors
730 are asymptotically dominated by the drift in (D.4). Accordingly, under these implicit error control
731 conditions, we obtain the interim convergence result below:

732 **Proposition D.1.** *Let \mathcal{S} be a club set, fix some confidence threshold $\eta > 0$, and let $X_t = Q(Y_t)$
733 be the sequence of play generated by (RL). If (Sub) and (Dom) hold, there exists an unbounded
734 initialization domain $\mathcal{D} \subseteq \mathcal{V}^*$ such that*

$$\mathbb{P}(X_t \text{ converges to } \mathcal{S} \mid Y_1 \in \mathcal{D}) \geq 1 - 2\eta. \quad (\text{D.19})$$

735 *Proof of Proposition D.1.* Fix some $z \in \mathcal{Z}$, let $E_t = E_z(Y_t)$, and pick $\alpha \in [0, 1)$ so that (Dom) holds
736 for some $C > 0$. In addition, set $c = -\sup_{x \in \mathcal{B}} \langle v(x), z \rangle > 0$, let $t_0 = \inf\{t : c\tau_t > C\tau_t^\alpha\}$, and write
737 $\Delta E = \max_t \{C\tau_t^\alpha - c\tau_t\}$. Then, if Y_1 is initialized in $\mathcal{D} \leftarrow \mathcal{D}(\epsilon + \Delta E)$ where ϵ is such that $\mathcal{D}(\epsilon) \subseteq \mathcal{B}$,
738 we will have $Y_t \in \mathcal{D}(\epsilon)$ for all t . Indeed, this being trivially the case for $t = 1$, assume it to be the
739 case for all $s = 1, 2, \dots, t$. Then, by (D.4) and our inductive hypothesis, we get

$$E_{t+1} \leq E_1 - \sum_{s=1}^t \gamma_s \langle v(X_s), z \rangle + I_t + \Pi_t \leq -\epsilon - \Delta E - c\tau_t + C\tau_t^\alpha / 2 + C\tau_t^\alpha / 2 \leq -\epsilon - \Delta E + \Delta E = -\epsilon \quad (\text{D.20})$$

740 i.e., $E_{t+1} \in \mathcal{D}(\epsilon)$, as claimed.

741 Now, since $E_t \in \mathcal{D}(\epsilon)$ for all t , we conclude that

$$E_{t+1} \leq E_1 - c\tau_t + I_t + \Pi_t \quad \text{for all } t = 1, 2, \dots \quad (\text{D.21})$$

742 Thus, if (Sub) holds, we readily get $E_t \rightarrow -\infty$ with probability 1 on the event that (Dom.I) and
743 (Dom.II) both hold. This implies that $E_t \rightarrow -\infty$, and since $z \in \mathcal{Z}$ above is arbitrary, we conclude
744 that $X_t \rightarrow \mathcal{S}$ with probability at least $1 - 2\eta$, as claimed. ■

745 We are now in a position to prove Theorem 2.

746 *Proof of Theorem 2.* Our proof will hinge on showing that (Sub) and (Dom) hold under the stated
747 step-size and sampling parameter schedules. Our claim will then follow by a direct application of
748 Proposition D.1 and a reduction to a suitable subface of \mathcal{X} .

749 First, regarding (Sub), the law of large numbers for martingale difference sequences [21, Theorem
750 2.18] shows that $I_t / \tau_t \rightarrow 0$ with probability 1 on the event $\{\sum_t \gamma_t^2 \mathbb{E}[\xi_t^2 \mid \mathcal{F}_t] / \tau_t^2 < \infty\}$. However

$$\mathbb{E}[\xi_t^2 \mid \mathcal{F}_t] \leq 2^2 \mathbb{E}[\|U_t\|_\infty^2 \mid \mathcal{F}_t] \leq 2^2 \sigma_t^2 = \mathcal{O}(t^{2\ell_\sigma}) \quad (\text{D.22})$$

751 so, in turn, we get

$$\sum_t \frac{\gamma_t^2 \mathbb{E}[\xi_t^2 \mid \mathcal{F}_t]}{\tau_t^2} = \mathcal{O}\left(\sum_t \frac{\gamma_t^2 \sigma_t^2}{\tau_t^2}\right) = \mathcal{O}\left(\sum_t \frac{t^{-2\ell_\gamma} t^{2\ell_\sigma}}{t^{2(1-\ell_\gamma)}}\right) = \mathcal{O}\left(\sum_t \frac{1}{t^{2-2\ell_\sigma}}\right) < \infty \quad (\text{D.23})$$

752 given that $\ell_\sigma < 1/2$. This establishes (Sub.I); the remaining requirement (Sub.II) follows trivially by
753 noting that $\sum_{s=1}^t \gamma_s B_s / \sum_{s=1}^t \gamma_s \rightarrow 0$ if and only if $B_t \rightarrow 0$, which is immediate from the theorem's
754 assumptions.

755 Second, regarding (Dom), since B_t is deterministic and $B_t = \mathcal{O}(1/t^{\ell_b})$ for some $\ell_b > 0$, it is always
 756 possible to find $C > 0$ and $\alpha \in (0, 1)$ so that (Dom.II) holds. We are thus left to establish (Dom.I).
 757 To that end, let $I_t^* = \sup_{1 \leq s \leq t} |I_t|$ and set $P_t := \mathbb{P}(I_t^* > C\tau_t^\alpha/2)$ so

$$P_t \leq \frac{\mathbb{E}[|I_t|^q]}{(C/2)^q \tau_t^{\alpha q}} \leq c_q \frac{\mathbb{E}\left[\left(\sum_{s=1}^t \gamma_s^2 \|U_s\|_\infty^2\right)^{q/2}\right]}{\tau_t^{\alpha q}} \quad (\text{D.24})$$

758 where c_q is a positive constant depending only on C and q , and we used Kolmogorov's inequality
 759 (Lemma A.4) in the first step and the Burkholder–Davis–Gundy inequality (Lemma A.6) in the
 760 second.

761 To proceed, we will require the following variant of Hölder's inequality [8, p. 15]:

$$\left(\sum_{s=1}^t a_s b_s\right)^\rho \leq \left(\sum_{s=1}^t a_s^{\frac{\lambda \rho}{\rho-1}}\right)^{\rho-1} \sum_{s=1}^t a_s^{(1-\lambda)\rho} b_s^\rho \quad (\text{D.25})$$

762 valid for all $a_s, b_s \geq 0$ and all $\rho > 1$, $\lambda \in [0, 1)$. Then, substituting $a_s \leftarrow \gamma_s^2$, $b_s \leftarrow \|U_s\|_\infty^2$,
 763 $\rho \leftarrow q/2$ and $\lambda \leftarrow 1/2 - 1/q$, (D.24) gives

$$P_t \leq c_q \frac{\left(\sum_{s=1}^t \gamma_s\right)^{q/2-1} \sum_{s=1}^t \gamma_s^{1+q/2} \mathbb{E}[\|U_s\|_\infty^q]}{\tau_t^{\alpha q}} \leq c_q \frac{\sum_{s=1}^t \gamma_s^{1+q/2} \sigma_s^q}{\tau_t^{1+(\alpha-1/2)q}} \quad (\text{D.26})$$

764 We now consider two cases, depending on whether the numerator of (D.26) is summable or not.

765 *Case 1:* $\ell_\gamma(1+q/2) \geq 1+q\ell_\sigma$. In this case, the numerator of (D.26) is summable under the theorem's
 766 assumptions, so the fraction in (D.26) behaves as $\mathcal{O}(1/t^{(1-\ell_\gamma)(1+(\alpha-1/2)q)})$.

767 *Case 2:* $\ell_\gamma(1+q/2) < 1+q\ell_\sigma$. In this case, the numerator of (D.26) is not
 768 summable under the theorem's assumptions, so the fraction in (D.26) behaves as
 769 $\mathcal{O}(t^{1-\ell_\gamma(1+q/2)+q\ell_\sigma} / t^{(1-\ell_\gamma)(1+(\alpha-1/2)q)})$.

770 Thus, working out the various exponents, a tedious – but otherwise straightforward – calculation
 771 shows that there exists some $\alpha \in (0, 1)$ such that P_t is summable as long as $\ell_\sigma < 1/2 - 1/q$ and
 772 $0 \leq \ell_\gamma < q/(2+q)$. Hence, if γ is sufficiently small relative to η , we conclude that

$$\mathbb{P}(I_t \leq C\tau_t^\alpha/2 \text{ for all } t) \geq 1 - \sum_t P_t \geq 1 - \eta/2. \quad (\text{D.27})$$

773 Finally, if $\ell_\gamma > 1/2 + \ell_\sigma$, (Dom.I) is a straightforward consequence of (D.24) for $q = 2$.

774 With all this in hand, the final steps of our proof proceed as follows:

775 **Closedness \implies Stability.** Our assertion follows by invoking Proposition D.1. ■

776 **Stability \implies Closedness.** Suppose that \mathcal{S} is not club. Then there exists some pure strategy $\alpha \in \mathcal{C}$
 777 and some deviation $\alpha' \notin \mathcal{C}$ such that the deviation from α to α' is not costly to the deviating player.
 778 Thus, if we consider the restriction of the game to the face spanned by α and α' (a single-player game
 779 with two strategies), the corresponding score difference will be

$$y_{\alpha',t} - y_{\alpha,t} \geq \sum_{s=1}^t \gamma_s b_s + \sum_{s=1}^t \gamma_s U_s \quad (\text{D.28})$$

780 By our standing assumptions for b_t and U_t (and Doob's martingale convergence theorem for the
 781 latter), both $\sum_{s=1}^t \gamma_s b_s$ and $\sum_{s=1}^t \gamma_s U_s$ will be bounded from below by some (a.s.) finite random
 782 variable A_0 . Since θ is steep, it follows that, with probability 1, $\liminf_{t \rightarrow \infty} (y_{\alpha,t}) > 0$, so \mathcal{C} cannot be
 783 stable. ■

784 **Minimality \implies Irreducible Stability.** Suppose that \mathcal{S} is m-club. Then, by our previous claim, \mathcal{S}
 785 is stochastically asymptotically stable. If \mathcal{S} contains a proper subface $\mathcal{S}' \subsetneq \mathcal{S}$ that is also stochastically
 786 asymptotically stable, \mathcal{S}' must be club by the converse implication of the first part of the theorem.
 787 However, in that case, \mathcal{S} would not be m-club, a contradiction which proves our claim. ■

788 **Irreducible Stability \implies Minimality.** For our last claim, assume that \mathcal{S} is irreducibly stable. By
789 the first part of our theorem, this implies that \mathcal{S} is club. Then, if it so happens that \mathcal{S} is not m-club, it
790 would contain a proper club subspace $\mathcal{S}' \subsetneq \mathcal{S}$; by the first part of our theorem, this set would be itself
791 stochastically asymptotically stable, in contradiction to the irreducibility assumption. This shows that
792 \mathcal{S} is m-club and concludes our proof. ■

793 We are only left to establish the convergence rate estimate of [Theorem 3](#).

794 *Proof of Theorem 3.* Going back to (D.21) and invoking [Lemma D.3](#) shows that there exist constants
795 $c_1 > 0$ and $c_2 \in \mathbb{R}$ such that, for all $\alpha_i \in \mathcal{A}_i \setminus \mathcal{C}_i, i \in \mathcal{N}$, we have

$$X_{i\alpha_i,t} \leq \varphi(\theta(1^-) + E_t) \leq \varphi(c_2 - c_1\tau_t) \quad (\text{D.29})$$

796 with probability 1 on the events of (Dom). We thus get

$$\text{dist}_1(X_t, \mathcal{S}) \leq \sum_{i \in \mathcal{N}} \sum_{\alpha_i \in \mathcal{A}_i \setminus \mathcal{C}_i} \varphi(c_2 - c_1\tau_t), \quad (\text{D.30})$$

797 and our proof is complete. ■

798 As for the rate estimates of [Corollary 2](#), the proof boils down to a simple derivation of the correspond-
799 ing rate functions:

800 *Proof of Corollary 2.* By a straightforward calculation, we have:

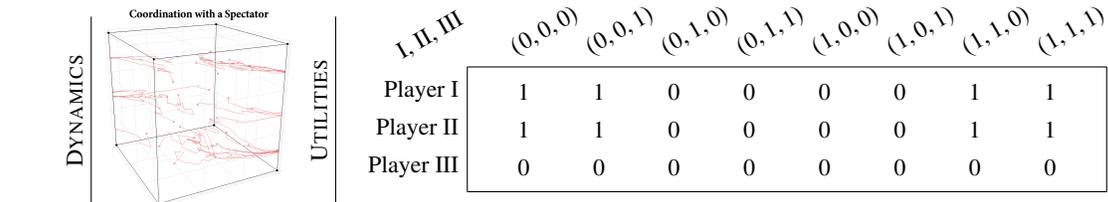
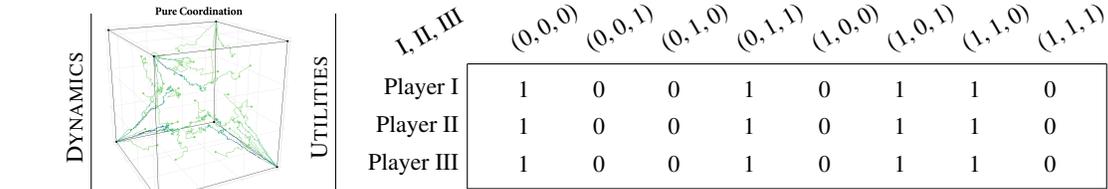
- 801 1. If $\theta(z) = z \log z$ then $\varphi(z) = \exp(1 + z)$.
- 802 2. If $\theta(z) = -4\sqrt{z}$ then $\varphi(z) = 4/z^2$.
- 803 3. If $\theta(z) = z^2/2$ then $\varphi(z) = [z]_0^1$.

804 Our claims then follow immediatly from the rate estimate (11) of [Theorem 2](#). ■

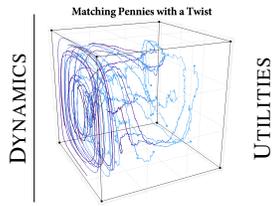
805 E Details on the numerics

806 In all our experiments, we ran the EXP3 variant of bandit FTRL (B-FTRL) (cf. [Algorithm 3](#)) with
807 step-size and sampling radius parameters $\gamma_t = 0.2 \times t^{-1/2}$ and $\delta_t = 0.1 \times t^{-0.15}$ respectively. The
808 algorithm was run for $T = 10^4$ iterations and, to reduce graphical clutter, we plotted only every third
809 point of each trajectory. Trajectories have been colored throughout with darker hues indicating later
810 times (e.g., light blue indicates that the trajectory is closer in time to its starting point, darker shades
811 of blue indicate proximity to the termination time). The algorithm's initial conditions were taken
812 from a uniform initialization grid of the form $y_1 \in \{-1, 0, 1\}^3$ and perturbed by a uniform random
813 number in $[-0.1, -0.1]$ to avoid non-generic initializations.

814 The payoffs of the chosen games were normalized to $[-1, 1]$ and players are assumed to choose
815 between two actions labeled “O” and “1”. The specific tableaus are shown in the table below, next to
816 the respective portrait (all taken from [Fig. 1](#)).

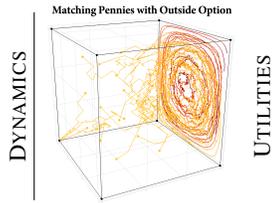


819



I, II, III	(0, 0, 0)	(0, 0, 1)	(0, 1, 0)	(0, 1, 1)	(1, 0, 0)	(1, 0, 1)	(1, 1, 0)	(1, 1, 1)
Player I	0	0	0	0	0.1	0.1	0.1	0.1
Player II	1	0	0	1	0	1	1	0
Player III	0	1	1	0	1	0	0	1

820



I, II, III	(0, 0, 0)	(0, 0, 1)	(0, 1, 0)	(0, 1, 1)	(1, 0, 0)	(1, 0, 1)	(1, 1, 0)	(1, 1, 1)
Player I	-1	1	-1	1	1	-1	1	-1
Player II	1	1	-1	-1	-1	1	1	-1
Player III	1	-1	-1	1	-1	1	1	-1