

On Gradient-Based Learning in Continuous Games*

Eric Mazumdar[†], Lillian J. Ratliff[‡], and S. Shankar Sastry[§]

Abstract.

We introduce a general framework for competitive gradient-based learning that encompasses a wide breadth of multi-agent learning algorithms, and analyze the limiting behavior of competitive gradient-based learning algorithms using dynamical systems theory. For both general-sum and potential games, we characterize a non-negligible subset of the local Nash equilibria that will be avoided if each agent employs a gradient-based learning algorithm. We also shed light on the issue of convergence to non-Nash strategies in general- and zero-sum games, which may have no relevance to the underlying game, and arise solely due to the choice of algorithm. The existence and frequency of such strategies may explain some of the difficulties encountered when using gradient descent in zero-sum games as, e.g., in the training of generative adversarial networks. To reinforce the theoretical contributions, we provide empirical results that highlight the frequency of linear quadratic dynamic games (a benchmark for multi-agent reinforcement learning) that admit global Nash equilibria that are almost surely avoided by policy gradient.

Key words. continuous games, gradient-based algorithms, multi-agent learning

AMS subject classifications.

1. Introduction. With machine learning algorithms increasingly being deployed in real world settings, it is crucial that we understand how the algorithms can interact, and the dynamics that can arise from their interactions. In recent years, there has been a resurgence in research efforts on multi-agent learning, and learning in games. The recent interest in adversarial learning techniques also serves to show how game theoretic tools can be being used to *robustify* and improve the performance of machine learning algorithms. Despite this activity, however, machine learning algorithms are still being treated as black-box approaches and being naïvely deployed in settings where other algorithms are actively changing the environment. In general, outside of highly structured settings, there exists no guarantees on the performance or limiting behaviors of learning algorithms in such settings.

Indeed, previous work on understanding the collective behavior of coupled learning algorithms, either in competitive or cooperative settings, has mainly looked at games where the global structure is well understood like bilinear games [19, 23, 25, 44], convex games [27, 40], or potential games [28], among many others. Such games are more conducive to the statement of global convergence guarantees since the assumed global structure can be exploited.

In games with fewer assumptions on the players' costs, however, there is still a lack of understanding of the dynamics and limiting behaviors of learning algorithms. Such settings are becoming increasingly prevalent as deep learning is increasingly being used in game theoretic settings [1, 15, 17, 49].

Gradient-based learning algorithms are extremely popular in a variety of these multi-agent settings due to their versatility, ease of implementation, and dependence on local information. There are nu-

*Submitted to the editors DATE.

Funding: This work was funded by the National Science Foundation Award CNS:1656873 and the Defense Advanced Research Projects Agency Award FA8750-18-C-0101

[†]University of California, Berkeley, Berkeley, CA (mazumdar@berkeley.edu).

[‡]University of Washington, Seattle, WA (ratliff@uw.edu).

[§]University of California, Berkeley, Berkeley, CA (sastry@eecs.berkeley.edu).

merous recent papers in multi-agent reinforcement learning that employ gradient-based methods (see, e.g. [1, 15, 49]), yet even within this well-studied class of learning algorithms, a thorough understanding of their convergence and limiting behaviors in general continuous games is still lacking.

Generally speaking, in both the game theory and the machine learning communities, two of the central questions when analyzing the dynamics of learning in games are the following:

Q1. *Are all attractors of the learning algorithms employed by agents equilibria relevant to the underlying game?*

Q2. *Are all equilibria relevant to the game also attractors of the learning algorithms agents employ?*

In this paper, we provide some answers to the above questions for the class of gradient-based learning algorithms by analyzing their limiting behavior in general continuous games. In particular, we leverage the continuous time limit of the more naturally discrete multi-agent learning algorithms. This allows us to draw on the extensive theory of dynamical systems and stochastic approximation to make statements about the limiting behaviors of these algorithms in both deterministic and stochastic settings. The latter is particularly relevant since it is common for stochastic gradient methods to be used in multi-agent machine learning contexts.

Analyzing gradient-based algorithms through the lens of dynamical systems theory has recently yielded new insights into their behavior in the classical optimization setting [22, 42, 48]. We show that a similar type of analysis can also help understand the limiting behaviors of gradient-based algorithms in games. We remark, however, that there is a *fundamental difference* between the dynamics that are analyzed in much of the single-agent, gradient-based learning and optimization literature and the ones we analyze in the competitive multi-agent case: the combined dynamics of gradient-based learning schemes in games *do not necessarily correspond to a gradient flow*. This may seem a subtle point, but it turns out to be extremely important.

Gradient flows admit desirable convergence guarantees—e.g., almost sure convergence to local minimizers—due to the fact that they preclude flows with the *worst geometries* [34]. In particular, they do not exhibit non-equilibrium limiting behavior such as periodic orbits. Gradient-based learning in games, on the other hand, does not preclude such behavior. Moreover, as we show, asymmetry in the dynamics of gradient-play in games can lead to surprising behaviors such as non-relevant limiting behaviors being attracting under the flow of the game dynamics and relevant limiting behaviors, such as a subset of the Nash equilibria being almost surely avoided.

1.1. Related Work. The study of continuous games is quite extensive (see e.g. [2, 30]), though in large part the focus has been on games admitting a fair amount of structure. The behavior of learning algorithms in games is also well-studied (see e.g. [16]). In this section, we comment on the most relevant prior work and defer a more comprehensive discussion of our results in the context of prior work to Section 6.

As we noted, previous work on learning in games in both the game theory literature, and more recently from the machine learning community, has largely focused on addressing (**Q1**) whether all attractors of the learning dynamics are game-relevant equilibria, and (**Q2**) whether all game-relevant equilibria are also attractors of the learning dynamics. The primary type of game-relevant equilibrium considered in the investigation of these two questions is a Nash equilibrium.

The majority of the existing work has focused on **Q1**. In fact, a large body of prior work focuses on games with structures that preclude the existence of non-Nash equilibria. Consequently, answering **Q1** reduces to analyzing the convergence of various learning algorithms (including gradient-play)

79 to the unique Nash equilibrium or the set of Nash equilibria. This is often shown by exploiting the
80 game structure. Examples of classes of games where gradient-play has been well-studied are potential
81 games [28], concave or monotone games [8, 27, 40], and gradient-play over the space of stochastic
82 policies in two-player finite-action bilinear games [44]. In the latter setting, other gradient-like algo-
83 rithms such as multiplicative weights have also been studied fairly extensively [19], and have been
84 shown to converge to cycling behaviors.

85 Some works have also attempted to address **Q1** in the context of gradient-play in two-player zero-
86 sum games. Concurrently with this paper, for a general class of “sufficiently smooth” two-player, zero-
87 sum games it was shown that there exists stationary points for gradient-play that are non-Nash [12]¹.
88 In such games, it has also been shown that gradient-play can converge to cycles (see, e.g., [19, 25, 47]).

89 There is also related work in more general games on the analysis of when Nash equilibria are
90 attracting for gradient-based approaches (i.e. **Q2**). Sufficient conditions for this to occur are the
91 conditions for stable differential Nash equilibria introduced in [35–37] and the condition for variational
92 stability later analyzed in [27]. We remark that these conditions are equivalent for the classes of games
93 we consider. Neither of these works give conditions under which Nash equilibria are avoided by
94 gradient-play or comment on other attracting behaviors.

95 Expanding on this rich body of literature (only the most relevant of which is covered in our short
96 review), in this paper we provide answers to **Q1** without imposing structure on the game outside
97 regularity conditions on the cost functions by exploiting the observation that gradient-based learning
98 dynamics are not gradient flows. We also provide answers to **Q2** by demonstrating that a non-trivial
99 set of games admit Nash equilibria that are almost surely avoided by gradient-play. We give explicit
100 conditions for when this occurs. Using similar analysis tools, we also provide new insights into the
101 behavior of gradient-based learning in structured classes of games such as zero-sum and potential
102 games.

103 **1.2. Contributions and Organization.** We present a general framework for modeling com-
104 petitive gradient-based learning that applies to a broad swath of learning algorithms. In Section 3,
105 we draw connections between the limiting behavior of this class of algorithms and game-theoretic
106 and dynamical systems notions of equilibria. In particular, we construct general-sum and zero-sum
107 games that admit non-Nash attracting equilibria of the gradient dynamics. Such points are attracting
108 under the learning dynamics, yet at least one player—and *potentially all of them*—has a direction in
109 which they could unilaterally deviate to decrease their cost. Thus, these non-Nash equilibria are of
110 questionable game theoretic relevance and can be seen as artifacts of the players’ algorithms.

111 In Section 4, we show that policy gradient multi-agent reinforcement learning (MARL), generative
112 adversarial networks (GANs), gradient-based multi-agent multi-armed bandits, among several other
113 common multi-agent learning settings, conform to this framework. The framework is amenable to
114 tools for analysis from dynamical systems theory.

115 Also in Section 4, we show that a subset of the local Nash equilibria in general-sum games and
116 potential games is avoided almost surely when each player employs a gradient-based algorithm. We
117 show that this holds in two broad settings: the full information setting when each player has oracle
118 access to their gradient but randomly initializes their first action, and a partial information setting
119 where each player has access to an unbiased estimate of their gradient.

¹This paper was under review at the time that [12] became publicly available. Our results show the existence of these non-Nash equilibria and attracting cycles in both general-sum and zero-sum games.

120 Thus, we provide a negative answer to both **Q1** and **Q2** for n -player general-sum games, and
 121 highlight the nuances present in zero-sum and potential games. We also show that the dynamics
 122 formed from the individual gradients of agents' costs are *not gradient flows*. This in turn implies that
 123 competitive gradient-based learning in general-sum games may converge to periodic orbits and other
 124 non-trivial limiting behaviors that arise in, e.g., chaotic systems.

125 To support the theoretical results, we present empirical results in Section 5 that show that policy
 126 gradient algorithms avoid global Nash equilibria in a large number of linear quadratic (LQ) dynamic
 127 games, a benchmark for MARL.

128 We conclude in Section 6 with a discussion of the implications of our results and some links with
 129 prior work as well as some comments on future directions.

130 **2. Preliminaries.** Consider n agents indexed by $\mathcal{I} = \{1, \dots, n\}$. Each agent $i \in \mathcal{I}$ has their
 131 own decision variable $x_i \in X_i$, where X_i is their finite-dimensional strategy space of dimension m_i .
 132 Define $X = X_1 \times \dots \times X_n$ to be the finite-dimensional joint strategy space with dimension $m =$
 133 $\sum_{i \in \mathcal{I}} m_i$. Each agent is endowed with a cost function $f_i \in C^s(X, \mathbb{R})$ with $s \geq 2$ and such that $f_i :$
 134 $(x_i, x_{-i}) \mapsto f_i(x_i, x_{-i})$ where we use the notation $x = (x_i, x_{-i})$ to make the dependence on the action
 135 of the agent x_i , and the actions of all agents excluding agent i , $x_{-i} = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$
 136 explicit. The agents seek to minimize their own cost, but only have control over their own decision
 137 variable x_i . In this setup, agents' costs are not necessarily aligned with one another, meaning they are
 138 competing.

139 Given the game $\mathcal{G} = (f_1, \dots, f_n)$, agents are assumed to update their strategies simultaneously
 140 according to a gradient-based learning algorithm of the form

$$141 \quad (2.1) \quad x_{i,t+1} = x_{i,t} - \gamma_{i,t} h_i(x_{i,t}, x_{-i,t}),$$

142 where $\gamma_{i,t}$ is agent i 's step-size at iteration t .

143 We analyze the following two settings:

- 144 1. Agents have *oracle access* to the gradient of their cost with respect to their own choice
 145 variable—i.e. $h_i(x_{i,t}, x_{-i,t}) = D_i f_i(x_{i,t}, x_{-i,t})$ where $D_i f_i \equiv \partial f_i / \partial x_i$ denotes the derivative
 146 of f_i with respect to x_i .
- 147 2. Agents have an *unbiased estimator* of their gradient—i.e., $h_i(x_{i,t}, x_{-i,t}) = D_i f_i(x_{i,t}, x_{-i,t}) +$
 148 $w_{i,t+1}$ where $\{w_{i,t}\}$ is a zero mean, finite variance stochastic process.

149 We refer to the former setting as *deterministic* gradient-based learning and the latter setting as *stochas-*
 150 *tic* gradient-based learning. Assuming that all agents are employing such algorithms, we aim to ana-
 151 lyze the limiting behavior of the agents' strategies. To do so, we leverage the following game-theoretic
 152 notion of a Nash equilibrium.

153 **Definition 2.1.** A strategy $x \in X$ is a *local Nash equilibrium* for the game (f_1, \dots, f_n) if, for
 154 each $i \in \mathcal{I}$, there exists an open set $W_i \subset X_i$ such that that $x_i \in W_i$ and $f_i(x_i, x_{-i}) \leq f_i(x'_i, x_{-i})$
 155 for all $x'_i \in W_i$. If the above inequalities are strict, then we say x is a *strict local Nash equilibrium*.

156 The focus on *local* Nash equilibria is due to our lack of assumptions on the agents' cost functions.
 157 If $W_i = X_i$ for each i , then a local Nash equilibrium x is a global Nash equilibrium. This holds in
 158 e.g the bimatrix games and the linear quadratic games we analyze in Section 5. Depending on the
 159 agents' costs, a game (f_1, \dots, f_n) may admit anywhere from one to a continuum of local or global
 160 Nash equilibria; or none at all.

161 **3. Linking Games and Dynamical Systems.** In this section, we draw links between the
 162 limiting behavior of dynamical systems and game-theoretic notions of equilibria in three broad classes
 163 of continuous games. For brevity, the proofs of the propositions in this section are supplied in Ap-
 164 pendix A. A high-level summary of the links we draw is shown in Figure 1.

165 Define $\omega(x) = (D_1 f_1(x), \dots, D_n f_n(x))$ to be the vector of player derivatives of their own cost
 166 functions with respect to their own choice variables. When each player is employing a gradient-based
 167 learning algorithm, the joint strategy of the players, (in the limit as the agents' step-sizes go to zero)
 168 follows the differential equation

$$169 \quad \dot{x} = -\omega(x).$$

171 A point $x \in X$ is said to be an equilibrium, critical point, or stationary point of the dynamics
 172 if $\omega(x) = 0$. Stationary points of $\dot{x} = -\omega(x)$ are joint strategies from which, under gradient-play,
 173 the agents do not move. We note that $\omega(x) = 0$ is a necessary condition for a point $x \in X$ to be a
 174 local Nash equilibrium [37]. Hence, all local Nash equilibria are critical points of the joint dynamics
 175 $\dot{x} = -\omega(x)$.

176 Central to dynamical systems theory is the study of limiting behavior and its stability properties.
 177 A classical result in dynamical systems theory allows us to characterize the stability properties of an
 178 equilibrium x^* by analyzing the Jacobian of the dynamics at x^* . The Jacobian of ω is defined by

$$179 \quad D\omega(x) = \begin{bmatrix} D_1^2 f_1(x) & \cdots & D_{n1} f_1(x) \\ \vdots & \ddots & \vdots \\ D_{1n} f_n(x) & \cdots & D_n^2 f_n(x) \end{bmatrix}.$$

180 Since $D\omega$ is a matrix of second derivatives, it is sometimes referred to as the ‘game Hessian’. Similar
 181 to the Hessian matrix of a gradient flow, $D\omega$ allows us to further characterize the critical points of ω by
 182 their properties under the flow of $\dot{x} = -\omega(x)$. Let $\lambda_i(x) \in \text{spec}(D\omega(x))$ for $i \in \{1, \dots, m\}$ denote
 183 the eigenvalues of $D\omega$ at x where $\text{Re}(\lambda_1(x)) \leq \dots \leq \text{Re}(\lambda_m(x))$ —that is, $\lambda_1(x)$ is the eigenvalue
 184 with the smallest real part. Of particular interest are asymptotically stable equilibria.

185 **Definition 3.1.** *A point $x \in X$ is a locally asymptotically stable equilibrium of the continuous*
 186 *time dynamics $\dot{x} = -\omega(x)$ if $\omega(x) = 0$ and $\text{Re}(\lambda) > 0$ for all $\lambda \in \text{spec}(D\omega(x))$.*

187 Locally asymptotically stable equilibria have two properties of interest. First, they are isolated,
 188 meaning that there exists a neighborhood around them in which no other equilibria exist. Second,
 189 they are exponentially attracting under the flow of $\dot{x} = -\omega(x)$, meaning that if agents initialize in a
 190 neighborhood of a locally asymptotically stable equilibrium x^* and follow the dynamics described by
 191 $\dot{x} = -\omega(x)$, they will converge to x^* exponentially fast [41]. This, in turn, implies that a discretized
 192 version of $\dot{x} = -\omega(x)$, namely

$$193 \quad (3.1) \quad x_{t+1} = x_t - \gamma\omega(x_t),$$

194 converges locally for appropriately selected step size γ at a rate of $O(1/t)$. Such results motivate
 195 the study of the continuous time dynamical system $\dot{x} = -\omega(x)$ in order to understand convergence
 196 properties of gradient-based learning algorithms of the form (2.1).

197 Another important class of critical points of a dynamical system are saddle points.

198 **Definition 3.2.** A point $x \in X$ is a saddle point of the dynamics $\dot{x} = -\omega(x)$ if $\omega(x) = 0$
 199 and $\lambda_1(x) \in \text{spec}(D\omega(x))$ is such that $\text{Re}(\lambda_1(x)) \leq 0$. A saddle point such that $\text{Re}(\lambda_i) < 0$ for
 200 $i \in \{1, \dots, \ell\}$ and $\text{Re}(\lambda_j) > 0$ for $j \in \{\ell + 1, \dots, m\}$ with $0 < \ell < m$ is a strict saddle point of the
 201 continuous time dynamics $\dot{x} = -\omega(x)$.

202 Strict saddle points are especially relevant to our analysis since their neighborhoods are character-
 203 ized by stable and unstable manifolds [41]. When the agents evolve according to the dynamics solely
 204 on the stable manifold, they converge exponentially fast to the critical point. However, when they
 205 evolve solely on the unstable manifold, they diverge from the equilibrium exponentially fast. Agents
 206 whose strategies lie on the union of the two manifolds asymptotically avoid the equilibrium. We make
 207 use of this general fact in Section 4.1.

208 To better understand the links between the critical points of the gradient dynamics and the Nash
 209 equilibria of the game, we make use of an equivalent characterization of strict local Nash that leverages
 210 first and second order conditions on player cost functions. This makes them simpler objects to link to
 211 the various dynamical systems notions of equilibria than local Nash equilibria.

212 **Definition 3.3** ([35, 37]). A point $x \in X$ is a differential Nash equilibrium for the game defined
 213 by (f_1, \dots, f_n) if $\omega(x) = 0$ and $D_i^2 f_i(x) \succ 0$ for each $i \in \mathcal{I}$.

214 In [36], it was shown that local Nash equilibria are generically differential Nash equilibria where
 215 $\det(D\omega(x)) \neq 0$ (i.e., $D\omega$ is non-degenerate). Thus, in the space of games where the agents' costs
 216 are at least twice differentiable, the set of games that admit local Nash equilibria that are not non-
 217 degenerate differential Nash equilibria is of measure zero [36]. In [36] it was also shown that non-
 218 degenerate Nash equilibria are structurally stable, meaning that small perturbations to the agents'
 219 costs functions will not change the fundamental nature of the equilibrium. This also implies that
 220 gradient-play with slightly biased estimators of the gradient will not have vastly different behaviors in
 221 neighborhoods of equilibria.

222 Given these different equilibrium notions of the learning dynamics and the underlying game, let us
 223 define the following sets which will be useful in stating the results in the following sections. For a game
 224 $\mathcal{G} = (f_1, \dots, f_n)$, denote the sets of strict saddle points and locally asymptotically stable equilibria
 225 of the gradient dynamics, $\dot{x} = -\omega(x)$, as $\text{SSP}(\omega)$ and $\text{LASE}(\omega)$, respectively, where we recall that
 226 $\omega(x) = (D_1 f_1(x), \dots, D_n f_n(x))$. Similarly, denote the set of local Nash equilibria, differential Nash
 227 equilibria, and non-degenerate differential Nash equilibria of \mathcal{G} as $\text{LNE}(\mathcal{G})$, $\text{DNE}(\mathcal{G})$, and $\text{NDDNE}(\mathcal{G})$,
 228 respectively. As previously mentioned, $\text{NDDNE}(\mathcal{G}) = \text{LNE}(\mathcal{G})$ in almost all continuous games. The key
 229 takeaways of this section are summarized in Figure 1.

230 **3.1. General-sum games.** We first analyze the properties of local Nash equilibria under the
 231 joint gradient dynamics in n -player general-sum games.

232 **Proposition 3.4.** A non-degenerate differential Nash equilibrium is either a locally asymptoti-
 233 cally stable equilibrium or a strict saddle point of $\dot{x} = -\omega(x)$ —i.e., $\text{NDDNE}(\mathcal{G}) \subset \text{SSP}(\omega) \cup \text{LASE}(\omega)$.
 234

235 Locally asymptotically stable differential Nash equilibria satisfy the notion of variational stability
 236 introduced in [27]. In fact, a simple analysis shows that the definitions of variationally stable equilibria
 237 and locally asymptotically stable differential Nash equilibria [35] are equivalent in the games we
 238 consider—i.e., games where each players' cost is at least twice continuously differentiable. We remark

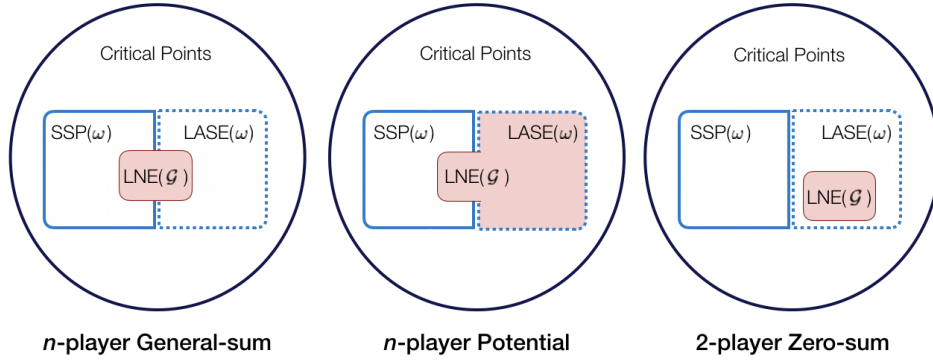


Figure 1: Links between the equilibria of generic continuous games \mathcal{G} and their properties under the gradient dynamics $\dot{x} = -\omega(x)$.

239 that, from the definition of asymptotic stability, the gradient dynamics have a $O(1/t)$ convergence rate
 240 in the neighborhood of such equilibria.

241 An important point to make is that not every locally asymptotically stable equilibrium of $\dot{x} =$
 242 $-\omega(x)$ is a non-degenerate differential Nash equilibrium. Indeed, the following proposition provides
 243 an entire class of games whose corresponding gradient dynamics admit locally asymptotically stable
 244 equilibria that are not local Nash equilibria.

245 **Proposition 3.5.** *In the class of general-sum continuous games, there exists a continuum of*
 246 *games containing games \mathcal{G} such that $\text{LASE}(\omega) \not\subset \text{NDDNE}(\mathcal{G})$, and moreover, $\text{LASE}(\omega) \not\subset \text{LNE}(\mathcal{G})$.*
 247

248 *Proof.* Consider a two player game $\mathcal{G} = (f_1, f_2)$ on \mathbb{R}^2 where

$$249 \quad f_1(x_1, x_2) = \frac{a}{2}x_1^2 + bx_1x_2, \quad \text{and} \quad f_2(x_1, x_2) = \frac{d}{2}x_2^2 + cx_1x_2$$

251 for constants $a, b, c, d \in \mathbb{R}$. The Jacobian of ω is given by

$$252 \quad (3.2) \quad D\omega(x_1, x_2) = \begin{bmatrix} a & b \\ c & d \end{bmatrix}, \quad \forall (x_1, x_2) \in \mathbb{R}^2.$$

254 If $a > 0$ and $d < 0$, then the unique stationary point $x = (0, 0)$ is neither a differential Nash nor a
 255 local Nash equilibria since the necessary conditions are violated (i.e., $d < 0$). However, if $a > -d$ and
 256 $ad > cb$, the eigenvalues of $D\omega$ have positive real parts and $(0, 0)$ is asymptotically stable. Further,
 257 this clearly holds for a continuum of games. Thus, the set of locally asymptotically stable equilibria
 258 that are not Nash equilibria may be arbitrarily large. ■

259 The preceding proposition shows that there exists attracting critical points of the gradient dynam-
 260 ics in general-sum continuous games that are not Nash equilibria and may not be even relevant to the
 261 game. Thus, this provides a negative answer to **Q2** (whether all attracting equilibria in general-games
 262 are game-relevant for the learning dynamics).

263 *Remark 3.6.* We note that, by definition, the non-Nash locally asymptotically stable equilibria (or
 264 non-Nash equilibria) do not satisfy the second-order conditions for Nash equilibria. Thus, at these
 265 joint strategies, at least one player – and maybe all of them – has a direction in which they would
 266 unilaterally deviate if they were not using gradient descent. As such, we view convergence to these
 267 points to be undesirable.

268 **3.2. Zero-sum games.** Let us now restrict our attention to two-player zero-sum games, which
 269 often arise when training GANs, in adversarial learning, and in MARL [9, 17, 29]. In such games, one
 270 player can be seen as minimizing f with respect to their decision variable and the other as minimizing
 271 $-f$ with respect to theirs. The following proposition shows that all differential Nash equilibria in
 272 two-player zero-sum games are locally asymptotically stable equilibria under the flow of $\dot{x} = -\omega(x)$.

273 **Proposition 3.7.** *For an arbitrary two-player zero-sum game, $(f, -f)$ on \mathbb{R}^m , if x is a differ-*
 274 *ential Nash equilibrium, then x is both a non-degenerate differential Nash equilibrium and a locally*
 275 *asymptotically stable equilibrium of $\dot{x} = -\omega(x)$ —that is, $\text{DNE}(\mathcal{G}) \equiv \text{NDDNE}(\mathcal{G}) \subset \text{LASE}(\omega)$.*

276 This result guarantees that the differential Nash equilibria of zero-sum games are isolated and
 277 exponentially attracting under the flow of $\dot{x} = -\omega(x)$. This in turn guarantees that simultaneous
 278 gradient-play has a local linear rate of convergence to all local Nash equilibria in all zero-sum con-
 279 tinuous games. Thus, the answer to **Q1** in the context of zero-sum games is “yes”, since all Nash
 280 equilibria are attracting for the gradient dynamics.

281 The converse of the preceding proposition, however, is not true. Not every locally asymptotically
 282 stable equilibrium in two-player zero-sum games are non-degenerate differential Nash equilibria. In-
 283 deed, there may be many locally asymptotically stable equilibria in a zero-sum game that are not local
 284 Nash equilibria. The following proposition highlights this fact.

285 **Proposition 3.8.** *In the class of zero-sum continuous games, there exists a continuum of games*
 286 *such that for each game \mathcal{G} , $\text{LASE}(\omega) \not\subset \text{DNE}(\mathcal{G}) \subset \text{LNE}(\mathcal{G})$.*

287 *Proof.* Consider the two-player zero-sum game $(f, -f)$ on \mathbb{R}^2 where

$$288 \quad f(x_1, x_2) = \frac{a}{2}x_1^2 + bx_1x_2 + \frac{c}{2}x_2^2;$$

289 and $a, b, c \in \mathbb{R}$. The Jacobian of ω is given by

$$291 \quad D\omega(x_1, x_2) = \begin{bmatrix} a & b \\ -b & -c \end{bmatrix}, \quad \forall (x_1, x_2) \in \mathbb{R}^2.$$

292 If $a > c > 0$ and $b^2 > ac$, then $D\omega(x_1, x_2)$ has eigenvalues with strictly positive real part, but the
 293 unique stationary point is not a differential Nash equilibrium—since $-c < 0$ —and, in fact, is not even
 294 a Nash equilibrium. Indeed,

$$295 \quad -f(0, 0) > -f(0, x_2) = -\frac{c}{2}x_2^2, \quad \forall x_2 \neq 0.$$

296 Thus, there exists a continuum of zero-sum games with a large set of locally asymptotically stable
 297 equilibria of the corresponding dynamics $\dot{x} = -\omega(x)$ that are not differential Nash. ■

298 The preceding proposition again shows that there exists non-Nash equilibria of the gradient dy-
 299 namics in zero-sum continuous games. Thus, this proposition also provides a negative answer to **Q2**
 300 in the context of zero-sum games.

301 **3.3. Potential Games.** One last set of games with interesting connections between the Nash
 302 equilibria and the critical points of the gradient dynamics is the class known as *potential games*. This
 303 particularly nice class of games are ones for which ω corresponds to a gradient flow under a coordinate
 304 transformation—that is, there exists a function ϕ (commonly referred to as the potential function)
 305 such that for each $i \in \mathcal{I}$, $D_i f_i \equiv D_i \phi$. We remark that due to the equivalence this class of games
 306 is sometimes referred to as an *exact* potential game. Note that a necessary and sufficient condition
 307 for (f_1, \dots, f_n) to be a potential game is that $D\omega$ is *symmetric* [28]—that is, $D_{ij} f_j \equiv D_{ji} f_i$. This
 308 gives potential games the desirable property that the only locally asymptotically stable equilibria of
 309 the gradient dynamics are local Nash equilibria.

310 **Proposition 3.9.** *For an arbitrary potential game, $\mathcal{G} = (f_1, \dots, f_n)$ on \mathbb{R}^m , if x is a locally*
 311 *asymptotically stable equilibrium of $\dot{x} = -\omega(x)$ (i.e., $x \in \text{LASE}(\omega)$), then x is a non-degenerate*
 312 *differential Nash equilibrium (i.e., $x \in \text{NDDNE}(\mathcal{G})$).*

313 The full proof of Proposition 3.9 is supplied in Appendix A. The preceding proposition rules out
 314 non-Nash locally asymptotically stable equilibria of the gradient dynamics in potential games, and
 315 implies that every local minimum of a potential game must be a local Nash equilibrium. Thus, in
 316 potential games, unlike in general-sum and zero-sum games, the answer to **Q2** is positive. However,
 317 the following proposition shows that the existence of a potential function is not enough to rule out
 318 local Nash equilibria that are saddle points of the dynamics.

319 **Proposition 3.10.** *In the class of continuous games, there exist a continuum of potential games*
 320 *containing games \mathcal{G} that admit Nash equilibria that are saddle points of the dynamics $\dot{x} = -\omega(x)$ —*
 321 *i.e., $\exists \mathcal{G}$ such that for some $x \in \text{LNE}(\mathcal{G})$, $x \in \text{SSP}(\omega)$.*

322 *Proof.* Consider the game (f, f) on $X = \mathbb{R}^2$ described by

$$323 \quad f(x_1, x_2) = \frac{a}{2}x_1^2 + bx_1x_2 + \frac{c}{2}x_2^2$$

324 where $a, b, d \in \mathbb{R}$. The Jacobian of ω is given by

$$325 \quad D\omega(x_1, x_2) = \begin{bmatrix} a & b \\ b & c \end{bmatrix}, \quad \forall (x_1, x_2) \in \mathbb{R}^2.$$

326 If $a, c > 0$, then $x = (0, 0)$ is a local Nash equilibrium. However, if $ac < b^2$, $D\omega(x)$ has one positive
 327 and one negative eigenvalue and $(0, 0)$ is a saddle point of the gradient dynamics. Thus, there exists a
 328 continuum of potential games where a large set of differential Nash equilibria are strict saddle points
 329 of $\dot{x} = -\omega(x)$. ■

330 Proposition 3.10 demonstrates a surprising fact about potential games. Even though all minimizers
 331 of the potential function must be local Nash equilibria, *not all local Nash equilibria are minimizers of*
 332 *the potential function.*

333 **3.4. Main Takeaways.** The main takeaways of this section are summarized in Figure 1. We
 334 note that for zero-sum games, Proposition 3.8 shows that $\text{LNE}(\mathcal{G}) \subset \text{LASE}(\omega)$. Since the inclusion is
 335 strict, the answer to **Q2** in such games is “no”. For general-sum games, Proposition 3.5 allows us to
 336 to conclude that there do exist attracting, non-Nash equilibria. Thus, the answer to **Q2** is also “no”. In
 337 potential games, since $\text{LASE}(\omega) \subset \text{LNE}(\mathcal{G})$ the answer is “yes”.

338 In the following sections, we provide answers to **Q1** by showing that all local Nash equilibria and
 339 $\text{LNE}(\mathcal{G}) \cap \text{SSP}(\omega)$ are avoided almost surely by gradient-based algorithms in both the deterministic and
 340 stochastic settings. In particular, since $\text{LNE}(\mathcal{G}) \cap \text{SSP}(\omega) \neq \emptyset$ in potential and general-sum games, one
 341 cannot give a positive answer to **Q1** in either of these classes of games.

342 **4. Convergence of Gradient-Based Learning.** In this section, we provide convergence and
 343 non-convergence results for gradient-based algorithms. We also include a high-level overview of well-
 344 known algorithms that fit into the class of learning algorithms we consider; more detail can be found
 345 in Appendix C.

346 **4.1. Deterministic Setting.** We first address convergence to equilibria in the *deterministic*
 347 setting in which agents have oracle access to their gradients at each time step. This includes the
 348 case where agents know their own cost functions f_i and observe their own actions as well as their
 349 competitors' actions—and hence, can compute the gradient of their cost with respect to their own
 350 choice variable.

351 Since we have assumed that each agent $i \in \mathcal{I}$ has their own *learning rate* (i.e. step sizes γ_i), the
 352 joint dynamics of all the players are given by

$$353 \quad (4.1) \quad x_{t+1} = g(x_t)$$

354 where $g : x \mapsto x - \gamma \odot \omega(x)$ with $\gamma = (\gamma_i)_{i \in \mathcal{I}}$ and $\gamma > 0$ element-wise. By a slight abuse of notation,
 355 $\gamma \odot \omega(x_t)$ is defined to be element-wise multiplication of γ and $\omega(\cdot)$ where γ_1 is multiplied by the first
 356 m_1 components of $\omega(\cdot)$, γ_2 is multiplied by the next m_2 components, and so on.

357 We remark that this update rule immediately distinguishes gradient-based learning in games from
 358 gradient descent. By definition, the dynamics of gradient descent in single-agent settings always
 359 correspond to gradient flows—i.e. x evolves according to an ordinary differential equation of the form
 360 $\dot{x} = -\nabla \phi(x)$ for some function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$. Outside of the class of *exact* potential games we defined
 361 in Section 3, the dynamics of players' actions in games are not afforded this luxury—indeed, $D\omega$ is
 362 not in general symmetric (which is a necessary condition for a gradient flow). This makes the potential
 363 limiting behaviors of $\dot{x} = -\omega(x)$ highly non-trivial to characterize in general-sum games.

364 The structure present in a gradient-flow implies strong properties on the limiting behaviors of x . In
 365 particular, it precludes the existence of limit cycles or periodic orbits (limiting behaviors of dynamical
 366 systems where the state of system cycles infinitely through a set of states with a finite period) and
 367 chaos (an attribute of nonlinear dynamical systems where the system's behavior can vary extremely
 368 due to slight changes in initial position) [41]. We note that both of these behaviors can occur in the
 369 dynamics of gradient-based learning algorithms in games².

370 Despite the wide breadth of behaviors that gradient dynamics can exhibit in competitive settings,
 371 we are still make statements about convergence (and non-convergence) to certain types of equilibria.
 372 To do so, we first make the following standard assumptions on the smoothness of the cost functions f_i
 373 and the magnitude of the agents' learning rates γ_i .

374 **Assumption 1.** For each $i \in \mathcal{I}$, $f_i \in C^s(X, \mathbb{R})$ with $s \geq 2$, $\sup_{x \in X} \|D\omega(x)\|_2 \leq L < \infty$, and
 375 $0 < \gamma_i < 1/L$ where $\|\cdot\|_2$ is the induced 2-norm.

²The Van der Pol oscillator and Lorenz system (see e.g. [41]) can be seen as the resulting gradient dynamics in a 2-player and 3-player general-sum game respectively. The first is a classic example of a system where players converge to cycles and the second is an example of a chaotic system.

376 Given these assumptions, the following result rules out converging to strict saddle points.

377 **Theorem 4.1.** *Let $f_i : X \rightarrow \mathbb{R}$ and γ satisfy Assumption 1. Suppose that $X = X_1 \times \cdots \times X_n \subseteq$
378 \mathbb{R}^m is open and convex. If $g(X) \subset X$, the set of initial conditions $x \in X$ from which competitive
379 gradient-based learning converges to strict saddle points is of measure zero.*

380 We remark that the above theorem holds for $X = X_1 \times \cdots \times X_n = \mathbb{R}^m$ in particular, since
381 $g(X) \subset X$ holds trivially in this case. It is also important to note that, as we point out in Section 3,
382 local Nash equilibria can be strict saddle points. Thus, all local Nash equilibria that are strict saddle
383 points for $\dot{x} = -\omega(x)$ are avoided almost surely by gradient-play even with oracle gradient access and
384 random initializations. This holds even when players randomly initialize uniformly in an arbitrarily
385 small ball around such Nash equilibria. In Section 5, we show that many linear quadratic dynamic
386 games have a strict saddle point as their global Nash equilibrium. For brevity, we provide the proof of
387 Theorem 4.1 in Appendix A, and provide a proof sketch below.

388 *Proof sketch of Theorem 4.1.* The core of the proof is the celebrated stable manifold theorem from
389 dynamical systems theory, presented in Theorem A.1. We construct the set of initial positions from
390 which gradient-play will converge to strict saddle points and then use the stable manifold theorem
391 to show that the set must have measure zero in the players' joint strategy space. Therefore, with a
392 random initialization players will never evolve solely on the stable manifold of strict saddles and they
393 will consequently diverge from such equilibria.

394 To be able to invoke the stable manifold theorem, we first show that the mapping $g : \mathbb{R}^m \rightarrow \mathbb{R}^m$
395 is a diffeomorphism, which is non-trivial due to the fact that we have allowed each agent to have their
396 own learning rate γ_i and $D\omega$ is not symmetric. We then iteratively construct the set of initializations
397 that will converge to strict saddle points under the game dynamics. By the stable manifold theorem,
398 and the fact that g is a diffeomorphism, the stable manifold of a strict saddle point must be measure
399 zero. Then, by induction we show that the set of all initial points that converge to a strict saddle point
400 must also be measure zero. ■

401 In potential games we can strengthen the above non-convergence result and give convergence
402 guarantees.

403 **Corollary 4.2.** *Consider a potential game (f_1, \dots, f_n) on open, convex $X = X_1 \times \cdots \times X_n \subseteq$
404 \mathbb{R}^m and where each $f_i \in C^s(X, \mathbb{R})$ for $s \geq 3$. Let ν be a prior measure with support X which is
405 absolutely continuous with respect to the Lebesgue measure and assume $\lim_{t \rightarrow \infty} g^t(x)$ exists. Then,
406 under Assumption 1, competitive gradient-based learning converges to non-degenerate differential
407 Nash equilibria almost surely. Moreover, the non-degenerate differential Nash to which it converges is
408 generically a local Nash equilibrium.*

409 Corollary 4.2 guarantees that in potential games, gradient-play will converge to a differential Nash
410 equilibrium. Combining this with Theorem 4.1 guarantees that the differential Nash equilibrium it
411 converges to is a local minimizer of the potential function. A simple implication of this result is that
412 gradient-based learning in potential games cannot exhibit limit cycles or chaos.

413 Of note is the fact that the agents *do not* need to be performing gradient-based learning on ϕ
414 to converge to Nash almost surely. That is, they do not need to know the function ϕ ; they simply
415 need to follow the derivative of their own cost with respect to their own choice variable, and they are
416 guaranteed to converge to a local Nash equilibrium that is a local minimizer of the potential function.

417 We note that convergence to Nash equilibria is a known characteristic of gradient-play in potential

418 games. However, our analysis also highlights that gradient-play will avoid a subset of the Nash equi-
 419 libria of the game. This is surprising given the particularly strong structural properties of such games.
 420 The proof for Corollary 4.2 is provided in Appendix A and follows from Proposition 3.9, Theorem 4.1,
 421 and the fact that $D\omega$ is symmetric in potential games.

422 **4.1.1. Implications and Interpretation of Convergence Analysis.** Both Theorem 4.1
 423 and Corollary 4.2 show that gradient-play in multi-agent settings avoids strict saddles almost surely
 424 even in the deterministic setting. Combined with the analysis in Section 3 which shows that (local)
 425 Nash equilibria can be strict saddles of the dynamics for general-sum games, this implies that a subset
 426 of the Nash equilibria are almost surely avoided by individual gradient-play, a potentially undesir-
 427 able outcome in view of **Q1** (whether all Nash equilibria are attracting for the learning dynamics). In
 428 Section 5, we show that the global Nash equilibrium is a saddle point of the gradient dynamics in a
 429 large number of randomly sampled LQ dynamic games. This suggests that policy gradient algorithms
 430 may fail to converge in such games, which is highly undesired. This is in stark contrast to the sin-
 431 gle agent setting where policy gradient has been shown to converge to the unique solution of LQR
 432 problems [13].

433 In Section 3, we also showed that local Nash equilibria of potential games can be strict saddles
 434 points of the potential function. Non-convergence to such points in potential games is not necessarily
 435 a bad result since this in turn implies convergence to a local minimizer of the potential function (as
 436 shown in [22,32]) which are guaranteed to be local Nash equilibria of the game. However, these results
 437 do imply that *one cannot answer “yes” to Q1 in potential games* since some of the Nash equilibria are
 438 not attracting under gradient-play.

439 In zero-sum games, where local Nash equilibria cannot be strict saddle points of the gradient
 440 dynamics, our result suggests that *eventually* gradient-based learning algorithms will escape saddle
 441 points of the dynamics.

442 The almost sure avoidance of all equilibria that are saddle points of the dynamics further implies
 443 that if (3) converges to a critical point x , then $x \in \text{LASE}(\omega)$ —i.e., x is locally asymptotically stable for
 444 $\dot{x} = -\omega(x)$. This may not be a desired property however, since we showed in Section 3 that zero-sum
 445 and general-sum games both admit non-Nash LASE.

446 Since gradient-play in games generally does not result in a gradient flow, other types of limiting
 447 behaviors such as limit cycles can occur in gradient-based learning dynamics. Theorem 4.1 says
 448 nothing about convergence to other limiting behaviors. In the following sections we prove that the
 449 results described in this section extend to the stochastic gradient setting. We also formally define
 450 periodic orbits in the context of dynamical systems and state stronger results on avoidance of some
 451 more complex limiting behaviors like linearly unstable limit cycles.

452 **4.2. Stochastic Setting.** We now analyze the stochastic case in which agents are assumed to
 453 have an unbiased estimator for their gradient. The results in this section allow us to extend the results
 454 from the deterministic setting to a setting where each agent builds an estimate of the gradient of their
 455 loss at the current set of strategies from potentially noisy observations of the environment. Thus, we
 456 are able to analyze the limiting behavior of a class of commonly used machine learning algorithms for
 457 competitive, multi-agent settings. In particular, we show that agents will almost surely not converge
 458 to strict saddle points. In Appendix B.1, we show that the gradient dynamics will actually avoid more
 459 general limiting behaviors called linearly unstable cycles which we define formally.

460 To perform our analysis, we make use of tools and ideas from the literature on stochastic approxi-

Class	Gradient Learning Rule
Gradient-Play	$x_i^+ = x_i - \gamma_i D_i f_i(x_i, x_{-i})$
GANs	$\theta^+ = \theta - \gamma \mathbb{E}[D_\theta L(\theta, w)]$ $w^+ = w + \gamma \mathbb{E}[D_w L(\theta, w)]$
MA Policy Gradient	$x_i^+ = x_i - \gamma_i \mathbb{E}[D_i J_i(x_i, x_{-i})]$
Individual Q-learning	$q_i^+(u_i) = q_i(u_i) + \gamma_i (r_i(u_i, \pi_{-i}(q_i, q_{-i})) - q_i(u_i))$
MA Gradient Bandits	$x_{i,\ell}^+ = x_{i,\ell} + \gamma_i \mathbb{E}[\beta_i R_i(u_i, u_{-i}) u_i = \ell], \ell = 1, \dots, m_i$
MA Experts	$x_{i,\ell}^+ = x_{i,\ell} + \gamma_i \mathbb{E}[R_i(u_i, u_{-i}) u_i = \ell], \ell = 1, \dots, m_i$

Table 1: Example problem classes that fit into competitive gradient-based learning rules. Details on the derivation of these update rules as gradient-based learning schemes is provided in Appendix C.

461 mations (see e.g [6]). We note that the convergence of stochastic gradient schemes in the single-agent
 462 setting has been extensively studied [7,26,33,38]. We extend this analysis to the behavior of stochastic
 463 gradient algorithms in games.

464 We assume that each agent updates their strategy using the update rule

465 (4.2)
$$x_{i,t+1} = x_{i,t} - \gamma_{i,t} (D_i f_i(x_{i,t}, x_{-i,t}) + w_{i,t+1})$$

466 for some zero-mean, finite-variance stochastic process $\{w_{i,t}\}$. Before presenting the results for the
 467 stochastic case, let us comment on the different learning algorithms that fit into this framework.

468 **4.2.1. Examples of Stochastic Gradient-Based Learning.** The stochastic gradient-based
 469 learning setting we study is general enough to include a variety of commonly used multi-agent learning
 470 algorithms. The classes of algorithms we include is hardly an exhaustive list, and indeed many exten-
 471 sions and altogether different algorithms exist that can be considered members of this class. In Table 1,
 472 we provide the gradient-based update rule for six different example classes of learning problems: (i)
 473 gradient-play in non-cooperative continuous games, (ii) GANs, (iii) multi-agent policy gradient, (iv)
 474 individual Q-learning, (v) multi-agent gradient bandits, and (vi) multi-agent experts. We provide a
 475 detailed analysis of these different algorithms including the derivation of the gradient-based update
 476 rules along with some interesting numerical examples in Appendix C. In each of these cases, one can
 477 view an agent employing the given algorithm as building an unbiased estimate of their gradient from
 478 their observation of the environment.

479 For example, in multi-agent policy gradient (see, e.g., [46, Chapter 13]), agents' costs are defined
 480 as functions of a parameter vector x_i that parameterize their policies $\pi_i(x_i)$. The parameters x_i are
 481 agent i 's choice variable. By following the gradient of their loss function, they aim to tune the param-
 482 eters in order to converge to an *optimal* policy π_i . Perhaps surprisingly, it is not necessary for agent
 483 i to have access to $\pi_{-i}(x_{-i})$ or even x_{-i} in order for them to construct an unbiased estimate of the
 484 gradient of their loss with respect to their own choice variable x_i as long as they observe the sequence
 485 of actions, say $u_{-i,t}$, of all other agents generated. These actions are implicitly determined by the other
 486 agents' policies $\pi_{-i}(x_{-i})(\cdot)$. Hence, in this case if agent i observes $\{(r_{j,t}, u_{j,t}, s_{j,t}), \forall j \in \mathcal{I}\}$ where

487 (r_j, u_j, s_j) are the reward, action, and state of agent j , then this is enough to construct an unbiased
 488 estimate of their gradient. We provide further details on multi-agent policy gradient in Appendix C.

489 **4.2.2. Stochastic Gradient Results.** Returning to the analysis of (4.2), we make the follow-
 490 ing standard assumptions on the noise processes [38, 39].

491 **Assumption 2.** *The stochastic process $\{w_{i,t+1}\}$ satisfies the assumptions $\mathbb{E}[w_{i,t+1} | \mathcal{F}_i^t] = 0$,
 492 $t \geq 0$ and $\mathbb{E}[\|w_{i,t+1}\|^2 | \mathcal{F}_i^t] \leq \sigma^2 < \infty$ a.s., for $t \geq 0$, where $\mathcal{F}_{i,t}$ is an increasing family of
 493 σ_i -fields—i.e. filtration, or history generated by the sequence of random variables—given by $\mathcal{F}_{i,t} =$
 494 $\sigma_i(x_{i,k}, w_{i,k}, k \leq t)$, $t \geq 0$.*

495 We also make new assumptions on the players’ step-sizes. These are standard assumptions in the
 496 stochastic approximation literature and are needed to ensure that the noise processes are asymptotically
 497 controlled.

498 **Assumption 3.** *For each $i \in \mathcal{I}$, $f_i \in C^s(X, \mathbb{R})$ with $s \geq 2$, $D_i f_i$ is L_i -Lipschitz with $0 <$
 499 $L_i < \infty$, the step-sizes satisfy $\gamma_{i,t} \equiv \gamma_t$ for all $i \in \mathcal{I}$ and $\sum_t \gamma_t = \infty$ and $\sum_t (\gamma_t)^2 < \infty$, and
 500 $\sup_t \|x_t\| < \infty$ a.s.*

501 Let $(a)^+ = \max\{a, 0\}$ and $a \cdot b$ denotes the inner product. The following theorem extends the results
 502 of Theorem 4.1 to the stochastic gradient dynamics in games.

503 **Theorem 4.3.** *Consider a game (f_1, \dots, f_n) on $X = X_1 \times \dots \times X_n = \mathbb{R}^m$. Suppose each agent
 504 $i \in \mathcal{I}$ adopts a stochastic gradient algorithm that satisfies Assumptions 2 and 3. Further, suppose that
 505 for each $i \in \mathcal{I}$, there exists a constant $b_i > 0$ such that $\mathbb{E}[(w_{i,t} \cdot v)^+ | \mathcal{F}_{i,t}] \geq b_i$ for every unit vector
 506 $v \in \mathbb{R}^{m_i}$. Then, competitive stochastic gradient-based learning converges to strict saddle points of
 507 the game on a set of measure zero.*

508 The proof follows directly from showing that (4.2) satisfies Theorem A.2, provided the assumptions
 509 of the theorem hold. The assumption that $\mathbb{E}[(w_{i,t} \cdot v)^+ | \mathcal{F}_{i,t}] \geq b_i$ rules out degenerate cases where the
 510 noise forces the stochastic dynamics onto the stable manifold of strict saddle points.

511 Theorem 4.3 implies that the dynamics of stochastic gradient-based learning defined in (4.2), have
 512 the same limiting properties as the deterministic dynamics vis-à-vis saddle points. Thus, the impli-
 513 cations described in Section 4.1.1 extend to the stochastic gradient setting. In particular, stochastic
 514 gradient-based algorithms will avoid a non-negligible subset of the Nash equilibria in general-sum
 515 and potential games. Further, in zero-sum and general-sum games, if the players do converge to a
 516 critical point, that point may be a non-Nash equilibrium.

517 **4.2.3. Further Convergence Results for Stochastic Gradient-Play in Games.** As we
 518 demonstrated in Section 4.1, outside of potential games, the dynamics of gradient-based learning al-
 519 gorithms in games are not gradient flows. As such, the players’ actions can converge to more complex
 520 sets than simple equilibria. A particularly prominent class of limiting behaviors for dynamical systems
 521 are known as limit cycles (see e.g [41]). Limit cycles (or periodic orbits) are sets of states \mathcal{S} such that
 522 each state $x \in \mathcal{S}$ is visited at periodic intervals *ad infinitum* under the dynamics. Thus, if the gradient-
 523 based algorithms converge to a limit cycle they will cycle infinitely through the same sequence of
 524 actions. Like equilibria, limit cycles can be stable or unstable under the dynamics $\dot{x} = -\omega(x)$, mean-
 525 ing that the dynamics can either converge to or diverge from them depending on their initializations.

526 We remark that the existence of oscillatory behaviors and limit cycles has been observed in the
 527 dynamics of of gradient-based learning in various settings like the training of Generative Adversarial

528 Networks [11], and multiplicative weights in finite action games [25]. We simply emphasize that the
 529 existence of such limiting behaviors is due to the fact that the dynamics are no longer gradient flows.
 530 This fact also allows for other complex limiting behaviors like chaos³ to exist in the dynamics of
 531 gradient-based learning in games. We also show in Appendix B.1 that gradient-based learning avoids
 532 some limit cycles.

533 In Appendix B.1, we formalize the notion of a limit cycle and its stability in the stochastic setting.
 534 Using these concepts, we then provide an analogous theorem to Theorem 4.3 which states that competi-
 535 tive stochastic gradient-based learning converges to linearly unstable limit cycles—a parallel notion
 536 to strict saddle points but pertaining to more general limit sets—on a set of measure zero, provided
 537 that analogous assumptions to those in the statement of Theorem 4.3 hold. Providing such guaran-
 538 tees requires a bit more mathematical formalism, and as such we leave the details of these results to
 539 Appendix B.

540 In pursuit of a more general class of games with desirable convergence properties, in Appendix B.2
 541 we also introduce a generalization of potential games, namely Morse-Smale games, for which the
 542 combined gradient dynamics correspond to a Morse-Smale vector field [18,31]. In such games players
 543 are guaranteed to converge to only (linearly stable) cycles or equilibria. In such games, however,
 544 players may still converge to non-Nash equilibria and avoid a subset of the Nash equilibria.

545 **5. Saddle Point LNE in LQ Dynamic Games.** In this section, we present empirical results
 546 that show that a non-negligible subset of two-player LQ games have local Nash equilibria that are
 547 strict saddle points of the gradient dynamics. LQ games serve as good benchmarks for analyzing the
 548 limiting behavior of gradient-play in a non-trivial setting since they are known to admit global Nash
 549 equilibria that can be found by solving a coupled set of Riccati equations [2]. LQ games can also
 550 be cast as multi-agent reinforcement learning problems where each agent has a policy that is a linear
 551 function of the state and a quadratic reward function. Gradient-play in LQ games can therefore be
 552 seen as a form of policy gradient.

553 The empirical results we now present imply that, even in the relatively straightforward case of lin-
 554 ear dynamics, linear feedback policies, and quadratic costs, policy gradient multi-agent reinforcement
 555 learning would be unable to find the local Nash equilibrium in a non-negligible subset of problems.

556 *LQ game setup.* For simplicity, we consider two-player LQ games in \mathbb{R}^2 . Consider a discrete
 557 time dynamical system defined by

$$558 \quad (5.1) \quad z(t+1) = Az(t) + B_1u_1(t) + B_2u_2(t)$$

560 where $z(t) \in \mathbb{R}^2$ is the state at time t , $u_1(t)$ and $u_2(t)$ are the control inputs of players 1 and 2,
 561 respectively, and A , B_1 , and B_2 are the system matrices. We assume that player i searches for a linear
 562 feedback policy of the form $u_i(t) = -K_i z(t)$ that minimizes their loss which is given by

$$563 \quad f_i(z_0, u_1, u_2) = \sum_{t=0}^{\infty} z(t)^T Q_i z(t) + u_i(t)^T R_i u_i(t)$$

564 where $Q_i \succ 0$ and $R_i \succ 0$ are the cost matrices on the state and input, respectively. We note that the
 565 two players are coupled through the dynamics since $z(t)$ is constrained to obey the update equation

³A general term used to characterize dynamical systems where arbitrarily small perturbations in the initial conditions lead to drastically different solutions to the differential equations

566 (5.1). The vector of player derivatives is given by $\omega(K_1, K_2) = (D_1 f_1(K_1, K_2), D_2 f_2(K_1, K_2))$
 567 where

$$568 \quad D_i f_i(K_1, K_2) = (R_{ii} K_i + B_i^T P_i (B_1 K_1 + B_2 K_2) - B_i^T P_i A) \sum_{t=0}^{\infty} z(t) z(t)^T, \quad i \in \{1, 2\}.$$

569 Note that there is a slight abuse of notation here as we are treating $D_i f_i$ as a matrix and as the vector-
 570 ization of a matrix. The matrices P_1 and P_2 can be found by solving the Riccati equations

$$571 \quad P_i = (A - B_1 K_1 - B_2 K_2)^T P_i (A - B_1 K_1 - B_2 K_2) + K_i^T R_i K_i + Q_i, \quad i \in \{1, 2\},$$

573 for a given (K_1, K_2) . As shown in [2], global Nash equilibria of LQ games can be found by solving
 574 coupled Riccati equations. Under the following assumption, this can be done using an analogous
 575 method to the method of Lyapunov iterations outlined in [24] for continuous time LQ games.

576 **Assumption 4.** *Either $(A, B_1, \sqrt{Q_1})$ or $(A, B_2, \sqrt{Q_2})$ is stabilizable-detectable.*

577 Further information on the uniqueness of Nash equilibria in LQ games and the method of Lyapunov
 578 iterations can be found in [2] and [24] respectively.

579 **Generating LQ games with strict saddle point Nash equilibria.** Without loss of generality,
 580 we assume $(A, B_1, \sqrt{Q_1})$ is stabilizable-detectable. Given that we have a method of finding the global
 581 Nash equilibrium of the LQ game, we now present our experimental setup.

582 We fix B_1, B_2, Q_1 , and R_1 and parametrize Q_2 , and R_2 by q and r respectively. The shared
 583 dynamics matrix A has entries that are sampled from the uniform distribution supported on $(0, 1)$. For
 584 each value of the parameters b, q , and r , we randomly sample 1000 different A matrices. Then, for each
 585 LQ game defined in terms of each of the sets of parameters, we find the optimal feedback matrices
 586 (K_1^*, K_2^*) using the method of Lyapunov iterations, and we numerically approximate $D\omega(K_1^*, K_2^*)$
 587 using auto-differentiation tools and check its eigenvalues.

588 The exact values of the matrices are defined as follows: $A \in \mathbb{R}^{2 \times 2}$ with each of the entries a_{ij}
 589 sampled from the uniform distribution on $(0, 1)$,

$$590 \quad B_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad B_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad Q_1 = \begin{bmatrix} 0.01 & 0 \\ 0 & 1 \end{bmatrix}, \quad Q_2 = \begin{bmatrix} 1 & 0 \\ 0 & q \end{bmatrix}, \quad R_1 = 0.01, \quad R_2 = r.$$

592 The results for various combinations of the parameters q and r are shown in Figure 2. For all of
 593 the different parameter configurations considered, we found that in anywhere from 0% – 25% of the
 594 randomly sampled LQ games, there was a global Nash equilibrium that was a strict saddle point of
 595 the gradient dynamics. Of particular interest is the fact that for all values of q and r we tested, at least
 596 5% of the LQ games had a global Nash equilibrium with the strict saddle property. In the worst case,
 597 around 25% of the LQ games for the given values of q and r admitted such Nash equilibria.

598 **Remark 5.1.** These empirical observations imply that multi-agent policy gradient, even in the rel-
 599 atively straightforward setting of linear dynamics, linear policies, and quadratic costs, has no guar-
 600 antees of convergence to the global Nash equilibria in a non-negligible number of games. Further
 601 investigation is warranted to validate this fact theoretically. This in turn supports the idea that for more
 602 complicated cost functions, policy classes, and dynamics, local Nash equilibria with the strict saddle
 603 property are likely to be very common.

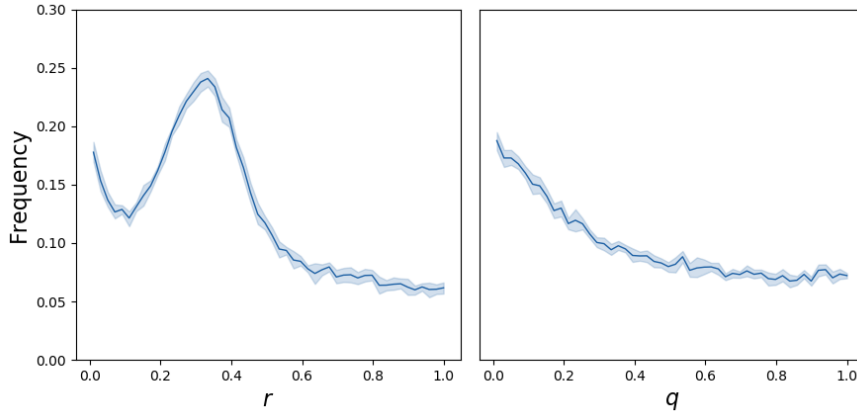


Figure 2: Frequency (out of 1000) of randomly sampled LQ games with global Nash equilibria that are avoided by policy-gradient. The experiment was run 10 times and the average frequency is shown by the solid line. The shaded region demarcates the 95% confidence interval of the experiment. (left) r is varied in $(0, 1)$, $q = 0.01$. (right) q is varied in $(0, 1)$, $r = 0.1$.

604 **6. Discussion and Future Directions.** In this paper we provided answers to the following
 605 two questions for classes of gradient-based learning algorithms:

606 **Q1.** *Are all attractors of the learning algorithms employed by agents equilibria relevant to the under-*
 607 *lying game?*

608 **Q2.** *Are all equilibria relevant to the game also attractors of the learning algorithms agents employ?*

609 We answered these questions in general-sum, zero-sum, and potential games without imposing
 610 structure on the game outside regularity conditions on the cost functions by exploiting the observation
 611 that gradient-based learning dynamics are not gradient flows. Our analysis, was shown in Section C to
 612 apply to a number of commonly used methods in multi-agent learning.

613 **6.1. Links with Prior Work.** As we noted, previous work on learning in games in both the
 614 game theory literature, and more recently from the machine learning community, has largely focused
 615 on **Q1**, though some recent work has analyzed **Q2** in the setting of zero-sum games.

616 In the seminal work by Rosen [40], n -player concave or monotone games are shown to either
 617 admit a unique Nash equilibrium or a continuum of Nash equilibria, all of which are attracting under
 618 gradient-play. The structure present in these games rules out the existence of non-Nash equilibria.

619 Two-player, finite-action bilinear games have also been extensively studied. In [44], the authors
 620 investigate the convergence of the gradient dynamics in such games. Additionally, the dynamics of
 621 other (non gradient-based) algorithms like multiplicative weights have been studied in [19] among
 622 many others. In such settings, the structure guarantees that there exists a unique global Nash equi-
 623 librium and no other critical points of the gradient dynamics. As such, non-Nash equilibria, cannot
 624 exist.

625 In the study of learning dynamics in the class of zero-sum games, it has been shown that cycles can
 626 be attractors of the dynamics (see, e.g., [19, 25, 47]). Concurrently with our results, [12] also showed
 627 the existence of non-Nash attracting equilibria in this setting.

628 In more general settings, there has been some analysis of the limiting behavior of gradient-play
 629 though the focus has been for the most part, on giving sufficient conditions under which Nash equilibria
 630 are attracting under gradient-play. For example, [35–37], introduced the notion of a differential Nash
 631 equilibrium which is characterized by first and second order conditions on the players’ individual cost
 632 functions and which we made extensive use of. Following this body of work, [27] also investigated
 633 the local convergence of gradient-play in continuous games. They showed that if a Nash equilibrium
 634 satisfies a property known as *variational stability*, the equilibrium is attracting under gradient play.
 635 In twice continuously differentiable games, this condition coincides exactly with the definition of
 636 stable differential Nash equilibria. Though these works analyze a general class of games, the focus of
 637 the analysis is solely on the local characterization and computation (via gradient play) of local Nash
 638 equilibria. As such, the issues of non-convergence that we show in this paper were not discussed.

639 **6.2. Open Questions.** Our results suggest that gradient-play in multi-agent settings has fun-
 640 damental problems. Depending on the players’ costs, in general games and even potential games,
 641 which have a particularly *nice* structure, a subset of the Nash equilibria will be almost surely avoided
 642 by gradient-based learning when the agents randomly initialize their first action. In zero-sum and
 643 general-sum games, even if the algorithms do converge, they may have converged to a point that has
 644 no game theoretic relevance, namely a non-Nash locally asymptotically stable equilibrium.

645 Lastly, these results show that limit cycles persist even under a stochastic update scheme. This
 646 explains the empirical observations of limit cycles in gradient dynamics presented in [11, 19, 23]. It
 647 also implies that gradient-based learning in multi-agent reinforcement learning, multi-armed bandits,
 648 generative adversarial networks, and online optimization all admit limit cycles under certain loss func-
 649 tions. Our empirical results show that these problems are not merely of theoretical interest, but also
 650 have great relevance in practice.

651 Which classes of games have all Nash being attracting for gradient-play and which classes pre-
 652 clude the existence of non-Nash equilibria is an open and particularly interesting question. Further, the
 653 question of whether gradient-based algorithms can be constructed for which only game-theoretically
 654 relevant equilibria are attracting is of particular importance as gradient-based learning is increasingly
 655 implemented in game theoretic settings. Indeed, more generally, as learning algorithms are increas-
 656 ingly deployed in markets and other competitive environments understanding and dealing with such
 657 theoretical issues will become increasingly important.

658 **Appendix A. Proofs of the Main Results.** This appendix contains the full proofs of the
 659 results in the paper.

660 **A.1. Proofs on Links Between Dynamical Systems and Games.** We begin with a proof
 661 of Proposition 3.4 that all differential Nash equilibria are either strict saddle points or asymptotically
 662 stable equilibria of the gradient dynamics. This relies mainly on the definitions of strict saddle points,
 663 locally asymptotically stable equilibria, and non-degenerate differential Nash equilibria and simple
 664 linear algebra.

665 *Proof of Proposition 3.4.* Suppose that $x \in X$ is a non-degenerate differential Nash equilibrium.
 666 We claim that $\text{tr}(D\omega(x)) > 0$. Since x is a differential Nash equilibrium, $D_i^2 f_i(x) \succ 0$ for each
 667 $i \in \mathcal{I}$; these are the diagonal blocks of $D\omega(x)$. Further $D_i^2 f_i(x) \succ 0$ implies that $\text{tr}(D_i^2 f_i(x)) > 0$.
 668 Since $\text{tr}(D\omega) = \sum_{i=1}^n \text{tr}(D_i^2 f_i(x))$, $\text{tr}(D\omega(x)) > 0$. Thus, it is not possible for all the eigenvalues
 669 to have negative real part. Since x is non-degenerate, $\det(D\omega(x)) \neq 0$ so that none of the eigenvalues

670 can have zero real part. Hence, at least one eigenvalue has strictly positive real part.

671 To complete the proof, we show that the conditions for non-degenerate differential Nash equilibri-
672 um are not sufficient to guarantee that x is locally asymptotically stable for the gradient dynamics—
673 that is, not all eigenvalues of $D\omega(x)$ have strictly positive real part. We do this by constructing a class
674 of games with the strict saddle point property. Consider a class of two player games $\mathcal{G} = (f_1, f_2)$ on
675 $\mathbb{R} \times \mathbb{R}$ defined as follows:

$$676 \quad (f_1(x_1, x_2), f_2(x_1, x_2)) = \left(\frac{a}{2}x_1^2 + bx_1x_2, \frac{d}{2}x_2^2 + cx_1x_2 \right).$$

677 In this game, the Jacobian of the gradient dynamics is given by

$$678 \quad (\text{A.1}) \quad D\omega(x) = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

679 with $a, b, c, d \in \mathbb{R}$. If x is a non-degenerate differential Nash equilibria, $a, d > 0$ and $\det(D\omega(x)) \neq 0$
680 which implies that $ad \neq cb$. Choosing c, d such that $ad < cb$ will guarantee that one of the eigenvalues
681 of $D\omega(x)$ is negative and the other is positive, making x a strict saddle point. This shows that non-
682 degenerate differential Nash equilibria can be strict saddle points of the combined gradient dynamics.

683 Hence, for any game (f_1, \dots, f_n) , a non-degenerate differential Nash equilibrium is either a lo-
684 cally asymptotically stable equilibrium or a strict saddle point, but it not strictly unstable or strictly
685 marginally stable (i.e. having eigenvalues all on the imaginary axis). ■

686 The proof of Proposition 3.7, which claims that all differential Nash equilibria in zero-sum games
687 are locally asymptotically stable, again just relies on basic linear algebra and the definition of a differ-
688 ential Nash equilibrium.

689 *Proof of Proposition 3.7.* Consider a two player game $(f, -f)$ on $X_1 \times X_2 = \mathbb{R}^m$ with $X_i = \mathbb{R}^{m_i}$.
690 For such a game,

$$691 \quad D\omega(x) = \begin{bmatrix} D_1^2 f(x) & D_{21} f(x) \\ -D_{12} f(x) & -D_2^2 f(x) \end{bmatrix}.$$

692 Note that $D_{21} f(x) = (D_{12} f(x))^T$. Suppose that $x = (x_1, x_2)$ is a differential Nash equilibrium
693 and let $v = [v_1, v_2] \in \mathbb{R}^m$ with $v_1 \in \mathbb{R}^{m_1}$ and $v_2 \in \mathbb{R}^{m_2}$. Then, $v^T D\omega(x) v = v_1^T D_1^2 f(x) v_1 -$
694 $v_2^T D_2^2 f(x) v_2 > 0$ since $D_1^2 f(x) \succ 0$ and $-D_2^2 f(x) \succ 0$ for x , a differential Nash equilibrium. Since
695 v is arbitrary, this implies that $D\omega(x)$ is positive definite and hence, clearly non-degenerate. Thus, for
696 two-player zero-sum games, all differential Nash equilibria are both non-degenerate differential Nash
697 equilibria and locally asymptotically stable equilibria of $\dot{x} = -\omega(x)$. ■

698 The proof that all locally asymptotically stable equilibria in potential games are differential Nash
699 equilibria relies on the symmetry of $D\omega$ in potential games.

700 *Proof of Proposition 3.9.* The proof follows from the definition of a potential game. Since $(f_1,$
701 $\dots, f_n)$ is a potential game, it admits a potential function ϕ such that $D_i f_i(x) = D_i \phi(x)$ for all x .
702 This, in turn, implies that at a locally asymptotically stable equilibrium of $\dot{x} = -\omega(x)$, $D\omega(x) =$
703 $D^2 \phi(x)$, where $D^2 \phi$ is the Hessian matrix of the function ϕ . Further $D^2 \phi(x)$ must have strictly
704 positive eigenvalues for x to be a locally asymptotically stable equilibrium of $\dot{x} = -\omega(x)$. Since the
705 Hessian matrix of a function must be symmetric, $D^2 \phi(x)$, must be positive definite, which through
706 Sylvester's criterion ensures that each of the diagonal blocks of $D^2 \phi(x)$ is positive definite. Thus, we

707 have that the existence of a potential function guarantees that the only locally asymptotically stable
708 equilibria of $\dot{x} = -\omega(x)$, are differential Nash equilibria. ■

709 **A.2. Proofs for Deterministic Setting.** We now present the proof of Theorem 4.1 and its
710 corollaries. The proof relies on the celebrated stable manifold theorem [43, Theorem III.7], [45].
711 Given a map ϕ , we use the notation $\phi^t = \phi \circ \dots \circ \phi$ to denote the t -times composition of ϕ .

712 **Theorem A.1 (Center and Stable Manifolds [43, Theorem III.7], [45]).** *Let x_0 be a fixed point
713 for the C^r local diffeomorphism $f : U \rightarrow \mathbb{R}^d$ where $U \subset \mathbb{R}^d$ is an open neighborhood of x_0 in \mathbb{R}^d
714 and $r \geq 1$. Let $E^s \oplus E^c \oplus E^u$ be the invariant splitting of \mathbb{R}^d into generalized eigenspaces of $D\phi(x_0)$
715 corresponding to eigenvalues of absolute value less than one, equal to one, and greater than one. To
716 the $D\phi(x_0)$ invariant subspace $E^s \oplus E^c$ there is an associated local ϕ -invariant C^r embedded disc
717 W_{loc}^{cs} called the local stable center manifold of dimension $\dim(E^s \oplus E^c)$ and ball B around x_0 such
718 that $\phi(W_{loc}^{cs}) \cap B \subset W_{loc}^{cs}$, and if $\phi^t(x) \in B$ for all $t \geq 0$, then $x \in W_{loc}^{sc}$.*

719 Some parts of the proof follow similar arguments to the proofs of results in [22, 32] which apply to
720 (single-agent) gradient-based optimization. Due to the different learning rates employed by the agents
721 and the introduction of the differential game form ω , the proof differs.

722 *Proof of Theorem 4.1.* The proof is composed of two parts: (a) the map g is a diffeomorphism, and
723 (b) application of the stable manifold theorem to conclude that the set of initial conditions is measure
724 zero.

725 *(a) g is diffeomorphism.* We claim the mapping $g : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is a diffeomorphism. If we can
726 show that g is invertible and a local diffeomorphism, then the claim follows. Consider $x \neq y$ and
727 suppose $g(y) = g(x)$ so that $y - x = \gamma \cdot (\omega(y) - \omega(x))$. The assumption $\sup_{x \in \mathbb{R}^m} \|D\omega(x)\|_2 \leq$
728 $L < \infty$ implies that ω satisfies the Lipschitz condition on \mathbb{R}^m . Hence, $\|\omega(y) - \omega(x)\|_2 \leq L\|y - x\|_2$.
729 Let $\Gamma = \text{diag}(\Gamma_1, \dots, \Gamma_n)$ where $\Gamma_i = \text{diag}((\gamma_i)_{j=1}^{m_i})$ —that is, Γ_i is an $m_i \times m_i$ diagonal matrix
730 with γ_i repeated on the diagonal m_i times. Then, $\|x - y\|_2 \leq L\|\Gamma\|_2\|y - x\|_2 < \|y - x\|_2$ since
731 $\|\Gamma\|_2 = \max_i |\gamma_i| < 1/L$.

732 Now, observe that $Dg = I - \Gamma D\omega(x)$. If Dg is invertible, then the implicit function theorem [21,
733 Theorem C.40] implies that g is a local diffeomorphism. Hence, it suffices to show that $\Gamma D\omega(x)$
734 does not have an eigenvalue of 1. Indeed, letting $\rho(A)$ be the spectral radius of a matrix A , we
735 know in general that $\rho(A) \leq \|A\|$ for any square matrix A and induced operator norm $\|\cdot\|$ so that
736 $\rho(\Gamma D\omega(x)) \leq \|\Gamma D\omega(x)\|_2 \leq \|\Gamma\|_2 \sup_{x \in \mathbb{R}^m} \|D\omega(x)\|_2 < \max_i |\gamma_i| L < 1$. Of course, the spectral
737 radius is the maximum absolute value of the eigenvalues, so that the above implies that all eigenvalues
738 of $\Gamma D\omega(x)$ have absolute value less than 1.

739 Since g is injective by the preceding argument, its inverse is well-defined and since g is a local
740 diffeomorphism on \mathbb{R}^m , it follows that g^{-1} is smooth on \mathbb{R}^m . Thus, g is a diffeomorphism.

741 *(b) Application of the stable manifold theorem.* Consider all critical points to the game—
742 i.e. $\mathcal{X}_c = \{x \in X \mid \omega(x) = 0\}$. For each $p \in \mathcal{X}_c$, let B_p be the open ball derived from Theorem A.1
743 and let $\mathcal{B} = \cup_p B_p$. Since $X \subseteq \mathbb{R}^m$, Lindelöf’s lemma [20]—every open cover has a countable
744 subcover—gives a countable subcover of \mathcal{B} . That is, for a countable set of critical points $\{p_i\}_{i=1}^\infty$ with
745 $p_i \in \mathcal{X}_c$, we have that $\mathcal{B} = \cup_{i=1}^\infty B_{p_i}$.

746 Starting from some point $x_0 \in X$, if gradient-based learning converges to a strict saddle point, then
747 there exists a t_0 and index i such that $g^t(x_0) \in B_{p_i}$ for all $t \geq t_0$. Again, applying Theorem A.1 and
748 using that $g(X) \subset X$ —which we note is obviously true if $X = \mathbb{R}^m$ —we get that $g^t(x_0) \in W_{loc}^{cs} \cap X$.

749 Using the fact that g is invertible, we can iteratively construct the sequence of sets defined by
 750 $W_1(p_i) = g^{-1}(W_{\text{loc}}^{cs} \cap X)$ and $W_{k+1}(p_i) = g^{-1}(W_k(p_i) \cap X)$. Then we have that $x_0 \in W_t(p_i)$ for
 751 all $t \geq t_0$. The set $\mathcal{X}_0 = \cup_{i=1}^{\infty} \cup_{t=0}^{\infty} W_t(p_i)$ contains all the initial points in X such that gradient-based
 752 learning converges to a strict saddle. Since p_i is a strict saddle, $I - \Gamma D\omega(p_i)$ has an eigenvalue greater
 753 than 1. This implies that the co-dimension of E^u is strictly less than m . (i.e. $\dim(W_{\text{loc}}^{cs}) < m$). Hence,
 754 $W_{\text{loc}}^{cs} \cap X$ has Lebesgue measure zero in \mathbb{R}^m .

755 Using again that g is a diffeomorphism, $g^{-1} \in C^1$ so that it is locally Lipschitz and locally
 756 Lipschitz maps are null set preserving. Hence, $W_k(p_i)$ has measure zero for all k by induction so that
 757 \mathcal{X}_0 is a measure zero set since it is a countable union of measure zero sets. ■

758 The proof of Corollary 4.2 follows from the symmetry of $D\omega$ in potential games, and our obser-
 759 vations in Section 3.

760 *Proof of Corollary 4.2.* Since the game admits a potential function ϕ , there is a transformation
 761 of coordinates such that agents following the dynamics $x_{t+1} = x_t - \gamma \odot \omega(x_t)$ converge to the
 762 same equilibria as the gradient dynamics $x_{t+1} = x_t - \gamma \odot D\phi(x_t)$. Hence, the analysis of the
 763 gradient-based learning scheme reduces to analyzing gradient-based optimization of ϕ . Moreover,
 764 existence of a potential function also implies that $D_{ij}f_j \equiv D_{ji}f_i$ so that $D\omega$ is symmetric. Indeed,
 765 writing $\omega(x)$ as the differential form $\sum_{i=1}^n D_i f_i(x) dx_i$ and noting that $d \circ d = 0$ for the differential
 766 operator d , we have that $d(\omega) = \sum_i d(D_i f_i) \wedge dx_i = \sum_{i,j:j>i} (D_{ij}f_j - D_{ji}f_i) dx_i \wedge dx_j = 0$
 767 where \wedge is the standard exterior product [21]. Symmetry of $D\omega$ implies that all periodic orbits are
 768 equilibria—i.e. the dynamics do not possess any limit cycles. By Theorem 4.1, the set of initial points
 769 that converge to strict saddle points is of measure zero. Since all the stable critical points of the
 770 dynamics are equilibria, with the assumption that $\lim_{t \rightarrow \infty} g^t(x)$ exists for all $x \in X$, we have that
 771 $P_\nu [\lim_{t \rightarrow \infty} g^t(x) = x^*] = 1$ where x^* is a non-degenerate differential Nash equilibrium which is
 772 generically a local Nash equilibrium [36].

773 **A.3. Classical Results from Dynamical Systems.** The remaining results use the following
 774 classical result from dynamical systems theory. Consider a general stochastic approximation frame-
 775 work $x_{t+1} = x_t + \gamma_t(h(x_t)) + \epsilon_t$ for $h : X \rightarrow TX$ with $h \in C^2$ and where $X \subset \mathbb{R}^d$ and where TX
 776 denotes the tangent space.

777 **Theorem A.2 (Theorem 1 [33]).** *Suppose γ_t is \mathcal{F}_t -measurable and $\mathbb{E}[w_t | \mathcal{F}_t] = 0$. Let the*
 778 *stochastic process $\{x_t\}_{t \geq 0}$ be defined as above for some sequence of random variables $\{\epsilon_t\}$ and $\{\gamma_t\}$.*
 779 *Let $p \in X$ with $h(p) = 0$ and let W be a neighborhood of p . Assume that there are constants*
 780 *$\eta \in (1/2, 1]$ and $c_1, c_2, c_3, c_4 > 0$ for which the following conditions are satisfied whenever $x_t \in$*
 781 *W and t sufficiently large: (i) p is a linear unstable critical point, (ii) $c_1/t^\eta \leq \gamma_t \leq c_2/t^\eta$, (iii)*
 782 *$\mathbb{E}[(w_t \cdot v)^+ | \mathcal{F}_t] \geq c_3/t^\eta$ for every unit vector $v \in TX$, and (iv) $\|w_t\|_2 \leq c_4/t^\eta$. Then $P(x_t \rightarrow p) = 0$.*
 783

784 **Appendix B. Expanded Results in the Stochastic Setting.** In this appendix, we provide
 785 extended results in the stochastic setting that require more mathematical formalism than the main body
 786 of the paper. In addition, we introduce a new class of games that generalize potential games and have
 787 stronger convergence guarantees than the broader class of general-sum continuous games.

788 **B.1. Avoidance of Repelling Sets.** To show that stochastic gradient-based learning avoids
 789 of more general limiting behaviors than saddle points, we need further assumptions on our underlying

790 space—i.e. we need the underlying decision spaces of each agent—i.e. X_i for each $i \in \mathcal{I}$ —to be
 791 *smooth, compact manifolds without boundary*⁴. The stochastic process $\{x_n\}$ which follows (4.2) is
 792 *defined on X* —that is, $x_n \in X$ for all $n \geq 0$. As before, it is natural to compare sample points $\{x_n\}$ to
 793 solutions of $\dot{x} = -\omega(x)$ where we think of (4.2) as a noisy approximation. The asymptotic behavior
 794 of $\{x_n\}$ can indeed be described by the asymptotic behavior of the flow generated by ω .

795 We also need a formal notion of *cycles*. A non-stationary periodic orbit of ω is called a *cycle*.
 796 Let $\xi \subset X$ be a cycle of period $T > 0$. Denote by Φ_T the flow corresponding to ω . For any
 797 $x \in \xi$, $\text{spec}(D\Phi_T(x)) = \{1\} \cup C(\xi)$ where $C(\xi)$ is the set of characteristic multipliers. We say ξ is
 798 *hyperbolic* if no element of $C(\xi)$ is on the complex unit circle. Further, if $C(\xi)$ is strictly inside the unit
 799 circle, ξ is called *linearly stable* and, on the other hand, if $C(\xi)$ has at least one element on the outside
 800 of the unit circle—that is, $D\Phi_T(x)$ for $x \in \xi$ has an eigenvalue with real part strictly greater than
 801 1—then ξ is called *linearly unstable*. The latter is the analog of strict saddle points in the context of
 802 periodic orbits. We denote by $\{x_t\}$ sample paths of the process (4.2) and $L(\{x_t\})$ is the *limit set* of any
 803 sequence $\{x_t\}_{t \geq 0}$ which is defined in the usual way as all $p \in X$ such that $\lim_{k \rightarrow \infty} x_{t_k} = p$ for some
 804 sequence $t_k \rightarrow \infty$. It was shown in [3] that under less restrictive assumptions than Assumptions 2
 805 and 3, $L(\{x_t\})$ is contained in the *chain recurrent set* of ω and $L(\{x_t\})$ is a non-empty, compact and
 806 connected set invariant under the flow of ω .

807 **Theorem B.1.** *Consider a game (f_1, \dots, f_n) where each X_i is a smooth, compact manifold with-*
 808 *out boundary. Suppose each agent $i \in \mathcal{I}$ adopts a stochastic gradient-based learning algorithm that*
 809 *satisfies Assumptions 2 and 3 and is such that sample points $x_t \in X$ for all $t \geq 0$. Further, suppose*
 810 *that for each $i \in \mathcal{I}$, there exist a constant $b_i > 0$ such that $\mathbb{E}[(w_{i,t} \cdot v)^+ | \mathcal{F}_{i,t}] \geq b_i$ for every unit vector*
 811 *$v \in \mathbb{R}^{m_i}$. Then competitive stochastic gradient-based learning converges to linearly unstable cycles*
 812 *on a set of measure zero—i.e. $P(L(x_t) = \xi) = 0$ where $\{x_t\}$ is a sample path.*

813 As we noted, periodic orbits are not necessarily excluded from the limiting behavior of gradient-based
 814 learning in games. We leave out the proof of Theorem B.1 since after some algebraic manipulation, it
 815 is a direct application of [4, Theorem 2.1] which is stated below.

816 **Theorem B.2 (Theorem 2.1 [4]).** *Let $\xi \subset X$ be a hyperbolic linearly unstable cycle of h .*
 817 *Assume the following (i) $h \in C^2$; (ii) $c_1/t^\eta \leq \gamma_t \leq c_2/t^\eta$ with $0 < c_1 \leq c_2$ and $0 < \eta \leq 1$;*
 818 *and (iii) there exists $b \geq 0$ such that for all unit vectors $v \in \mathbb{R}^m$, $\mathbb{E}[(w_t \cdot v)^+ | \mathcal{F}_t] \geq b$. Then*
 819 *$P(L(\{x_t\}) = \xi) = 0$.*

820 **B.2. Morse-Smale Games.** For a class of games admitting *gradient-like* vector fields we can
 821 go beyond non-convergence results and give convergence guarantees. Following [4], we introduce a
 822 new class of games, which we call *Morse-Smale games*, that are a generalization of potential games.
 823 Such games represent an important class since the vector field of ω corresponds to Morse-Smale vector
 824 field which is known to be generic in \mathbb{R}^2 and are otherwise structurally stable [18, 31].

825 **Definition B.3.** *A game (f_1, \dots, f_n) with $f_i \in C^r$ for some $r \geq 3$ and where strategy spaces X_i*
 826 *is a smooth, compact manifold without boundary for each $i \in \mathcal{I}$ is a Morse-Smale game if the vector*
 827 *field corresponding to the differential ω is Morse-Smale—that is, the following hold: (i) all periodic*
 828 *orbits ξ (i.e. equilibria and cycles) are hyperbolic and $W^s(\xi) \cap W^u(\xi) = \emptyset$ (i.e. the stable and unstable*

⁴The torus $\mathbb{T} = \mathbb{S}^1 \times \mathbb{S}^1$ is an example. The interested reader can consult, e.g., [21] for more details on differential geometry.

829 manifolds of ξ intersect transversally), (ii) every forward and backward omega limit set is a periodic
830 orbit, (iii) and ω has a global attractor.

831 The conditions of Morse-Smale in the above definition ensure that there are only finitely many periodic
832 orbits. The dynamics of games with more general vector fields, on the other hand, can admit chaos (e.g.
833 the classic Lorentz attractor can be cast as gradient-play in a 3-player game). Hyperbolic equilibria and
834 periodic orbits are the only types of limiting behavior that have been shown to correspond to strategies
835 relevant to the underlying game [5]. The simplest example of a Morse-Smale vector field is a gradient
836 flow. However, not all Morse-Smale vector fields are gradient flows and hence, not all Morse-Smale
837 games are potential games.

838 **Example 1.** Consider the n -player game with $X_i = \mathbb{R}$ for each $i \in \mathcal{I}$ and $f_n(x) = x_n(x_1^2 -$
839 $1)$, $f_i(x) = x_i x_{i+1}$, $\forall i \in \mathcal{I}/\{n\}$. This is a Morse-Smale game that is not a potential game. Indeed,
840 $\dot{x} = -\omega(x)$ where $\omega = [x_2, x_3, \dots, x_{n-1}, x_1^2 - 1]$ is a dynamical system with a Morse-Smale vector
841 field that is not a gradient vector field [10].

842 Essentially, in a neighborhood of a critical point for a Morse-Smale game, the game behavior can
843 be described by a Morse function ϕ such that near critical points ω can be written as $D\phi$ and away
844 from critical points ω points in the same direction as $D\phi$ —i.e. $\omega \cdot D\phi > 0$. Specializing the class of
845 Morse-Smale games, we have stronger convergence guarantees.

846 **Theorem B.4.** Consider a Morse-Smale game (f_1, \dots, f_n) on smooth boundaryless compact
847 manifold X . Suppose Assumptions 2 and 3 hold and that $\{x_t\}$ is defined on X . Let $\{\xi_i, i = 1, \dots, l\}$
848 denote the set of periodic orbits in X . Then $\sum_{i=1}^l P(L(\{x_t\}) = \xi_i) = 1$ and $P(L(\{x_t\}) = \xi_i) > 0$
849 implies ξ_i is linearly stable. Moreover, if the periodic orbit ξ_i with $P(L(\{x_t\}) = \xi_i) > 0$ is an equi-
850 librium, then it is either a non-degenerate differential Nash equilibrium—which is generically a local
851 Nash—or a non-Nash locally asymptotically stable equilibrium.

852 The proof of Theorem B.4 follows by invoking Corollary B.5 which is stated below.

853 **Corollary B.5 (Corollary 2.2 [4]).** Assume that there exists $\delta \geq 1$ such that $\sum_{n \geq 0} \gamma_n^{1+\delta} < \infty$
854 and that h is a Morse-Smale vector field. If we denote by $\{\xi_i, i = 1, \dots, l\}$ the set of periodic orbits
855 in X , then $\sum_{i=1}^l P(L(\{x_t\}) = \xi_i) = 1$. Further, if conditions (i)–(iii) of Theorem B.2 hold, then
856 $P(L(\{x_t\}) = \xi_i) > 0$ implies ξ_i is linearly stable.

857 Thus, in Morse-Smale games, with probability one, the limit sets of competitive gradient-based
858 learning with stochastic updates are attractors (i.e., periodic orbits, which includes limit cycles and
859 equilibria) of $\dot{x} = -\omega(x)$ and if any attractor has positive probability of being a limit set of the
860 players' collective update rule, then it is (linearly) stable. Moreover, attractors that are equilibria
861 are either non-degenerate differential Nash equilibria (generically local Nash equilibria) or non-Nash
862 locally asymptotically stable equilibria, but not saddle points.

863 If we further restrict the class of games to potential games, the results for Morse-Smale games
864 imply convergence to Nash almost surely, a particularly strong convergence guarantee.

865 **Corollary B.6.** Consider the game (f_1, \dots, f_n) on smooth boundaryless compact manifold $X =$
866 $X_1 \times \dots \times X_n$ admitting potential function ϕ . Suppose each agent $i \in \mathcal{I}$ adopts a stochastic gradient-
867 based learning algorithm that satisfies Assumptions 2 and 3 and such that $\{x_t\}$ evolves on X . Further,
868 suppose that for each $i \in \mathcal{I}$, there exist a constant $b_i > 0$ such that $\mathbb{E}[(w_{i,t} \cdot v)^+ | \mathcal{F}_{i,t}] \geq b_i$ for every unit
869 vector $v \in \mathbb{R}^{m_i}$. Then, competitive stochastic gradient-based learning converges to a non-degenerate

870 *differential Nash equilibrium almost surely.*

871 The proof of Corollary B.6 follows from the fact that potential games are trivially Morse-Smale
872 games that admit no periodic cycles as we showed in the proof of Corollary 4.2.

873 *Proof of Corollary B.6.* Consider a potential game (f_1, \dots, f_n) where each X_i is a smooth, com-
874 pact boundaryless manifold. Then $\omega = D\phi$ for some $\phi \in C^r$ which implies that ω is a gradient flow
875 and hence, does not admit limit cycles. Let $\{\xi_i, i = 1, \dots, l\}$ be the set of equilibrium points in X .
876 Under the assumptions of Theorem B.4, $\sum_{i=1}^l P(L(\{x_t\}) = \xi_i) = 1$ and, if $P(L(\{x_t\}) = \xi_i) > 0$,
877 then ξ_i is a linearly stable equilibrium point which is a non-degenerate differential Nash equilibrium
878 of the game due to the fact that $D\omega(x)$ is symmetric in potential games. Hence, a sample path $\{x_t\}$
879 converges to a non-degenerate differential Nash equilibrium with probability one. Moreover, by [36],
880 we know it is generically a local Nash. ■

881 We note, that even though a potential function is enough to guarantee convergence to a local Nash
882 equilibrium, potential games can still admit local Nash equilibria that are strict saddle points as shown
883 in Section 3. Thus, even this relatively well-behaved class of games has fundamental problems when
884 applying a gradient-based learning scheme.

885 **Appendix C. Classes of Gradient-Based Learning Algorithms.** In this section, we
886 provide derivation of the gradient-based learning rules provided in Table 1. We note that the deriva-
887 tion of gradient-based approaches for multi-armed bandits can be found in [46] among other classic
888 references on reinforcement learning.

889 **C.1. Online Optimization: Gradient Play in Non-Cooperative Games.** We first show
890 that classical online optimization algorithms fit into the framework we describe. In this case, each
891 agent is directly trying to minimize their own function $f_i(x_i, x_{-i})$, which can depend on the current
892 iterate of the other agents. There are many examples in the optimization literature of this type of
893 setup. We note that in the full information case, the competitive gradient-based learning framework
894 we describe here is simply *gradient play* [16], a very well-studied game-theoretic learning rule.

895 Of more interest are some gradient-free online optimization algorithms that also fit into the frame-
896 work we describe. The game can be described as follows. At each iteration, t of the game, every
897 player publishes their current iterate $x_{i,t}$. Player i , implementing this algorithm, then updates their
898 iterate by taking a random unit vector u , and querying $f_i(x_i + \delta_i u, x_{-i})$. The update map is given
899 by $x_{i,t+1} = x_{i,t} - \gamma_i f_i(x_i + \delta_i u, x_{-i})u$. It is shown in [14] that $f_i(x_i + \delta_i u, x_{-i})u$ is an unbiased
900 estimate of the gradient of a smoothed version of f_i —i.e. $\hat{f}_i(x_i, x_{-i}) = \mathbb{E}_v[f_i(x + \delta v, x_{-i})]$. Thus the
901 loss function being minimized by the agent is \hat{f}_i . In this case, the results on characterizing limiting
902 behavior presented in Section 4.2 apply.

903 **C.2. Generative Adversarial Networks.** Generative adversarial networks take a game theo-
904 retic approach to fitting a generative model in complex structured spaces. Specifically, they approach
905 the problem of fitting a generative model from a data set of samples from some distribution $Q \in \Delta(Y)$
906 as a zero-sum game between a *generator* and a *discriminator*. In general, both the generator and
907 the discriminator are modeled as deep neural networks. The generator network outputs a sample
908 $G_\theta(z) \in Y$ in the same space Y as the sampled data set given a random noise signal $z \sim F$ as an
909 input. The discriminator $D_w(y)$ tries to discriminate between a true sample and a sample generated
910 by the generator—that is, it takes as input a sample y drawn from Q or the generator and tries to de-

911 termine if its *real* or *fake*. The goal, is to find a Nash equilibrium of the zero-sum game under which
 912 the generator will learn to generate samples that are indistinguishable from the true samples—i.e. in
 913 equilibrium, the generator has learned the underlying distribution.

914 To prevent instabilities in the training of GANs with zero-one discriminators, the Wasserstein
 915 GAN attempts to approximate the Wasserstein-1 metric between the true distribution and the distribu-
 916 tion of the generator. In this setting, $D_w(\cdot)$ is a 1-Lipschitz function leading to the problem

$$917 \quad \inf_{\theta} \sup_w \mathbb{E}_{y \sim Q} [D_w(y)] - \mathbb{E}_{z \sim F} [D_w(G_{\theta}(z))]$$

919 which has corresponding dynamics $w_{t+1} = w_t + \gamma \nabla_w L(\theta_t, w_t)$ and $\theta_{t+1} = \theta_t - \gamma \nabla_{\theta} L(\theta_t, w_t)$ where
 920 $L(\theta, w) = \mathbb{E}_{y \sim Q} [D_w(y)] - \mathbb{E}_{z \sim F} [D_w(G_{\theta}(z))]$ and where γ is the learning rate.

921 GANs are notoriously difficult to train. The typical approach is to allow each player to perform
 922 (stochastic) gradient descent on the derivative of their cost with respect to their own choice variable.
 923 There are two important observations about gradient-based learning approaches to GANs relevant to
 924 this paper. First, the equilibrium that is sought is generally a saddle point and second, the dynamics
 925 of GANs are complex enough to admit limit cycles [25]. None-the-less, training GANs with gradient
 926 descent is still very common. We note that our results suggest that, on top of periodic orbits and
 927 oscillations, training GANs with gradient descent can result in convergence to non-Nash equilibria.

928 **C.3. Multi-Agent Reinforcement Learning Algorithms.** Consider a setting in which all
 929 agents are operating in an MDP. There is a shared state space \mathcal{S} . Each agent, indexed by $\mathcal{I} =$
 930 $\{1, \dots, n\}$ has their own action space U_i and reward function $R_i : \mathcal{S} \times U \rightarrow \Delta_{\mathbb{R}}$ where $U =$
 931 $U_1 \times \dots \times U_n$. We note the reward functions could themselves be random, but for illustrative pur-
 932 poses we suppose they are deterministic. Finally, the dynamics of the MDP are described by a state
 933 transition kernel $P : \mathcal{S} \times U \rightarrow \Delta_{\mathcal{S}}$ and an initial state distribution P_0 . Each agent i also has a policy,
 934 π_i , that returns a distribution over U_i for each state $s \in \mathcal{S}$. We define a trajectory of the MDP, τ as
 935 $\tau = \{(s_t, u_{i,t}, u_{-i,t})\}_{t=0}^{T-1}$. Thus, a trajectory is a finite sequence of states, the actions of each player
 936 in that state, and the reward agent i received in that state, where T is the time horizon. Given fixed
 937 policies we can define a distribution over the space of all trajectories Γ , namely $P_{\Gamma}(\pi)$, by

$$938 \quad P_{\Gamma}(\tau; \pi) = P_0(s_0) \prod_{i \in \mathcal{I}} \pi_i(u_{i,0} | s_0) \cdots P(s_t | s_{t-1}, u_{t-1}) \prod_{i \in \mathcal{I}} \pi_i(u_{i,t} | s_t) \cdots$$

940 The goal of each single agent in this setup is to maximize their cumulative expected reward over a
 941 time horizon T . That is, the agent is trying to find a policy π_i so as to maximize some function,
 942 which in keeping with our general formulation in Section 2, we write as $-f_i$ since this problem is
 943 a maximization. When an agent is employing policy gradient in this MARL setup, we assume that
 944 their policy comes from a parametric class of policies parametrized by $x_i \in X_i \subset \mathbb{R}^{m_i}$. To simplify
 945 notation, we write the parametric policy as $\pi_i(x_i)$ where for each x_i , given an state s , $\pi_i(x_i)$ is a
 946 probability distribution on actions u_i which we denote by $\pi_i(x_i)(\cdot | s)$.

947 The policy gradient MARL algorithm can be reformulated in the competitive gradient-based learn-
 948 ing framework. An agent i using policy gradient is trying to tune the parameters x_i of their policy to
 949 maximize their expected reward over a trajectory of length T . We define the reward of agent i over a
 950 trajectory of the MDP, $\tau \in \Gamma$, to be $\mathbf{R}_i(\tau) = \sum_{t=0}^{T-1} R_i(s_t, u_{i,t}, u_{-i,t})$. Thus, each agent's loss function
 951 f_i , in keeping with our notation, is given by $f_i(x_i, x_{-i}) = -J_i(\pi_i(x_i), \pi_{-i}) = -\mathbb{E}_{\tau \sim P_{\Gamma}(\pi)} [\mathbf{R}_i(\tau)]$.
 952 The actions of agent i in the continuous game framework described in previous sections are the pa-
 953 rameters of their policy, and thus their action space is $X_i \subset \mathbb{R}^{m_i}$. We note that we have made no

954 assumptions on the other player’s actions x_{-i} . That is, they do not need to be employing the same
 955 parameterized policy class or exactly the same gradient-based update procedure; the only requirement
 956 is that they also be using a gradient based multi-agent learning algorithm, and that their actions give
 957 rise to a set of policies π_{-i} that govern the way they choose their actions in the MDP.

958 In the full information case, at each round, t of the game, a player plays according to $\pi_i(x_{i,t})$ for
 959 a time horizon T , and then performs a gradient update on their parameters where $D_i J_i(\pi_i(x_i), \pi_{-i}) =$
 960 $D_i J_i(\pi_i(x_i), \pi_{-i,t})$ is given by

$$961 \quad (C.1) \quad D_i J_i(\pi_i(x_i), \pi_{-i}) = \mathbb{E}_{\tau \sim P_{\Gamma}(\pi)} \left[\sum_{t=0}^{T-1} R_i(s_t, u_t) \sum_{j=0}^t \nabla_{x_i} \log \pi_i(x_i)(u_{i,j} | s_j) \right]$$

963 The derivation of this gradient is exactly the same as that of classic policy gradient. From (C.1) it is
 964 clear that an unbiased estimate of the gradient can be constructed. At each time t in the policy gradient
 965 update procedure, agent i receives a T horizon roll-out, say $z_{i,t} = \{(s_k, u_{i,k}, r_{i,k})\}_{k=0}^{T-1}$, and constructs
 966 the unbiased estimate of the gradient—i.e. $\widehat{D_i J_i} = \sum_{k=0}^{T-1} r_{i,k} (\sum_{j=0}^k \nabla_{x_i} \log \pi_i(x_{i,t})(u_{i,j} | s_j))$. We
 967 note that in this case, the agent does not need to know the policies of the other agents, or anything
 968 about the dynamics of the MDP. The agent can construct the estimator solely from the sequence of
 969 states, the reward they received in those states, and their own actions. With these two derivations of
 970 the gradient for the full information and gradient-free cases, policy gradient for MARL conforms to
 971 the competitive gradient-based learning framework and hence, the results of Section 4 apply under
 972 appropriate assumptions.

973 REFERENCES

- 974 [1] S. ABDALLAH AND V. LESSER, *A multiagent reinforcement learning algorithm with non-linear dynamics*, Journal
 975 of Artificial Intelligence Research, (2008).
- 976 [2] T. BASAR AND G. OLSDER, *Dynamic Noncooperative Game Theory*, Society for Industrial and Applied Mathemat-
 977 ics, 1998.
- 978 [3] M. BENAÏM, *A dynamical system approach to stochastic approximations*, SIAM Journal on Control and Optimization,
 979 (1996).
- 980 [4] M. BENAÏM AND M. HIRSCH, *Dynamics of morse-smale urn processes*, Ergodic Theory and Dynamical Systems, 15
 981 (1995).
- 982 [5] M. BENAÏM AND M. W. HIRSCH, *Learning processes, mixed equilibria and dynamical systems arising from repeated*
 983 *games*, Games and Economic Behavior, (1997).
- 984 [6] V. BORKAR, *Stochastic Approximation: A Dynamical Systems Viewpoint*, 2008.
- 985 [7] L. BOTTOU, *Large-scale machine learning with stochastic gradient descent*, Proceedings in Computational Statistics,
 986 (2010).
- 987 [8] M. BRAVO, D. LESLIE, AND P. MERTIKOPOULOS, *Bandit learning in concave n-person games*, in Advances in
 988 Neural Information Processing Systems, 2018.
- 989 [9] A. S. CHIVUKULA AND W. LIU, *Adversarial learning games with deep learning models*, International Joint Confer-
 990 ence on Neural Networks, (2017).
- 991 [10] C. CONLEY, *Isolated invariant sets and the morse index*, in CBMS Regional Conference Series in Mathematics,
 992 1978.
- 993 [11] C. DASKALAKIS, A. ILYAS, V. SYRGKANIS, AND H. ZENG, *Traning GANs with Optimism*, Arxiv, (2017).
- 994 [12] C. DASKALAKIS AND I. PANAGEAS, *The limit points of (optimistic) gradient descent in min-max optimization*, in
 995 Proceedings of the 32Nd International Conference on Neural Information Processing Systems, 2018.
- 996 [13] M. FAZEL, R. GE, S. M. KAKADE, AND M. MESBAHI, *Global convergence of policy gradient methods for the*
 997 *linear quadratic regulator*, in International Conference of Machine Learning, 2018.
- 998 [14] A. FLAXMAN, A. KALAI, AND B. MCMAHAN, *Online convex optimization in the bandit setting: Gradient descent*
 999 *without a gradient*, in Proceedings of the sixteenth annual ACM-SIAM symposium on Discrete algorithms, 2005.

- 1000 [15] J. FOERSTER, R. Y. CHEN, M. AL-SHEDIVAT, S. WHITESON, P. ABBEEL, AND I. MORDATCH, *Learning with*
1001 *opponent-learning awareness*, in Proceedings of the 17th International Conference on Autonomous Agents and
1002 MultiAgent Systems, International Foundation for Autonomous Agents and Multiagent Systems, 2018.
- 1003 [16] D. FUDENBERG AND D. K. LEVINE, *The theory of learning in games*, vol. 2, MIT press, 1998.
- 1004 [17] I. GOODFELLOW, J. POUGET-ABADIE, M. MIRZA, B. XU, D. WARDE-FARLEY, S. OZAIR, A. COURVILLE, AND
1005 Y. BENGIO, *Generative adversarial networks*, in Advances in Neural Information Processing Systems 27, 2014.
- 1006 [18] M. W. HIRSCH, *Differential topology*, Springer-Verlag, 1976.
- 1007 [19] C. H. HOMMES AND M. I. OCHEA, *Multiple equilibria and limit cycles in evolutionary games with logit dynamics*,
1008 Games and Economic Behavior, (2012).
- 1009 [20] J. KELLEY, *General Topology*, Van Nostrand Reinhold Company, 1955.
- 1010 [21] J. LEE, *Introduction to smooth manifolds*, Springer, 2012.
- 1011 [22] J. D. LEE, M. SIMCHOWITZ, M. I. JORDAN, AND B. RECHT, *Gradient descent only converges to minimizers*, in
1012 29th Annual Conference on Learning Theory, 2016.
- 1013 [23] D. S. LESLIE AND E. J. COLLINS, *Individual Q-Learning in Normal Form Games*, SIAM J. Control and Optimiza-
1014 tion, (2005).
- 1015 [24] T.-Y. LI AND Z. GAJIC, *Lyapunov Iterations for Solving Coupled Algebraic Riccati Equations of Nash Differential*
1016 *Games and Algebraic Riccati Equations of Zero-Sum Games*, in New Trends in Dynamic Games and Applica-
1017 tions, 1995.
- 1018 [25] P. MERTIKOPOULOS, C. H. PAPANITRIOU, AND G. PILIOURAS, *Cycles in adversarial regularized learning*, in
1019 Proceedings of the 29th annual ACM-SIAM symposium on discrete algorithms, 2018.
- 1020 [26] P. MERTIKOPOULOS AND M. STAUDIGL, *On the convergence of gradient-like flows with noisy gradient input*, SIAM
1021 Journal on Optimization, (2018).
- 1022 [27] P. MERTIKOPOULOS AND Z. ZHOU, *Learning in games with continuous action sets and unknown payoff functions*,
1023 Mathematical Programming, (2019).
- 1024 [28] D. MONDERER AND L. S. SHAPLEY, *Potential games*, Games and Economic Behavior, 14 (1996).
- 1025 [29] S. OMIDSHAFIEI, J. PAZIS, C. AMATO, J. HOW, AND J. VIAN, *Deep decentralized multi-task multi-agent reinforce-*
1026 *ment learning under partial observability*, Arxiv, (2017).
- 1027 [30] M. OSBORNE, *A Course in Game Theory*, MIT Press, 1994.
- 1028 [31] J. PALIS AND S. SMALE, *Structural stability theorems*, Proceedings of the Symposium on Pure Mathematics, (1970).
- 1029 [32] I. PANAGEAS AND G. PILIOURAS, *Gradient descent only converges to minimizers: Non-isolated critical points and*
1030 *invariant regions*, in Innovations in Theoretical Computer Science, 2016.
- 1031 [33] R. PEMANTLE, *Nonconvergence to unstable points in urn models and stochastic approximations*, Annals Probability,
1032 (1990).
- 1033 [34] R. PEMANTLE, *A survey of random processes with reinforcement*, Probability Surveys, 4 (2007).
- 1034 [35] L. J. RATLIFF, S. A. BURDEN, AND S. S. SASTRY, *Characterization and computation of local Nash equilibria*
1035 *in continuous games*, in Proceedings of the 51st Annual Allerton Conference on Communication, Control, and
1036 Computing, 2013.
- 1037 [36] L. J. RATLIFF, S. A. BURDEN, AND S. S. SASTRY, *Genericity and Structural Stability of Non-Degenerate Differ-*
1038 *ential Nash Equilibria*, in Proceedings of the American Control Conference, 2014.
- 1039 [37] L. J. RATLIFF, S. A. BURDEN, AND S. S. SASTRY, *On the Characterization of Local Nash Equilibria in Continuous*
1040 *Games*, IEEE Transactions on Automatic Control, (2016).
- 1041 [38] J. W. ROBBIN, *A structural stability theorem*, Annals of Mathematics, (1971).
- 1042 [39] H. ROBBINS AND D. SIEGMUND, *A convergence theorem for non negative almost supermartingales and some ap-*
1043 *plications*, in Herbert Robbins Selected Papers, Springer New York, 1985.
- 1044 [40] J. B. ROSEN, *Existence and uniqueness of equilibrium points for concave n-person games*, Econometrica, (1965).
- 1045 [41] S. SASTRY, *Nonlinear Systems*, Springer New York, 1999.
- 1046 [42] D. SCIEUR, V. ROULET, F. BACH, AND A. D'ASPROMONT, *Integration methods and optimization algorithms*, in
1047 Advances in Neural Information Processing Systems 30.
- 1048 [43] M. SHUB, *Global Stability of Dynamical Systems*, Springer-Verlag, 1978.
- 1049 [44] S. P. SINGH, M. J. KEARNS, AND Y. MANSOUR, *Nash convergence of gradient dynamics in general-sum games*, in
1050 Proceedings of the 16th Conference on Uncertainty in Artificial Intelligence, 2000.
- 1051 [45] S. SMALE, *Differentiable dynamical systems*, Bulletin of the American Mathematical Society, (1967).
- 1052 [46] R. S. SUTTON AND A. G. BARTO, *Reinforcement Learning: An Introduction*, MIT press, 2017.
- 1053 [47] E. WESSON AND R. RAND, *Hopf bifurcations in delayed rock-paper-scissors replicator dynamics*, Dynamic Games

- 1054 and Applications, (2016).
1055 [48] A. C. WILSON, B. RECHT, AND M. I. JORDAN, *A Lyapunov analysis of momentum methods in optimization*, Arxiv,
1056 (2016).
1057 [49] C. ZHANG AND V. LESSER, *Multi-agent learning with policy prediction*, in Proceedings of the Twenty-Fourth AAAI
1058 Conference on Artificial Intelligence, 2010.