



**HAL**  
open science

# On the convergence of policy gradient methods to Nash equilibria in general stochastic games

Angeliki Giannou, Kyriakos Lotidis, Panayotis Mertikopoulos,  
Emmanouil-Vasileios Vlatakis-Gkaragkounis

► **To cite this version:**

Angeliki Giannou, Kyriakos Lotidis, Panayotis Mertikopoulos, Emmanouil-Vasileios Vlatakis-Gkaragkounis. On the convergence of policy gradient methods to Nash equilibria in general stochastic games. NeurIPS 2022 - 36th International Conference on Neural Information Processing Systems, Nov 2022, New Orleans, United States. pp.1-43. hal-03874018

**HAL Id: hal-03874018**

**<https://hal.science/hal-03874018>**

Submitted on 27 Nov 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# ON THE CONVERGENCE OF POLICY GRADIENT METHODS TO NASH EQUILIBRIA IN GENERAL STOCHASTIC GAMES

ANGELIKI GIANNOU\*, KYRIAKOS LOTIDIS<sup>‡</sup>, PANAYOTIS MERTIKOPOULOS<sup>◊,\*</sup>,  
AND EMMANOUIL V. VLATAKIS-GKARAGKOUNIS<sup>§</sup>

ABSTRACT. Learning in stochastic games is a notoriously difficult problem because, in addition to each other’s strategic decisions, the players must also contend with the fact that the game itself evolves over time, possibly in a very complicated manner. Because of this, the convergence properties of popular learning algorithms – like policy gradient and its variants – are poorly understood, except in specific classes of games (such as potential or two-player, zero-sum games). In view of this, we examine the long-run behavior of policy gradient methods with respect to Nash equilibrium policies that are second-order stationary (SOS) in a sense similar to the type of sufficiency conditions used in optimization. Our first result is that SOS policies are locally attracting with high probability, and we show that policy gradient trajectories with gradient estimates provided by the REINFORCE algorithm achieve an  $\mathcal{O}(1/\sqrt{n})$  distance-squared convergence rate if the method’s step-size is chosen appropriately. Subsequently, specializing to the class of *deterministic* Nash policies, we show that this rate can be improved dramatically and, in fact, policy gradient methods converge within a *finite* number of iterations in that case.

## 1. INTRODUCTION

Ever since they were introduced by Shapley [50] in the 1950’s, stochastic games have been one of the staples of non-cooperative game theory, with a range of pioneering applications to multi-agent reinforcement learning [51], unmanned vehicles [49], general game-playing [37, 52, 57], etc. Informally, a stochastic game unfolds in discrete time as follows: At each point in time, the players are at a given state which determines the rules of the game for that stage. The actions of the players in this state determine not only their instantaneous payoffs (as defined by the stage game), but also the transition probabilities towards the next state of the process. In this way, each player has to balance two distinct – and often competing – objectives: optimizing the payoffs of *today* versus picking a possibly suboptimal action which could yield significant benefits *tomorrow* (i.e., by influencing the transitions of the process towards a more favorable state for the player).

Since all players in the game are involved in a similar dilemma, the decision-making problem for each player is a very complicated affair. In particular, in addition to their changing strategic decisions, the players of the game must also contend with the fact that the

---

\* UNIVERSITY OF WISCONSIN-MADISON.

‡ DEPARTMENT OF MANAGEMENT SCIENCE & ENGINEERING, STANFORD UNIVERSITY.

◊ UNIV. GRENOBLE ALPES, CNRS, INRIA, GRENOBLE INP, LIG, 38000 GRENOBLE, FRANCE.

\* CRITEO AI LAB.

§ UNIVERSITY OF CALIFORNIA, BERKELEY.

*E-mail addresses:* [giannou@wisc.edu](mailto:giannou@wisc.edu), [klotidis@stanford.edu](mailto:klotidis@stanford.edu), [panayotis.mertikopoulos@imag.fr](mailto:panayotis.mertikopoulos@imag.fr),  
[emvlatakis@berkeley.edu](mailto:emvlatakis@berkeley.edu).

2020 *Mathematics Subject Classification.* Primary 91A15, 91A26; secondary 68Q32, 68T05, 90C40.

*Key words and phrases.* Nash equilibrium; stochastic games; policy gradient; stationary policies; strict equilibria.

stage game itself evolves over time. Because of this, even the existence of a Nash equilibrium policy – viz. a stationary Markovian policy that is stable to unilateral deviations [16] – is far more difficult to prove compared to standard, stateless normal form games; for a comprehensive survey, cf. [41, 53] and references therein.

The question we seek to address in this paper is whether an ensemble of boundedly rational players can reach an equilibrium policy in a stochastic game. Specifically, if players do not have sufficient information – or the computational resources required – to solve a high-dimensional Bellman equation [15, 54], it is not at all clear if they would somehow end up playing a Nash policy in the long run. After all, the complexity of most games increases exponentially with the number of players, so the identification of a game’s equilibria quickly becomes prohibitively difficult [27].

**Our contributions in the context of related work.** This issue has sparked a vigorous literature with important ramifications for the range of applications mentioned above. Nevertheless, these efforts must grapple with a series of strong lower bounds for computing even weaker solution concepts like coarse correlated equilibria in turn-based stochastic games [12, 27]. On that account, a recent line of work has focused on establishing convergence in *specific* subclasses of stochastic games, such as *min-max* [7, 11, 32, 47, 48, 58] and common interest *potential* games [13, 31, 61]. However, despite these encouraging results, the general case remains particularly elusive.

Our paper takes a complementary approach to the above and seeks to study the convergence landscape of a class of *equilibrium policies* – not *games*. For concreteness, we focus on the general class of policy gradient methods as pioneered by [28, 29, 55, 59], and we examine the methods’ convergence properties in general random stopping games – as opposed to ergodic stochastic games with an infinite horizon [32, 42]. Concretely, this means that the sequence of play evolves episode-by-episode: within each episode, the players commit a policy and play the game, and from one episode to the next, they use an iterative gradient step to update their policy and continue playing.

Our main contributions in this general context may be summarized as follows:

- (1) We introduce a flexible algorithmic template for the analysis of policy gradient methods which accounts for different information and update frameworks – from perfect policy gradients to value-based estimates obtained on a per-episode basis, e.g., via the REINFORCE algorithm [4, 55, 59].
- (2) Within this framework, we show that Nash policies that satisfy a certain strategic stability condition are locally attracting with arbitrarily high probability. Moreover, to estimate the method’s rate of convergence, we focus on Nash policies that satisfy a second-order sufficiency condition similar to the type of sufficiency conditions used in optimization, and we show that such policies enjoy an  $\mathcal{O}(1/\sqrt{n})$  squared distance convergence rate.
- (3) Finally, we also consider the method’s convergence to *deterministic* Nash policies – a special case of second-order stationary (SOS) policies – and we show that, generically, the above rate can be improved dramatically. In particular, by a simple tweak to the method’s projection step, the induced sequence of play converges to equilibrium in a *finite* number of iterations, despite all the noise and uncertainty.

It is also worth noting that our analysis focuses squarely on the actual, episode-by-episode trajectory of play, not any “best-iterate” or time-averaged variant thereof. In regards to the latter class of guarantees, the recent work of Jin et al. [26] proposed an algorithm (called V-learning) which updates the policy  $\pi_n$  of the  $n$ -th episode based on the observed rewards so far. Thanks to the algorithm’s regret guarantees, Jin et al. [26] showed that (a) in min-max

games, the time-averaged policy  $\bar{\pi}_n = (1/n) \sum_{k=1}^n \pi_k$  converges to equilibrium at a rate of  $\mathcal{O}(1/\sqrt{n})$ ; whereas (b) in *general* stochastic games, the empirical frequency of play converges to the game’s set of coarse correlated equilibria (a substantial relaxation of the notion of Nash equilibrium) at a rate of  $\mathcal{O}(1/\sqrt{n})$ .

By contrast, as we mentioned above, our paper focuses on the *actual* sequence of play, i.e., the policy  $\pi_n$  employed at each episode of the game. Moreover, the rates that we obtain all concern the convergence of the players’ policies to a *Nash* equilibrium – not a correlated equilibrium or other relaxation thereof. In this regard, the best-iterate / ergodic convergence rates are incomparable to our own as they concern a weaker type of convergence (time-averaged instead of the actual sequence), and to a weaker solution concept (correlated equilibria instead of Nash equilibria). This aspect of our results is especially relevant for multi-agent reinforcement learning scenarios where agents learn “on the fly”, and it has important ramifications for many of the practical applications of stochastic games.

From a technical standpoint, our analysis is based on mapping the problem of multi-agent policy learning to the problem of equilibrium learning in a class of continuous games characterized by the fact that first-order stationary points are necessarily Nash (itself a consequence of the so-called “gradient dominance” property of stochastic games). By means of this reframing, we are able to leverage a series of recent techniques for establishing local convergence in (non-monotone) continuous games and variational inequalities [3, 8, 23, 24, 33, 45], which ultimately also yield convergence in our setting. As a result, even though the unbounded variance of the REINFORCE estimator is a source of considerable complications, the resulting link between stochastic and continuous games is of particular technical interest because it opens up a wide array of stochastic approximation tools and techniques that can be used for the analysis of multi-agent learning in stochastic games.

## 2. PRELIMINARIES

**2.1. Setup of the game.** Throughout this work we consider  $N$ -player generic stochastic games where players repeatedly select actions in a shared Markov decision process (MDP) with the goal of maximizing their individual value functions. Formally, we study the tabular version with random stopping of general stochastic games, which is specified by a tuple  $\mathcal{G} = (\mathcal{S}, \mathcal{N}, \{\mathcal{A}_i, R_i\}_{i \in \mathcal{N}}, P, \zeta, \rho)$  with the following primitives:

- A finite set of *agents*  $i \in \mathcal{N} = \{1, 2, \dots, N\}$  and a finite set of *states*  $\mathcal{S} = \{1, \dots, S\}$ .
- For each  $i \in \mathcal{N}$ , a finite space of *actions* (or *pure strategies*)  $\mathcal{A}_i$  indexed by  $\alpha_i = 1, \dots, A_i = |\mathcal{A}_i|$ . We will write  $\mathcal{A} = \prod_{i \in \mathcal{N}} \mathcal{A}_i$  and  $\mathcal{A}_{-i} = \prod_{j \neq i} \mathcal{A}_j$  for the action space of all agents and that of all agents other than  $i$  respectively. In a similar vein, we will also write  $\alpha = (\alpha_i, \alpha_{-i})$  when we want to highlight the action  $\alpha_i$  of player  $i$  against the action profile  $\alpha_{-i}$  of  $i$ ’s opponents.
- For each  $i \in \mathcal{N}$ , we will write  $R_i: \mathcal{S} \times \mathcal{A} \rightarrow [-1, 1]$  for the *reward function* of agent  $i \in \mathcal{N}$ , i.e.,  $R_i(s, \alpha_i, \alpha_{-i})$  will denote the value of the reward of agent  $i$  when the game is at state  $s \in \mathcal{S}$ , the focal agent  $i \in \mathcal{N}$  plays  $\alpha_i \in \mathcal{A}_i$ , and all other agents take actions  $\alpha_{-i} \in \mathcal{A}_{-i}$ .
- The game transits from one state to another according to a Markov transition process, so that  $P(s' | s, \alpha)$  denotes the probability of transitioning from  $s$  to  $s'$  when  $\alpha \in \mathcal{A}$  is the action profile chosen by the agents.
- Given an action profile  $\alpha$  at state  $s$ , the process terminates with probability  $\zeta_{s, \alpha} > 0$ , i.e.,  $\zeta_{s, \alpha} = 1 - \sum_{s' \in \mathcal{S}} P(s' | s, \alpha)$ ; for convenience, we will write  $\zeta := \min_{s, \alpha} \{\zeta_{s, \alpha}\}$ .
- $\rho \in \Delta(\mathcal{S})$  is the distribution for the initial state of the game.

**Episodic Setting.** We consider an episodic setting, where in each episode a realization of the game is completed. At every time step  $t \geq 0$  of each episode, all agents observe the common state  $s_t \in \mathcal{S}$ , select actions  $\alpha_t$  and receive rewards  $\{R_i(s_t, \alpha_t)\}_{i \in \mathcal{N}}$ . Then, with probability  $\zeta_{s_t, \alpha_t}$  the game terminates, and with probability  $1 - \zeta_{s_t, \alpha_t}$ , it moves to the state  $s_{t+1}$ , which is drawn according to  $P(\cdot | s_t, \alpha_t)$ . Denoting the realized reward of player  $i$  at time  $t$  as  $r_{i,t} := R_i(s_t, \alpha_t)$ , we will write  $\tau = (s_t, \alpha_t, r_t)_{t \leq T(\tau)}$  to denote the trajectory of the episode, where  $r_t := (r_{i,t})_{i \in \mathcal{N}}$ , and  $T(\tau)$  the time the episode terminates.

**Policies and value functions.** We consider *stationary Markovian* policies, i.e., policies that do not depend on the time-step and the history, given the current state of the game. More specifically, for each agent  $i \in \mathcal{N}$ , a *policy*  $\pi_i: \mathcal{S} \rightarrow \Delta(\mathcal{A}_i)$  specifies a probability distribution over the actions of agent  $i$  in state  $s \in \mathcal{S}$ , i.e.,  $\alpha_i \sim \pi_i(\cdot | s)$  denotes the (random) action drawn by agent  $i$  at state  $s \in \mathcal{S}$  according to  $\pi_i$ , viewed here as an element of  $\Pi_i := \Delta(\mathcal{A}_i)^{\mathcal{S}}$ . In addition, we will also write  $\pi = (\pi_i)_{i \in \mathcal{N}} \in \Pi := \prod_i \Pi_i$  and  $\pi_{-i} = (\pi_j)_{j \neq i} \in \Pi_{-i} := \prod_{j \neq i} \Pi_j$  for the policy profile of all agents and all agents other than  $i$ , respectively.

The expected reward of agent  $i \in \mathcal{N}$  if agents follow policy  $\pi$ , starting from initial state  $s \in \mathcal{S}$ , defines the *value function* of agent  $i$ , denoted as  $V_{i,s}(\pi)$ , and is equal to

$$V_{i,s}(\pi) := \mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} R_i(s_t, \alpha_t) \middle| s_0 = s \right] \quad (1)$$

where  $\tau \sim \text{MDP}$  denotes the randomness induced by the policy profile  $\pi$ , and the state-transition probabilities of the MDP. Overloading the notation, we set  $V_{i,\rho}(\pi) := \mathbb{E}_{s \sim \rho} [V_{i,s}(\pi)]$ . Although value functions are, in general, non-convex, they share similar smoothness properties with the payoff functions of normal form games, namely bounded and Lipschitz gradients. For precise statements, we defer to the paper’s supplement.

**Visitation distribution and the mismatch coefficient.** For a policy profile  $\pi \in \Pi$  and an arbitrary initial state distribution  $s_0 \sim \rho$ , we define the discounted state visitation measure/distribution as

$$\tilde{d}_\rho^\pi(s) = \mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} \mathbb{1}\{s_t = s\} \middle| s_0 \sim \rho \right], \quad d_\rho^\pi(s) := \tilde{d}_\rho^\pi(s) / Z_\rho^\pi$$

In the appendix, we prove formally that the above definition is well-posed for the random stopping episodic framework described above, i.e.,  $\tilde{d}_\rho^\pi(s) < \infty$ , so  $Z_\rho^\pi := \sum_{s \in \mathcal{S}} \tilde{d}_\rho^\pi(s)$  is well-defined. In our proofs, we will leverage a standard property of visitation distributions, namely the equivalence of the expected value of state-action function and the expected cumulative value over a random trajectory. More precisely, we have:

**Lemma 1.** [Conversion Lemma] *For an arbitrary state-action function  $f: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , a policy profile  $\pi$  and an initial state distribution  $s_0 \sim \rho$ , we have*

$$\mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} f(s_t, \alpha_t) \right] = Z_\rho^\pi \mathbb{E}_{s \sim d_\rho^\pi} \mathbb{E}_{\alpha \sim \pi(\cdot | s)} [f(s, \alpha)] \quad (2)$$

Finally, to quantify the difficulty of hard-to-reach states via a policy gradient method, we will follow the standard approach of [9, 14, 38, 39, 61] and use an appropriately-defined distribution “mismatch coefficient”, generalizing the single-agent counterpart of Agarwal et al. [1]. More precisely, for a stochastic game  $\mathcal{G}$ , we define the *mismatch coefficient* as  $\mathcal{C}_\mathcal{G} := \max_{\pi, \pi' \in \Pi} \{ \|\tilde{d}_\rho^\pi / \tilde{d}_\rho^{\pi'}\|_\infty \}$  or, more simply, as  $\mathcal{C}_\mathcal{G} := \max_{\pi, \pi' \in \Pi} \{ \frac{1}{\zeta} \|\tilde{d}_\rho^\pi / \tilde{d}_\rho^{\pi'}\|_\infty \}$ . Similar to prior work in this direction [1, 5, 11], we will assume  $\mathcal{C}_\mathcal{G}$  is finite, which, equivalently, means that  $d_\rho^\pi(s) > 0$  for any policy  $\pi$  and state  $s$ .

**2.2. Solution concepts.** The most widely used solution concept in game theory is that of a Nash equilibrium i.e., a strategy profile  $\pi^* \in \Pi$  that discourages unilateral deviations. However, in stochastic games, the definition of a Nash policy is much more involved because of the existence of multiple states and steps, cf. [16, 50, 53, 56] and references therein. Formally, we have:

**Definition 1** (Nash policies). A policy  $\pi^* = (\pi_i^*)_{i \in \mathcal{N}} \in \Pi$  is said to be a *Nash policy* for a given distribution of initial states  $\rho \in \Delta(\mathcal{S})$  if, for every player  $i \in \mathcal{N}$ , we have

$$V_{i,\rho}(\pi_i^*; \pi_{-i}^*) \geq V_{i,\rho}(\pi_i; \pi_{-i}^*) \quad \text{for all } i \in \mathcal{N} \text{ and all } \pi_i \in \Delta(\mathcal{A}_i)^{\mathcal{S}}. \quad (\text{NE})$$

In contrast to general non-convex continuous games, stochastic games satisfy a version of the well-known Polyak-Łojasiewicz condition [43] but with linear gradient growth, also known as a *gradient dominance property* (GDP) [1, 5]. For the multi-agent case, Zhang et al. [61] and Daskalakis et al. [11] showed that a similar property holds even in the episodic setting:

**Lemma 2** (Gradient dominance property). *For any policy profile  $\pi = (\pi_i)_{i \in \mathcal{N}} \in \Pi$ , we have that*

$$V_{i,\rho}(\pi_i'; \pi_{-i}) - V_{i,\rho}(\pi_i; \pi_{-i}) \leq \mathcal{C}_G \max_{\bar{\pi}_i \in \Pi_i} \langle \nabla_i V_{i,\rho}(\pi), \bar{\pi}_i - \pi_i \rangle \quad (\text{GDP})$$

for any unilateral deviation  $\pi_i' \in \Pi_i$  of player  $i \in \mathcal{N}$ .

*Remark.* In the above and throughout our paper, we will write  $\nabla_i$  to denote the gradient of the quantity in question with respect to  $\pi_i$ , i.e., when  $\pi_{-i}$  is kept fixed and only  $\pi_i$  is varied. For concision, we will write  $v_i(\pi) = \nabla_i V_{i,\rho}(\pi)$  for the individual gradient of player  $i$ 's value function, and  $v(\pi) = (v_i(\pi))_{i \in \mathcal{N}}$  for the ensemble thereof.  $\blacksquare$

Thanks to (GDP), it is straightforward to check that first-order stationary (FOS) points of  $V$  are Nash. Formally, as in [11, 31, 61], we have the following characterization:

**Lemma 3** (First-order stationary policies are Nash). *A policy  $\pi^* = (\pi_i^*)_{i \in \mathcal{N}} \in \Pi$  is Nash if and only if it satisfies the first-order stationary condition*

$$\langle v(\pi^*), \pi - \pi^* \rangle \leq 0 \quad \text{for all } \pi \in \Pi. \quad (\text{FOS})$$

Leonardos et al. [31] and Zhang et al. [61] proved a relaxation of the above lemma to the effect that policies that satisfy (FOS) up to  $\varepsilon$  (i.e., in lieu of 0 in the RHS) are  $\mathcal{O}(\varepsilon)$ -Nash. Going in the other direction, we will consider the following series of refinements of Nash policies which are particularly important from a learning standpoint [30, 53]:

**Definition 2.** Let  $\pi^* = (\pi_i^*)_{i \in \mathcal{N}} \in \Pi$  be a Nash policy. We then say that:

- $\pi^*$  is *stable* if  $\langle v(\pi), \pi - \pi^* \rangle < 0$  for all  $\pi \neq \pi^*$  sufficiently close to  $\pi^*$ .
- $\pi^*$  is *second-order stationary* if it satisfies the sufficiency condition

$$(\pi - \pi^*)^\top \text{Jac}_v(\pi^*)(\pi - \pi^*) < 0 \quad \text{for all } \pi \in \Pi \setminus \{\pi^*\}, \quad (\text{SOS})$$

where  $\text{Jac}_v(\pi^*) = (\nabla_j v_i(\pi^*))_{i,j \in \mathcal{N}} = (\nabla_j \nabla_i V_i(\pi^*))_{i,j \in \mathcal{N}}$  denotes the Jacobian of  $v$  at  $\pi^*$ .

- $\pi^*$  is *deterministic* if it induces a deterministic selection rule  $\pi_i^*: \mathcal{S} \rightarrow \mathcal{A}_i$  for all  $i \in \mathcal{N}$ .
- $\pi^*$  is *strict* if it is deterministic and (FOS) holds as a strict inequality whenever  $\pi \neq \pi^*$ .

*Remark 1.* In the above and what follows, “sufficiently close” means that there exists a neighborhood  $\mathcal{U}$  of  $\pi^*$  in  $\Pi$  such that the stated inequality holds for all  $\pi \in \mathcal{U}$ . Unless mentioned otherwise, we will measure distances on  $\Pi$  relative to the Euclidean norm, but this choice does not impact our results.

Intuitively, the condition for equilibrium stability is a game-theoretic analogue of first-order KKT sufficiency condition, while the condition for second-order stationarity is the second-order version thereof. In this regard, the distinction between first-order stationary, stable and second-order stationary points is formally analogous to the distinction between critical points, minimizers, and second-order minimum points in optimization. As for deterministic policies, we should mention that, generically, deterministic policies are also strict, so we will use the two terms interchangeably.<sup>1</sup>

Importantly, as we show in [Appendix F](#), these refinements admit the following characterizations:

**Proposition 1.** *Let  $\pi^* = (\pi_i^*)_{i \in \mathcal{N}} \in \Pi$  be a Nash policy. Then:*

a) *If  $\pi^*$  is second-order stationary, there exists some  $\mu > 0$  such that*

$$\langle v(\pi), \pi - \pi^* \rangle \leq -\mu \|\pi - \pi^*\|^2 \quad \text{for all } \pi \text{ sufficiently close to } \pi^*. \quad (3a)$$

b) *If  $\pi^*$  is strict, there exists some  $\mu > 0$  such that*

$$\langle v(\pi), \pi - \pi^* \rangle \leq -\mu \|\pi - \pi^*\| \quad \text{for all } \pi \text{ sufficiently close to } \pi^*. \quad (3b)$$

In view of all the above, we get the following string of implications for equilibria in generic games:

$$\text{strict/deterministic} \implies \text{SOS} \implies \text{stable} \implies \text{FOS} = \text{Nash} \quad (4)$$

For posterity, we should clarify here that, due to the highly complicated structure of the game’s value functions, it is not trivial to construct a concrete example where [\(3a\)](#) holds but [\(3b\)](#) does not. Examples of strict Nash policies abound in the literature [\[30, 53\]](#), but we are not otherwise aware of an argument that could be used to close the gap between [\(3a\)](#) and [\(3b\)](#). In view of this, our analysis will treat both cases concurrently (with the obvious anticipation that more refined solution concepts should enjoy stronger convergence guarantees).

### 3. POLICY GRADIENT METHODS

We now proceed to describe our general model for episodic learning in stochastic games. To that end, we will consider a framework where agents follow a specific policy  $\pi_n$  within each episode, and update it from one episode to the next with the objective of increasing their individual rewards. Formally, our approach will adhere to the following inter-episode sequence of events:

- (1) At the beginning of each episode  $n = 1, 2, \dots$ , every agent  $i \in \mathcal{N}$  chooses a policy  $\pi_{i,n} \in \Pi_i$ .
- (2) Within the  $n$ -th episode, each player executes their chosen policy  $\pi_{i,n}$ , inducing in this way an intra-episode trajectory of play  $\tau_n = (s_t^{(n)}, \alpha_t^{(n)}, r_t^{(n)})_{t \leq T(\tau_n)}$ .
- (3) Once the episode terminates, agents update their policies and the process repeats.

In terms of feedback, we will treat several models, depending on what type of information is available to the agents during play. More precisely, we will focus on the generic policy gradient (PG) template

$$\pi_{n+1} = \text{proj}_{\Pi}(\pi_n + \gamma_n \hat{v}_n) \quad (\text{PG})$$

where:

- (1)  $\pi_n = (\pi_{i,n})_{i \in \mathcal{N}} \in \Pi$  denotes the player’s policy profile at each episode  $n = 1, 2, \dots$
- (2)  $\hat{v}_n = (\hat{v}_{i,n})_{i \in \mathcal{N}} \in \prod_i \mathbb{R}^{\mathcal{A}_i \times \mathcal{S}}$  is an estimate for the agents’ individual policy gradients.

<sup>1</sup>The notion of genericity is stated here in the sense of Baire, i.e., the stated property holds for all but a “meager” set of games (i.e., a countable union of nowhere dense sets in the space of all games).

- (3)  $\text{proj}_{\Pi}: \prod_i \mathbb{R}^{A_i \times S} \rightarrow \Pi$  denotes the Euclidean projection to the agents' policy space  $\Pi$ .
- (4)  $\gamma_n > 0$  is the method's step-size, for which we will assume throughout that  $\sum_n \gamma_n = \infty$ ; typically, (PG) is run with a step-size of the form  $\gamma_n = \gamma/(n+m)^p$  for some  $\gamma > 0$ ,  $m \geq 0$  and  $p \geq 0$ .

Regarding the gradient signal  $\hat{v}_n$ , we will decompose it as

$$\hat{v}_n = v(\pi_n) + U_n + b_n \quad (5)$$

where

$$U_n = \hat{v}_n - \mathbb{E}[\hat{v}_n | \mathcal{F}_n] \quad \text{and} \quad b_n = \mathbb{E}[\hat{v}_n | \mathcal{F}_n] - v(\pi_n). \quad (6)$$

In the above, we treat  $\pi_n$ ,  $n = 1, 2, \dots$ , as a stochastic process on some complete probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ , and we write  $\mathcal{F}_n := \mathcal{F}(\pi_1, \dots, \pi_n) \subseteq \mathcal{F}$  for the history (adapted filtration) of  $\pi_n$  up to – and including – stage  $n$ . By definition,  $\mathbb{E}[U_n | \mathcal{F}_n] = 0$  and  $b_n$  is  $\mathcal{F}_n$ -measurable, so  $U_n$  can be interpreted as a random, zero-mean error relative to  $v(\pi_n)$ , whereas  $b_n$  captures all systematic (non-zero-mean) errors. To make this precise, we will further assume that  $b_n$  and  $U_n$  are bounded as

$$\mathbb{E}[\|b_n\| | \mathcal{F}_n] \leq B_n \quad \text{and} \quad \mathbb{E}[\|U_n\|^2 | \mathcal{F}_n] \leq \sigma_n^2 \quad (7)$$

where the sequences  $B_n$  and  $\sigma_n$ ,  $n = 1, 2, \dots$ , are to be construed as deterministic upper bounds on the bias, fluctuations, and magnitude of the gradient signal  $\hat{v}_n$ .

Depending on these bounds, a gradient signal with  $B_n = 0$  will be called *unbiased*, and an unbiased signal with  $\sigma_n = 0$  will be called *perfect*. More generally, we will assume that the above statistics are bounded as

$$B_n = \mathcal{O}(1/n^{\ell_b}) \quad \text{and} \quad \sigma_n = \mathcal{O}(n^{\ell_\sigma}) \quad (8)$$

for some  $\ell_b, \ell_\sigma > 0$  which depend on the specific model under consideration. For concreteness, we describe below three basic models that adhere to the above template for  $\hat{v}_n$  in order of decreasing information requirements:

**Model 1** (Full gradient information). The first model we will consider assumes that agents observe their *full policy gradients*, i.e.,

$$\hat{v}_n = v(\pi_n) \quad (9)$$

implying in particular that  $U_n = b_n = 0$ . This model is fully deterministic across episodes (though intra-episode play remains stochastic). In particular, it tacitly assumes that agents know the game (and can observe their opponents' policies) so as to calculate the full gradients of their individual value functions  $V_{i,\rho}$ , cf. [2, 31, 61] and references therein. ¶

**Model 2** (Learning with stochastic gradients). A relaxation of the above model which is particularly relevant for applications to deep reinforcement learning concerns the case where the player have access to stochastic policy gradients [60], i.e., unbiased gradient estimates of the form

$$\hat{v}_n = v(\pi_n) + U_n \quad (10)$$

with  $\mathbb{E}[U_n | \mathcal{F}_n] = 0$  (so we can formally take  $\ell_b = \infty$  and  $\ell_\sigma = 0$  in Eq. (8) above). ¶

**Model 3** (Value-based learning). The last model we will consider concerns the case where agents only have access to their instantaneous rewards and need to reconstruct their individual gradients based on this information. A widely used method to achieve this is via the REINFORCE subroutine, which we describe in pseudocode form in Algorithm 1. In words, when employing REINFORCE, each agent  $i \in i$  commits to a sampling policy  $\hat{\pi}_i \in \Pi_i$  and executes it in an episode of the stochastic game in play. Then, at the end of the episode,

**Algorithm 1:** REINFORCE

---

```

1: Input:  $\hat{\pi} \in \Pi, \tau = (s_t, \alpha_t, r_t)_{t \leq T(\tau)} \in \mathcal{T}$ 
2: for  $i = 1, \dots, N$  do
3:    $R_i(\tau) \leftarrow \sum_{t=0}^{T(\tau)} r_{i,t}$ 
4:    $\Lambda_i(\tau) \leftarrow \sum_{t=0}^{T(\tau)} \nabla_i(\log \hat{\pi}_i(\alpha_{i,t}|s_t))$ 
5:    $\hat{v}_i \leftarrow R_i(\tau) \cdot \Lambda_i(\tau)$ 
6: end for
7: return  $\{\hat{v}_i\}_{i \in \mathcal{N}}$ 

```

---

**Algorithm 2:**  $\varepsilon$ -GREEDY POLICY GRADIENT

---

```

1: Input:  $\pi_1, \{\gamma_n\}_{n \in \mathbb{N}}, \{\varepsilon_n\}_{n \in \mathbb{N}}$ 
2: for  $n = 1, 2, \dots$  do
3:    $\hat{\pi}_n \leftarrow (1 - \varepsilon_n)\pi_n + \frac{\varepsilon_n}{|\mathcal{A}|}$ 
4:   Sample  $\tau_n \sim \text{MDP}(\hat{\pi}_n|s_0)$ 
5:    $\hat{v}_n \leftarrow \text{REINFORCE}(\hat{\pi}_n, \tau_n)$ 
6:    $\pi_{n+1} \leftarrow \text{proj}_{\Pi}(\pi_n + \gamma_n \hat{v}_n)$ 
7: end for

```

---

players gather the total reward  $R_i(\tau) \leftarrow \sum_{t=0}^{T(\tau)} r_{i,t}$  associated to the intra-episode trajectory of play  $\tau$ , and they estimate their policy gradients via the so-called “log-trick” [59] as

$$\hat{v}_i = R_i(\tau) \cdot \sum_{t=0}^{T(\tau)} \nabla_i(\log \hat{\pi}_i(\alpha_{i,t}|s_t)). \quad (11)$$

Lemma 4 below provides the vital statistics of the REINFORCE estimator:

**Lemma 4.** *Suppose that each agent  $i \in \mathcal{N}$  follows a stationary policy  $\pi_i \in \Pi_i$ . Then:*

$$a) \quad \mathbb{E}_{\tau \sim \text{MDP}}[\text{REINFORCE}(\pi)] = v(\pi) \quad (12a)$$

$$b) \quad \mathbb{E}_{\tau \sim \text{MDP}}[\|\text{REINFORCE}_i(\pi) - v_i(\pi)\|^2] \leq \frac{24A_i}{\kappa_i \zeta^4} \quad (12b)$$

where  $\kappa_i = \min_{s \in \mathcal{S}, \alpha_i \in \mathcal{A}_i} \pi_i(\alpha_i|s)$ .

Therefore, if REINFORCE is executed at  $\hat{\pi} \leftarrow \pi_n$  at each episode  $n = 1, 2, \dots$ , we will have

$$\mathbb{E}[\hat{v}_{i,n}] = v_i(\pi_n) \quad \text{and} \quad \mathbb{E}[\|U_{i,n}\|^2 | \mathcal{F}_n] \leq \frac{24A_i}{\zeta^4 \min_{s \in \mathcal{S}, \alpha_i \in \mathcal{A}_i} \pi_{i,n}(\alpha_i|s)}. \quad (13)$$

In particular, this means that we will always have  $B_n = 0$  for the bias of the estimator, but its variance could be unbounded if  $\pi_n$  gets close to the boundary of  $\Pi$ . To avoid this, REINFORCE can be paired with an explicit exploration step that modifies the sampling policy of the  $n$ -th episode to

$$\hat{\pi}_{i,n} = (1 - \varepsilon_n)\pi_{i,n} + \varepsilon_n \text{Unif}_{\mathcal{A}_i} \quad \text{for all } s \in \mathcal{S} \quad (14)$$

i.e.,  $\hat{\pi}_{i,n}$  is the mixture between  $\pi_{i,n}$  and the uniform distribution  $\text{Unif}_{\mathcal{A}_i}$  over  $\mathcal{A}_i$ . The resulting algorithm is known as  $\varepsilon$ -GREEDY POLICY GRADIENT; for a pseudocode representation, see Algorithm 2.

Importantly, by calling REINFORCE at  $\hat{\pi}_n$  instead of  $\pi_n$ ,  $\hat{v}_n$  becomes biased (because of the difference between  $\hat{\pi}_n$  and  $\pi_n$ ), but its variance is bounded; in particular, by invoking Lemma 4, we have

$$\mathbb{E}[\|b_{i,n}\| | \mathcal{F}_n] \leq G\varepsilon_n \quad \text{and} \quad \mathbb{E}[\|U_{i,n}\|^2 | \mathcal{F}_n] \leq \frac{24A_i^2}{\varepsilon_n \zeta^4} \quad (15)$$

where  $G$  is a constant that depends on the smoothness of  $V$  and the cardinalities of  $\mathcal{A}$  and  $\mathcal{S}$ .<sup>2</sup> In this way, Algorithm 2 can be seen as a special case of (PG) with  $B_n = \mathcal{O}(\varepsilon_n)$  and  $\sigma_n^2 = \mathcal{O}(1/\varepsilon_n)$ .  $\blacksquare$

<sup>2</sup>Specifically, from Lemma D.7 we know that  $\|v_i(\hat{\pi}_n) - v_i(\pi_n)\| \leq 3\sqrt{A}/\zeta^3 \cdot \sum_j \sqrt{A_j} \cdot \|\hat{\pi}_{j,n} - \pi_{j,n}\|$ . Moreover,  $|\pi_{i,n}(\alpha|s) - \hat{\pi}_{i,n}(\alpha|s)| \leq \varepsilon_n$  for all  $s \in \mathcal{S}, \alpha \in \mathcal{A}_i$ , so  $\|\pi_{i,n} - \hat{\pi}_{i,n}\| \leq \sqrt{SA_i}\varepsilon_n$ . Combining the above, it follows that we can take  $G = 3NA^3/2\sqrt{S}/\zeta^3$ .

## 4. CONVERGENCE ANALYSIS AND RESULTS

We are now in a position to state and discuss our main results. For convenience, we will present our results in order of increasing structure, starting with stable policies, and then moving on to second-order stationary and deterministic Nash policies. All proofs are deferred to the appendix.

**4.1. Asymptotic convergence to stable Nash policies.** Our first convergence result concerns Nash policies that satisfy the stability requirement  $\langle v(\pi), \pi - \pi^* \rangle < 0$  of [Definition 2](#). In this case, we have the following guarantee:

**Theorem 1.** *Let  $\pi^*$  be a stable Nash policy, and let  $\pi_n$  be the sequence of play generated by (PG) with step-size  $\gamma_n = \gamma/(n+m)^p$ ,  $p \in (1/2, 1]$ , and policy gradient estimates such that  $p + \ell_b > 1$  and  $p - \ell_\sigma > 1/2$  as per [\(8\)](#). Then there exists a neighborhood  $\mathcal{U}$  of  $\pi^*$  in  $\Pi$  such that, for any given  $\delta > 0$ , we have*

$$\mathbb{P}(\pi_n \text{ converges to } \pi^* \mid \pi_1 \in \mathcal{U}) \geq 1 - \delta \tag{16}$$

provided that  $\gamma$  is small enough (or  $m$  large enough) relative to  $\delta$ .

**Corollary 1.** *Suppose that [Models 1–3](#) are run with a step-size of the form  $\gamma_n = \gamma/(n+m)^p$ ,  $p > 1/2$ , and if applicable, an exploration parameter  $\varepsilon_n = \varepsilon/(n+m)^r$  such that  $1 - p < r < 2p - 1$ . Then:*

- For [Models 1 and 2](#): the conclusions of [Theorem 1](#) hold as stated.
- For [Model 3](#): the conclusions of [Theorem 1](#) hold as long as  $p > 2/3$ .

We note here that [Theorem 1](#) provides a trajectory convergence guarantee which is otherwise quite difficult to obtain even in structured stochastic games. For example, if we zoom in on the class of stochastic potential (or min-max) games, the existing guarantees in the literature concern the “best iterate” of the algorithm, cf. [\[31, 61\]](#) and references therein. Because of this, said guarantees do not apply to the actual trajectory of play generated by (PG); this makes them less suitable for agent-based learning where the players involved are learning “as they go”, as opposed to *simulating* the game in order to approximately compute an equilibrium policy offline.

We should also note that the convergence guarantees of [Theorem 1](#) hold locally with arbitrarily high probability. Without further assumptions, it is not possible to obtain global trajectory convergence guarantees that hold with probability 1, even in single-state games – that is, the case of learning in finite normal form games. The reason for this locality is twofold: First, equilibrium policies are not unique in general, and gradient-based dynamics may also admit non-equilibrium attractors, such as limit cycles and the like [\[25, 34–36\]](#). As a result, in the presence of multiple equilibria/attractors, the best one can hope for is a local equilibrium convergence result, conditioned on the basin of attraction of said equilibrium (as per [Theorem 1](#)).

The second obstruction to a global, unconditional convergence result is probabilistic in nature, and has to do with the randomness that enters the learning process (e.g., in the estimation of policy gradients via the REINFORCE). In this case, no matter how close one starts to an equilibrium policy, there is always a finite, non-zero probability that an unlucky realization of the noise can drive the process away from its basin, possibly never to return. This issue can only be overcome in games where  $\Pi$  is partitioned (up to a set of measure zero) into basins of attraction of equilibrium policies. However, this can only occur in games with a sufficiently strong global structure, like potential stochastic games, two-player zero-sum games and the like; in complete generality, locality cannot be lifted, even in single-state problems [\[17, 19\]](#).

**4.2. Convergence to second-order stationary policies.** Albeit valuable as an asymptotic convergence guarantee, [Theorem 1](#) does not provide an indication of how long it will take players to actually converge to a Nash policy. Of course, in full generality, it is not plausible to expect to be able to derive such a convergence rate because the stability requirement provides no indication on how fast the players' policy gradients stabilize near a solution. This kind of estimate is provided by the second-order sufficient condition (SOS), which allows us to establish sufficient control over the sequence of play as indicated by the following theorem.

**Theorem 2.** *Let  $\pi^*$  be a second-order stationary policy, let  $\mathcal{B}$  be a neighborhood of  $\pi^*$  such that [\(3a\)](#) holds on  $\mathcal{B}$ , and let  $\pi_n$  be the sequence of play generated by (PG) with step-size  $\gamma_n = \gamma/(n+m)^p$ ,  $p \in (1/2, 1]$ , and policy gradient estimates such that  $p + \ell_b > 1$  and  $p - \ell_\sigma > 1/2$  as per [\(8\)](#). Then:*

- (1) *There exists a neighborhood  $\mathcal{U}$  of  $\pi^*$  in  $\Pi$  such that, for any confidence level  $\delta > 0$ , the event*

$$\mathcal{E} = \{\pi_n \in \mathcal{B} \text{ for all } n = 1, 2, \dots\} \quad (17)$$

*occurs with probability  $\mathbb{P}(\mathcal{E} \mid \pi_1 \in \mathcal{U}) \geq 1 - \delta$  if  $m$  is large enough relative to  $\delta$ .*

- (2) *The sequence  $\pi_n$  converges to  $\pi^*$  with probability 1 on  $\mathcal{E}$ ; in particular, we have*

$$\mathbb{P}(\pi_n \text{ converges to } \pi^* \mid \pi_1 \in \mathcal{U}) \geq 1 - \delta \quad (18)$$

*if  $m$  is large relative to  $\delta$ . Moreover, conditioned on  $\mathcal{E}$  and taking  $q = \min\{\ell_b, p - 2\ell_\sigma\}$ , we have*

$$\mathbb{E}[\|\pi_n - \pi^*\|^2 \mid \mathcal{E}] = \begin{cases} \mathcal{O}(1/n^{2\mu\gamma}) & \text{if } p = 1 \text{ and } 2\mu\gamma < q, \\ \mathcal{O}(1/n^q) & \text{otherwise.} \end{cases} \quad (19)$$

**Corollary 2.** *Suppose that [Models 1–3](#) are run with a step-size of the form  $\gamma_n = \gamma/(n+m)^p$ ,  $p > 1/2$ , and if applicable, an exploration parameter  $\varepsilon_n = \varepsilon/(n+m)^{p/2}$ . Then:*

- *For [Models 1 and 2](#): the conclusions of [Theorem 2](#) hold with  $q = p$ ; in particular, [\(19\)](#) gives an  $\mathcal{O}(1/n)$  rate of convergence if  $p = 1$  and  $2\mu\gamma > q$ .*
- *For [Model 3](#): the conclusions of [Theorem 2](#) hold for  $p > 2/3$  with  $q = p/2$ ; in particular, [\(19\)](#) gives an  $\mathcal{O}(1/\sqrt{n})$  rate of convergence if  $p = 1$  and  $2\mu\gamma > q$ .*

*Remark 2.* Getting an explicit estimate for the constant in the  $\mathcal{O}(\cdot)$  guarantee of [Theorem 2](#) is quite involved but, up to logarithmic and subleading factors, Chung's lemma [[10](#), [44](#)] can be used to show that a) if  $2\mu\gamma > q$ , it scales as  $(C_b + C_\sigma)/[(2\mu\gamma - q)(1 - \delta)]$  where  $C_b = \sup_n \gamma_n B_n$  and  $C_\sigma = \sup_n \gamma_n^2 \sigma_n^2$ ; b) if  $2\mu\gamma = q$ , it scales as  $(C_b + C_\sigma)(1 + \max\{(2\mu\gamma)^2, 4\mu\gamma\})/(1 - \delta)$ ; and c) if  $2\mu\gamma < q$  as  $(C_b + C_\sigma)(1 + \max\{(2\mu\gamma)^2, 4\mu\gamma\})/[(q - 2\mu\gamma)(1 - \delta)]$ .

Besides providing a general framework for achieving trajectory convergence, [Theorem 2](#) gives the rates of convergence of the sequence of play to the Nash policy in question. In particular, with this result in hand, one can confidently argue about the distance of the iterates of (PG) from equilibrium in a series of different environments. More to the point, this convergence guarantee allows the algorithm designer to adapt the parameters of the learning process according to the complexity and limitations of the environment, a feature which further highlights the significance of this result.

We should also note the delicate interplay between the method's step-size and the achieved convergence rate. In the case of [Model 1](#), [Corollary 2](#) suggests a step-size of the form  $\gamma_n = \Theta(1/n)$ , leading to a  $\mathcal{O}(1/n)$  convergence rate. As we show in the appendix, this rate can be improved: in the deterministic case with perfect gradient information, (PG) with a suitably chosen constant step-size achieves a *geometric* convergence rate, i.e.,

$\|\pi_n - \pi^*\| = \mathcal{O}(\exp(-\rho n))$  for some  $\rho > 0$  (cf. [Proposition B.1](#) in [Appendix B](#)). By contrast, in the case of [Model 2](#), the  $\mathcal{O}(1/n)$  rate we provide cannot be improved, even if the quadratic minorant [\(3a\)](#) that characterizes SOS policies holds *globally* – and this because the learning process is running against standard lower bounds from convex optimization [\[6, 40\]](#).

Perhaps the most significant guarantee from a practical point of view is the  $\mathcal{O}(1/\sqrt{n})$  convergence rate attained in [Model 3](#) (cf. [Algorithms 1](#) and [2](#)). This guarantee amounts to a  $\mathcal{O}(1/n^{1/4})$  convergence rate in terms of the (non-squared) distance to equilibrium which, mutatis mutandis, represents a notable improvement over the  $\mathcal{O}(1/n^{1/6})$  guarantee of Leonardos et al. [\[31\]](#) (expressed in norm values). Of course, the latter guarantee is global – because the focus of [\[31\]](#) is stochastic *potential* games – but it also concerns the “best iterate” of the process (not its “last iterate”), so the two results are not immediately comparable. However, a useful “best-of-both-worlds” heuristic that can be inferred by the combination of these works is that, given a budget of training episodes, [Algorithm 2](#) can be run with a constant step-size as per [\[31\]](#) for a sufficient fraction of this budget, and then with a  $\mathcal{O}(1/n)$  “cooldown” schedule for the rest. In this way, after an aggressive “exploration” phase, the algorithm’s  $\mathcal{O}(1/n^{1/4})$  rate would kick in and supply faster stabilization to an SOS policy.

**4.3. Convergence to deterministic Nash policies.** Our last series of results concerns the rate of convergence to deterministic Nash policies in generic stochastic games. As we discussed in [Section 2](#), deterministic Nash policies also satisfy (SOS), so the rate of convergence of (PG) to such policies can be harvested directly from [Theorem 2](#). However, as we show below, a simple projection tweak in (SOS) can improve this rate dramatically.

The tweak in question is inspired by the geometry of  $\Pi$  around a deterministic policy: by definition, such policies are corner points of  $\Pi$ , so any consistent drift towards them will cause  $\pi_n$  to hit the boundary of  $\Pi$  in finite time. Of course, under (PG), the process may rebound from the boundary and return to the interior of  $\Pi$  if the policy gradient estimate is not particularly good at a given iteration of the algorithm. However, if we replace the projection step of (PG) with a “lazy projection” in the spirit of Zinkevich [\[62\]](#), the aggregation of gradient steps will eventually push the process far inside the normal cone of  $\Pi$  at  $\pi^*$ , so rebounds of this type can no longer occur.

Formally, we will consider the following *lazy policy gradient* (LPG) scheme:

$$y_{n+1} = y_n + \gamma_n \hat{v}_n \quad \pi_{n+1} = \text{proj}_{\Pi}(y_{n+1}) \quad (\text{LPG})$$

where  $y_n = (y_{i,n})_{i \in \mathcal{N}} \in \prod_i \mathbb{R}^{\mathcal{A}_i \times \mathcal{S}}$  is an auxiliary variable that maintains an aggregate of gradient steps *before* projecting them back to  $\Pi$ . We then have the following convergence result:

**Theorem 3.** *Let  $\pi_n$  be the sequence of play under (LPG) with step-size and policy gradient estimates such that  $p + \ell_b > 1$  and  $p - \ell_\sigma > 1/2$  as per [\(8\)](#). If  $\pi^*$  is a deterministic Nash policy, there exists an unbounded open set  $\mathcal{W} \subseteq \prod_i \mathbb{R}^{\mathcal{A}_i \times \mathcal{S}}$  of initializations such that, for any  $\delta > 0$ , we have*

$$\mathbb{P}(\pi_n \text{ converges to } \pi^* \mid y_1 \in \mathcal{W}) \geq 1 - \delta, \quad (20)$$

*provided that  $\gamma > 0$  is small enough. Moreover, conditioned on this event,  $\pi_n$  converges to  $\pi^*$  at a finite number of iterations, i.e., there exists some  $n_0$  such that  $\pi_n = \pi^*$  for all  $n \geq n_0$ .*

**Corollary 3.** *Suppose that [Models 1–3](#) are run with parameters  $\gamma_n = \gamma/n^p$ ,  $p \in (1/2, 1]$ , and if applicable,  $\varepsilon_n = \varepsilon/n^r$  with  $1 - p < r < 2p - 1$ . Then the conclusions of [Theorem 3](#) hold.*

*Remark 3.* Getting an explicit bound for  $n_0$  is quite complicated, but the last part of the proof of [Theorem 3](#) shows that  $n_0$  scales in terms of the parameters of the game and the algorithm as  $n_0 = \mathcal{O}\left(\left(\frac{MSA}{c\gamma}\right)^{1/(1-p)}\right)$  where  $c > 0$  measures the minimum payoff difference

between equilibrium and non-equilibrium strategies at  $\pi^*$ ,  $M$  is a measure of the initial distance from  $\pi^*$ , and  $S$  and  $A$  is the number of states and pure strategies respectively.

**Theorem 3** – and, by extension, **Corollary 3** – are fairly unique because they provide a guarantee for convergence to an *exact* Nash equilibrium in a *finite* number of iterations. To the best of our knowledge, the only comparable result in the literature is that of [61], where the authors provide a finite-time convergence guarantee to strict equilibria with *perfect* policy gradients (as per **Model 1**). The result of Zhang et al. [61] echoes the convergence properties of deterministic first-order algorithms around sharp minima of convex functions [44], but the fact that **Theorem 3** applies to models with *stochastic* gradient feedback of *unbounded* variance (**Models 2** and **3** respectively) is a major difference. As far as we are aware, this is the first guarantee of its kind in the literature on learning in stochastic games.

## 5. CONCLUDING REMARKS

A key roadblock encountered by practical applications of multi-agent reinforcement learning is the lack of universal equilibrium convergence guarantees. While the impossibility results of [21, 22] imply that unconditional convergence is not a reasonable aspiration without further assumptions on the game, the existence of local convergence results mitigates this deficiency as it provides a range of theoretically grounded stability and runtime guarantees. In this regard, deterministic policies acquire particular importance, as the convergence of policy gradient methods is especially rapid and robust in this case. Of course, this leaves open the question of non-tabular settings and parametrically encoded policies, e.g., as in the case of deep reinforcement learning; we defer these investigations to future work.

Another open issue of high practical relevance concerns policy gradient methods that do not rely on Euclidean projections to  $\Pi$ . In the single-state case (i.e., learning in finite normal form games), the use of methods relying on softmax choice / exponential weights is very widely used because of its regret guarantees. Whether the use of similar softmax techniques can lead to finer convergence guarantees in the context of general stochastic games is an important and intriguing question for future research.

## APPENDIX A. ASYMPTOTIC CONVERGENCE TO STABLE NASH POLICIES

Our goal in this appendix is to prove **Theorem 1** and **Corollary 1**, which we restate below for convenience:

**Theorem 1.** *Let  $\pi^*$  be a stable Nash policy, and let  $\pi_n$  be the sequence of play generated by (PG) with step-size  $\gamma_n = \gamma/(n+m)^p$ ,  $p \in (1/2, 1]$ , and policy gradient estimates such that  $p + \ell_b > 1$  and  $p - \ell_\sigma > 1/2$  as per (8). Then there exists a neighborhood  $\mathcal{U}$  of  $\pi^*$  in  $\Pi$  such that, for any given  $\delta > 0$ , we have*

$$\mathbb{P}(\pi_n \text{ converges to } \pi^* \mid \pi_1 \in \mathcal{U}) \geq 1 - \delta \quad (16)$$

*provided that  $\gamma$  is small enough (or  $m$  large enough) relative to  $\delta$ .*

**Corollary 1.** *Suppose that **Models 1–3** are run with a step-size of the form  $\gamma_n = \gamma/(n+m)^p$ ,  $p > 1/2$ , and if applicable, an exploration parameter  $\varepsilon_n = \varepsilon/(n+m)^r$  such that  $1 - p < r < 2p - 1$ . Then:*

- For **Models 1** and **2**: the conclusions of **Theorem 1** hold as stated.
- For **Model 3**: the conclusions of **Theorem 1** hold as long as  $p > 2/3$ .

Our proof strategy will comprise the following basic steps:

- (1) To begin with, we will show that the squared distance

$$D(\pi) = \frac{1}{2} \|\pi - \pi^*\|^2 \quad (\text{A.1})$$

can be seen as a “local Lyapunov function” for (PG) in the sense that it is locally decreasing near  $\pi^*$ , up to a series of error terms – both zero-mean and non-zero-mean.

- (2) Due to these errors, the evolution of the iterates  $D_n := D(\pi_n)$  of  $D$  over time may exhibit *significant* jumps: in particular, a single “bad” realization of the noise could carry  $\pi_n$  out of the basin of attraction of  $\pi^*$ , possibly never to return. To exclude this event, our second step will be to show that the aggregation of these errors can be controlled with probability at least  $1 - \delta$ .
- (3) Conditioned on the above, we will show that, with probability at least  $1 - \delta$ , the iterates  $D_n$  cannot grow more than a token value. As a result, if (PG) is initialized close to  $\pi^*$ , it will remain in a neighborhood thereof for all  $n$  (again, with probability at least  $1 - \delta$ ).
- (4) Thanks to this “stochastic Lyapunov stability” result, we employ a series of martingale limit theory arguments to extract a subsequence converging to  $\pi^*$ .
- (5) Finally, we show that the increments of  $D_n$  are summable; hence, by invoking the Gladyshev’s lemma [44, p. 49], we conclude that  $D_n$  converges to some (finite) random variable  $D_\infty$ . Combining this fact with the existence of a convergent subsequence, we obtain the desired conclusion that  $\pi_n$  converges to  $\pi^*$  with probability at least  $1 - \delta$ .

In the sequel, we make the above precise in a series of intermediate results.

**A.1. Energy inequality.** We begin by establishing a “quasi-Lyapunov” inequality for the iterates  $D_n = \|\pi_n - \pi^*\|^2/2$  of (A.1).

**Lemma A.1.** *Let  $D_n := D(\pi_n)$ . Then, for all  $n = 1, 2, \dots$ , we have*

$$D_{n+1} \leq D_n + \gamma_n \langle v(\pi_n), \pi_n - \pi^* \rangle + \gamma_n \xi_n + \gamma_n \chi_n + \gamma_n^2 \psi_n^2, \quad (\text{A.2})$$

where the error terms  $\xi_n$ ,  $\chi_n$ , and  $\psi_n$  are given by

$$\xi_n = \langle U_n, \pi_n - \pi^* \rangle, \quad \chi_n = \|\Pi\| B_n \quad \text{and} \quad \psi_n^2 = \frac{1}{2} \|\hat{v}_n\|^2. \quad (\text{A.3})$$

with  $\|\Pi\| := \max_{\pi, \pi' \in \Pi} \|\pi - \pi'\|$ .

*Proof.* By the definition of the iterates of (PG), we have

$$\begin{aligned} D_{n+1} &= \frac{1}{2} \|\pi_{n+1} - \pi^*\|^2 = \frac{1}{2} \|\text{proj}_\Pi(\pi_n + \gamma_n \hat{v}_n) - \text{proj}_\Pi(\pi^*)\|^2 \\ &\leq \frac{1}{2} \|\pi_n + \gamma_n \hat{v}_n - \pi^*\|^2 \\ &= \frac{1}{2} \|\pi_n - \pi^*\|^2 + \gamma_n \langle \hat{v}_n, \pi_n - \pi^* \rangle + \frac{1}{2} \gamma_n^2 \|\hat{v}_n\|^2 \\ &= D_n + \gamma_n \langle v(\pi_n) + U_n + b_n, \pi_n - \pi^* \rangle + \frac{1}{2} \gamma_n^2 \|\hat{v}_n\|^2 \\ &\leq D_n + \gamma_n \langle v(\pi_n), \pi_n - \pi^* \rangle + \gamma_n \xi_n + \gamma_n \chi_n + \gamma_n^2 \psi_n^2 \end{aligned} \quad (\text{A.4})$$

where we used the Cauchy-Schwarz inequality to bound the bias term as  $\langle b_n, \pi_n - \pi^* \rangle \leq \|b_n\| \cdot \|\pi_n - \pi^*\| \leq \|\Pi\| B_n = \chi_n$ . ■

**A.2. Error control and stability.** The second major step in our proof (and the most challenging one from a technical standpoint) is to establish a suitable measure of control over the error increments in (A.1), with the aim of showing that the process  $\pi_n$  never leaves a neighborhood of  $\pi^*$ .

To make this idea precise, let  $\mathcal{B} = \{\pi \in \Pi : \|\pi - \pi^*\| \leq \varrho\}$  be a ball of radius  $\varrho$  based on  $\pi^*$  in  $\Pi$  so that  $\langle v(\pi), \pi - \pi^* \rangle < 0$  for all  $\pi \in \mathcal{B} \setminus \{\pi^*\}$  (without loss of generality, we can assume that  $\mathcal{B}$  is maximal in that regard). We will then examine the event that the aggregation of the error terms in (A.1) is not sufficient to drive  $\pi_n$  to escape from  $\mathcal{B}$ .

To that end, we will begin by aggregating the errors in (A.1) as

$$M_n = \sum_{k=1}^n \gamma_k \xi_k \quad \text{and} \quad S_n = \sum_{k=1}^n [\gamma_k \chi_k + \gamma_k^2 \psi_k^2]. \quad (\text{A.5})$$

Since  $\mathbb{E}[\xi_n | \mathcal{F}_n] = 0$ , we have  $\mathbb{E}[M_n | \mathcal{F}_n] = M_{n-1}$ , so  $M_n$  is a martingale; likewise,  $\mathbb{E}[S_n | \mathcal{F}_n] \geq S_{n-1}$ , so  $S_n$  is a submartingale. Then, using a technique of Hsieh et al. [23] that builds on an earlier idea by Mertikopoulos & Zhou [33], we will also consider the “mean square” error process

$$W_n = M_n^2 + S_n, \quad (\text{A.6})$$

and the associated indicator events

$$\mathcal{E}_n = \{\pi_k \in \mathcal{B} \text{ for all } k = 1, 2, \dots, n\} \quad \text{and} \quad H_n = \{W_k \leq a \text{ for all } k = 1, 2, \dots, n\}, \quad (\text{A.7a})$$

where, with a fair amount of hindsight, the error tolerance level  $a > 0$  is such that  $2a + \sqrt{a} < \varrho$ , and we are employing the convention  $\mathcal{E}_0 = H_0 = \Omega$  (since every statement is true for the elements of the empty set). We will then assume that  $\pi_1$  is initialized in a ball of radius  $\sqrt{2a}$  centered at  $\pi^*$ , viz.

$$\mathcal{U} = \{\pi \in \Pi : D(\pi) \leq a\} = \{\pi \in \Pi : \|\pi - \pi^*\|^2 / 2 \leq a\}. \quad (\text{A.8})$$

With all this in hand, the key to showing that  $\pi_n$  remains close to  $\pi^*$  with high probability is the following conditional estimate:

**Lemma A.2.** *Let  $\pi_n$  be the sequence of play generated by (PG) initialized at  $\pi_1 \in \mathcal{U}$ . We then have:*

- (1)  $\mathcal{E}_{n+1} \subseteq \mathcal{E}_n$  and  $H_{n+1} \subseteq H_n$  for all  $n = 1, 2, \dots$
- (2)  $H_{n-1} \subseteq \mathcal{E}_n$  for all  $n = 1, 2, \dots$
- (3) Consider the “bad realization” event

$$\tilde{H}_n := H_{n-1} \setminus H_n = \{W_k \leq a \text{ for } k = 1, 2, \dots, n-1 \text{ and } W_n > a\}, \quad (\text{A.9})$$

and let  $\tilde{W}_n = W_n \mathbf{1}_{H_{n-1}}$  be the cumulative error subject to the noise being “small”. Then we have:

$$\mathbb{E}[\tilde{W}_n] \leq \mathbb{E}[\tilde{W}_{n-1}] + \gamma_n \|\Pi\| B_n + \gamma_n^2 \|\Pi\|^2 \sigma_n^2 + \frac{3}{2} \gamma_n^2 (G^2 + B_n^2 + \sigma_n^2) - a \mathbb{P}(\tilde{H}_{n-1}), \quad (\text{A.10})$$

where, by convention,  $\tilde{H}_0 = \emptyset$  and  $\tilde{W}_0 = 0$ .

*Remark.* In the above (and what follows), the notation  $\mathbf{1}_A$  is used to indicate the logical indicator of an event  $A \subseteq \Omega$ , i.e.,  $\mathbf{1}_A(\omega) = 1$  if  $\omega \in A$  and  $\mathbf{1}_A(\omega) = 0$  otherwise.

The proof of Lemma A.2 is quite technical, so we first proceed to derive an important stability result based on this estimate.

**Proposition A.1.** Fix some confidence threshold  $\delta > 0$  and let  $\pi_n$  be the sequence of play generated by (PG) with step-size and policy gradient estimates as per Theorem 1. We then have:

$$\mathbb{P}(H_n \mid \pi_1 \in \mathcal{U}) \geq 1 - \delta \quad \text{for all } n = 1, 2, \dots \quad (\text{A.11})$$

provided that  $\gamma$  is small enough (or  $m$  large enough) relative to  $\delta$ .

*Proof.* We begin by bounding the probability of the “bad realization” event  $\tilde{H}_n = H_{n-1} \setminus H_n$ . Indeed, if  $\pi_1 \in \mathcal{U}$ , we have:

$$\mathbb{P}(\tilde{H}_n) = \mathbb{P}(H_{n-1} \setminus H_n) = \mathbb{E}[\mathbb{1}_{H_{n-1}} \times \mathbb{1}\{W_n > a\}] \leq \mathbb{E}[\mathbb{1}_{H_{n-1}} \times (W_n/a)] = \mathbb{E}[\tilde{W}_n]/a \quad (\text{A.12})$$

where, in the penultimate step, we used the fact that  $W_n \geq 0$  (so  $\mathbb{1}\{W_n > a\} \leq W_n/a$ ). Telescoping (A.10) then yields

$$\mathbb{E}[\tilde{W}_n] \leq \mathbb{E}[\tilde{W}_0] + \|\Pi\| \sum_{k=1}^n \gamma_k B_k + \sum_{k=1}^n \gamma_k^2 \varrho_k^2 - a \sum_{k=1}^n \mathbb{P}(\tilde{H}_{k-1}) \quad (\text{A.13})$$

where we set

$$\varrho_n^2 = \|\Pi\|^2 \sigma_n^2 + \frac{3}{2}(G^2 + B_n^2 + \sigma_n^2). \quad (\text{A.14})$$

Hence, combining (A.12) and (A.13) and invoking our stated assumptions for  $\gamma_n$ ,  $B_n$  and  $\sigma_n$ , we get

$$\sum_{k=1}^n \mathbb{P}(\tilde{H}_k) \leq \frac{1}{a} \sum_{k=1}^n [\gamma_k B_k \|\Pi\| + \gamma_k^2 \varrho_k^2] \leq \frac{C}{a} \quad (\text{A.15})$$

for some  $C \equiv C(\gamma, m) > 0$  with  $\lim_{\gamma \rightarrow 0^+} C(\gamma, m) = \lim_{m \rightarrow \infty} C(\gamma, m) = 0$  (since  $\gamma_n = \gamma/(n+m)^p$  and  $p > 0$ ).

Now, by choosing  $\gamma$  sufficiently small (or  $m$  sufficiently large), we can ensure that  $C/a < \delta$ ; thus, given that the events  $\tilde{H}_k$  are disjoint for all  $k = 1, 2, \dots$ , we get  $\mathbb{P}(\bigcup_{k=1}^n \tilde{H}_k) = \sum_{k=1}^n \mathbb{P}(\tilde{H}_k) \leq \delta$ . In turn, this implies that  $\mathbb{P}(H_n) = \mathbb{P}(\tilde{H}_1^c \cap \dots \cap \tilde{H}_n^c) \geq 1 - \delta$ , and our assertion follows. ■

We conclude this appendix with the proof of our technical result on the events  $\mathcal{E}_n$  and  $H_n$ :

*Proof of Lemma A.2.* The first claim of the lemma is obvious. For the second, we proceed inductively:

- (1) For the base case  $n = 1$ , we have  $\mathcal{E}_1 = \{\pi_1 \in \mathcal{B}\} \supseteq \{\pi_1 \in \mathcal{U}\} = \Omega$  (recall that  $\pi_1$  is initialized in  $\mathcal{U} \subseteq \mathcal{B}$ ). Since  $H_0 = \Omega$ , our claim follows.
- (2) Inductively, assume that  $H_{n-1} \subseteq \mathcal{E}_n$  for some  $n \geq 1$ . To show that  $H_n \subseteq \mathcal{E}_{n+1}$ , suppose that  $W_k \leq a$  for all  $k = 1, 2, \dots, n$ . Since  $H_n \subseteq H_{n-1}$ , this implies that  $\mathcal{E}_n$  also occurs, i.e.,  $\pi_k \in \mathcal{B}$  for all  $k = 1, 2, \dots, n$ ; as such, it suffices to show that  $\pi_{n+1} \in \mathcal{B}$ . To do so, given that  $\pi_k \in \mathcal{U} \subseteq \mathcal{B}$  for all  $k = 1, 2, \dots, n$ , we readily obtain

$$D_{k+1} \leq D_k + \gamma_k \xi_k + \gamma_k \chi_k + \gamma_k^2 \psi_k^2, \quad \text{for all } k = 1, 2, \dots, n, \quad (\text{A.16})$$

and hence, after telescoping over  $k = 1, 2, \dots, n$ , we get

$$D_{n+1} \leq D_1 + M_n + S_n \leq D_1 + \sqrt{W_n} + W_n \leq a + \sqrt{a} + a = 2a + \sqrt{a}. \quad (\text{A.17})$$

We conclude that  $D(\pi_{n+1}) \leq 2a + \sqrt{a}$ , i.e.,  $\pi_{n+1} \in \mathcal{B}$ , as required for the induction.

For our third claim, note first that

$$\begin{aligned} W_n &= (M_{n-1} + \gamma_n \xi_n)^2 + S_{n-1} + \gamma_n \chi_n + \gamma_n^2 \psi_n^2 \\ &= W_{n-1} + 2\gamma_n \xi_n M_{n-1} + \gamma_n^2 \xi_n^2 + \gamma_n \chi_n + \gamma_n^2 \psi_n^2, \end{aligned} \quad (\text{A.18})$$

so, after taking expectations, we get

$$\mathbb{E}[W_n | \mathcal{F}_n] = W_{n-1} + 2M_{n-1}\gamma_n \mathbb{E}[\xi_n | \mathcal{F}_n] + \mathbb{E}[\gamma_n^2 \xi_n^2 + \gamma_n \chi_n + \gamma_n^2 \psi_n^2 | \mathcal{F}_n] \geq W_{n-1}, \quad (\text{A.19})$$

i.e.,  $W_n$  is a submartingale. To proceed, let  $\tilde{W}_n = W_n \mathbb{1}_{H_{n-1}}$  so

$$\begin{aligned} \tilde{W}_n &= W_n \mathbb{1}_{H_{n-1}} = W_{n-1} \mathbb{1}_{H_{n-1}} + (W_n - W_{n-1}) \mathbb{1}_{H_{n-1}} \\ &= W_{n-1} \mathbb{1}_{H_{n-2}} - W_{n-1} \mathbb{1}_{\tilde{H}_{n-1}} + (W_n - W_{n-1}) \mathbb{1}_{H_{n-1}}, \\ &= \tilde{W}_{n-1} + (W_n - W_{n-1}) \mathbb{1}_{H_{n-1}} - W_{n-1} \mathbb{1}_{\tilde{H}_{n-1}}, \end{aligned} \quad (\text{A.20})$$

where we used the fact that  $H_{n-1} = H_{n-2} \setminus \tilde{H}_{n-1}$  so  $\mathbb{1}_{H_{n-1}} = \mathbb{1}_{H_{n-2}} - \mathbb{1}_{\tilde{H}_{n-1}}$  (since  $H_{n-1} \subseteq H_{n-2}$ ). Then, (A.18) yields

$$W_n - W_{n-1} = 2M_{n-1}\gamma_n \xi_n + \gamma_n^2 \xi_n^2 + \gamma_n \chi_n + \gamma_n^2 \psi_n^2 \quad (\text{A.21})$$

and hence, given that  $H_{n-1}$  is  $\mathcal{F}_n$ -measurable, we get:

$$\mathbb{E}[(W_n - W_{n-1}) \mathbb{1}_{H_{n-1}}] = 2 \mathbb{E}[\gamma_n M_{n-1} \xi_n \mathbb{1}_{H_{n-1}}] \quad (\text{A.22a})$$

$$+ \mathbb{E}[\gamma_n^2 \xi_n^2 \mathbb{1}_{H_{n-1}}] \quad (\text{A.22b})$$

$$+ \mathbb{E}[(\gamma_n \chi_n + \gamma_n^2 \psi_n^2) \mathbb{1}_{H_{n-1}}]. \quad (\text{A.22c})$$

However, since  $H_{n-1}$  and  $M_{n-1}$  are both  $\mathcal{F}_n$ -measurable, we have the following estimates:

- (1) For the noise term in (A.22a), we have:

$$\mathbb{E}[M_{n-1} \xi_n \mathbb{1}_{H_{n-1}}] = \mathbb{E}[M_{n-1} \mathbb{1}_{H_{n-1}} \mathbb{E}[\xi_n | \mathcal{F}_n]] = 0. \quad (\text{A.23})$$

- (2) The term (A.22b) is where the reduction to  $H_{n-1}$  kicks in; indeed, we have:

$$\begin{aligned} \mathbb{E}[\xi_n^2 \mathbb{1}_{H_{n-1}}] &= \mathbb{E}[\mathbb{1}_{H_{n-1}} \mathbb{E}[|\langle \pi_n - \pi^*, U_n \rangle|^2 | \mathcal{F}_n]] \\ &\leq \mathbb{E}[\mathbb{1}_{H_{n-1}} \|\pi_n - \pi^*\|^2 \mathbb{E}[\|U_n\|^2 | \mathcal{F}_n]] \quad \# \text{ by Cauchy-Schwarz} \\ &\leq \|\Pi\|^2 \sigma_n^2. \end{aligned} \quad (\text{A.24})$$

- (3) Finally, for the term (A.22c), we have:

$$\mathbb{E}[\psi_n^2 \mathbb{1}_{H_{n-1}}] \leq \frac{3}{2} [G^2 + B_n^2 + \sigma_n^2] \quad (\text{A.25})$$

where we used the bound  $\|v(\pi)\| \leq G$ . Likewise,  $\chi_n \mathbb{1}_{H_{n-1}} \leq \|\Pi\| B_n$ , so

$$(\text{A.22c}) \leq \gamma_n \|\Pi\| B_n + \frac{3}{2} \gamma_n^2 (G^2 + B_n^2 + \sigma_n^2) \quad (\text{A.26})$$

Thus, putting together all of the above, we obtain:

$$\mathbb{E}[(W_n - W_{n-1}) \mathbb{1}_{H_{n-1}}] \leq \gamma_n \|\Pi\| B_n + \gamma_n^2 \|\Pi\|^2 \sigma_n^2 + \frac{3}{2} \gamma_n^2 (G^2 + B_n^2 + \sigma_n^2) \quad (\text{A.27})$$

Going back to (A.20), we have  $W_{n-1} > a$  if  $\tilde{H}_{n-1}$  occurs, so the last term becomes

$$\mathbb{E}[W_{n-1} \mathbb{1}_{\tilde{H}_{n-1}}] \geq a \mathbb{E}[\mathbb{1}_{\tilde{H}_{n-1}}] = a \mathbb{P}(\tilde{H}_{n-1}). \quad (\text{A.28})$$

Our claim then follows by combining Eqs. (A.20), (A.25), (A.26) and (A.28).  $\blacksquare$

**A.3. Extraction of a convergent subsequence.** Our next step is to show that any realization  $\pi_n$  of (PG) that is contained in  $\mathcal{B}$  admits a subsequence  $\pi_{n_k}$  converging to  $\pi^*$ .

**Proposition A.2.** *Let  $\pi^*$  be a stable Nash policy, and let  $\pi_n$  be the sequence of play generated by (PG) with step-size and policy gradient estimates such that  $p + \ell_b > 1$  and  $p - \ell_\sigma > 1/2$  as per (8). Then  $\pi_n$  admits a subsequence  $\pi_{n_k}$  that converges to  $\pi^*$  with probability 1 on the event  $\mathcal{E} = \bigcap_n \mathcal{E}_n = \{\pi_n \in \mathcal{B} \text{ for all } n = 1, 2, \dots\}$ .*

*Proof.* Let  $\mathcal{Q} = \{\pi_n \in \mathcal{B} \text{ for all } n\} \cap \{\liminf_n \|\pi_n - \pi^*\| > 0\}$  denote the event that  $\pi_n$  is contained in  $\mathcal{B}$  but the sequence  $\pi_n$  does not admit a subsequence converging to  $\pi^*$ . We will show that  $\mathbb{P}(\mathcal{Q}) = 0$ .

Indeed, assume ad absurdum that  $\mathbb{P}(\mathcal{Q}) > 0$ . Hence, with probability 1 on  $\mathcal{Q}$ , there exists some positive constant  $c > 0$  (again, possibly random) such that  $\langle v(\pi_n), \pi_n - \pi^* \rangle \leq -c < 0$  for all  $n$ . Thus, going back to (A.1), we get

$$D_{n+1} \leq D_n - \gamma_n c + \gamma_n \xi_n + \gamma_n \chi_n + \gamma_n^2 \psi_n^2, \quad (\text{A.29})$$

so if we let  $\tau_n = \sum_{k=1}^n \gamma_k$  and telescope the above, we obtain the bound

$$D_{n+1} \leq D_1 - \tau_n \left[ c - \frac{M_n}{\tau_n} - \frac{S_n}{\tau_n} \right] \quad (\text{A.30})$$

with  $\xi_n, \chi_n$  and  $\psi_n$  given by (A.3), and  $M_n = \sum_{k=1}^n \gamma_k \xi_k$ ,  $S_n = \sum_{k=1}^n [\gamma_k \chi_k + \gamma_k^2 \psi_k^2]$  defined as in (C.9). Also, (7) readily gives

$$\sum_{n=1}^{\infty} \mathbb{E}[\gamma_n^2 \xi_n^2 | \mathcal{F}_n] \leq \sum_{n=1}^{\infty} \gamma_n^2 \mathbb{E}[\|\pi_n - \pi^*\|^2 \|U_n\|^2 | \mathcal{F}_n] \leq \|\Pi\|^2 \sum_{n=1}^{\infty} \gamma_n^2 \sigma_n^2 < \infty \quad (\text{A.31})$$

so, by the strong law of large numbers for martingale difference sequences [20, Theorem 2.18], we conclude that  $M_n/\tau_n$  converges to 0 with probability 1. In a similar vein, for the submartingale  $S_n$  we have

$$\mathbb{E}[S_n] = \sum_{k=1}^n \gamma_k \chi_k + \sum_{k=1}^n \gamma_k^2 \mathbb{E}[\psi_k^2] \leq \|\Pi\| \sum_{k=1}^n \gamma_k B_k + \frac{3}{2} \sum_{k=1}^n \gamma_k^2 [G^2 + B_k^2 + \sigma_k^2], \quad (\text{A.32})$$

so, by (7) and the stated conditions for the method's step-size and bias/noise parameters, it follows that  $S_n$  is bounded in  $L^1$ . Therefore, by Doob's submartingale convergence theorem [20, Theorem 2.5], we further deduce that  $S_n$  converges with probability 1 to some (finite) random variable  $S_\infty$ .

Going back to (A.30) and letting  $n \rightarrow \infty$ , the above shows that  $D_n \rightarrow -\infty$  with probability 1 on  $\mathcal{Q}$ . Since  $D$  is nonnegative by construction and  $\mathbb{P}(\mathcal{Q}) > 0$  by assumption, we obtain a contradiction and our proof is complete.  $\blacksquare$

**A.4. Convergence of the energy values.** Our last auxiliary result concerns the convergence of the values of the dual energy function  $D$ . We encode this as follows.

**Proposition A.3.** *If (PG) is run with assumptions as in Proposition A.1, there exists a finite random variable  $D_\infty$  such that*

$$\mathbb{P}(D_n \rightarrow D_\infty \text{ as } n \rightarrow \infty \mid \pi_n \in \mathcal{B} \text{ for all } n) = 1. \quad (\text{A.33})$$

*Proof.* Let  $\mathcal{E}_n = \{\pi_k \in \mathcal{B} \text{ for all } k = 1, 2, \dots, n\}$  be defined as in (A.7), and let  $\tilde{D}_n = \mathbb{1}_{\mathcal{E}_n} D_n$ . Then, by the energy inequality (A.2) and the fact that  $\mathcal{E}_{n+1} \subseteq \mathcal{E}_n$ , we get

$$\begin{aligned} \tilde{D}_{n+1} &= \mathbb{1}_{\mathcal{E}_{n+1}} D_{n+1} \leq \mathbb{1}_{\mathcal{E}_n} D_{n+1} \\ &\leq \mathbb{1}_{\mathcal{E}_n} D_n + \mathbb{1}_{\mathcal{E}_n} \gamma_n \langle v(\pi_n), \pi_n - \pi^* \rangle + (\gamma_n \xi_n + \gamma_n \chi_n + \gamma_n^2 \psi_n^2) \mathbb{1}_{\mathcal{E}_n} \\ &\leq \tilde{D}_n + \gamma_n \mathbb{1}_{\mathcal{E}_n} \xi_n + (\gamma_n \chi_n + \gamma_n^2 \psi_n^2) \mathbb{1}_{\mathcal{E}_n}, \end{aligned} \quad (\text{A.34})$$

where we used the fact that  $\langle v(\pi_k), \pi_k - \pi^* \rangle \leq 0$  for all  $k = 1, 2, \dots, n$  if  $\mathcal{E}_n$  occurs. Since  $\mathcal{E}_n$  is  $\mathcal{F}_n$ -measurable, conditioning on  $\mathcal{F}_n$  and taking expectations yields

$$\begin{aligned} \mathbb{E}[\tilde{D}_{n+1} | \mathcal{F}_n] &\leq \tilde{D}_n + \gamma_n \mathbb{1}_{\mathcal{E}_n} \mathbb{E}[\xi_n | \mathcal{F}_n] + \mathbb{1}_{\mathcal{E}_n} \gamma_n \chi_n + \mathbb{1}_{\mathcal{E}_n} \mathbb{E}[\gamma_n^2 \psi_n^2 | \mathcal{F}_n] \\ &\leq \tilde{D}_n + \gamma_n \|\Pi\| B_n + \gamma_n \chi_n + \mathbb{E}[\gamma_n^2 \psi_n^2 | \mathcal{F}_n] \end{aligned}$$

$$\leq \tilde{D}_n + \gamma_n \|\Pi\| B_n + \frac{3}{2} [G^2 + B_n^2 + \sigma_n^2]. \quad (\text{A.35})$$

By our step-size assumptions, we have  $\sum_n \gamma_n^2 (1 + B_n^2 + \sigma_n^2) < \infty$  and  $\sum_n \gamma_n B_n < \infty$ , which means that  $\tilde{D}_n$  is an almost supermartingale with almost surely summable increments, i.e.,

$$\sum_{n=1}^{\infty} \left[ \mathbb{E}[\tilde{D}_{n+1} | \mathcal{F}_n] - \tilde{D}_n \right] < \infty \quad \text{with probability 1} \quad (\text{A.36})$$

Therefore, by Gladyshev's lemma [44, p. 49], we conclude that  $\tilde{D}_n$  converges almost surely to some (finite) random variable  $D_\infty$ . Since  $\mathbf{1}_{\mathcal{E}_n} = 1$  for all  $n$  if and only if  $\pi_n \in \mathcal{B}$  for all  $n$ , we conclude that  $\mathbb{P}(D_n \text{ converges} | \pi_n \in \mathcal{B} \text{ for all } n) = \mathbb{P}(\tilde{D}_n \text{ converges}) = 1$ , and our claim follows.  $\blacksquare$

**A.5. Putting everything together.** We are now in a position to prove [Theorem 1](#) and [Corollary 1](#).

*Proof of Theorem 1.* Let  $\mathcal{E} = \bigcap_n \mathcal{E}_n = \{\pi_n \in \mathcal{B} \text{ for all } n\}$  denote the event that  $\pi_n$  lies in  $\mathcal{B}$  for all  $n$ . By [Proposition A.1](#), if  $\pi_1$  is initialized within the neighborhood  $\mathcal{U}$  defined in [\(A.8\)](#), we have  $\mathbb{P}(\mathcal{E} | \pi_1 \in \mathcal{U}) \geq 1 - a$ , noting also that the neighborhood  $\mathcal{U}$  is independent of the required confidence level  $a$ . Then, by [Propositions A.2](#) and [A.3](#), it follows that a)  $\liminf_n \|\pi_n - \pi^*\| = 0$ ; and b)  $D_n$  converges, both events occurring with probability 1 on the set  $\mathcal{E} \cap \{\pi_1 \in \mathcal{U}\}$ . We thus conclude that  $\lim_{n \rightarrow \infty} D_n = 0$  and hence

$$\begin{aligned} \mathbb{P}(\pi_n \rightarrow \pi^* | \pi_1 \in \mathcal{U}) &\geq \mathbb{P}(\mathcal{E} \cap \{\pi_n \rightarrow \pi^*\} | \pi_1 \in \mathcal{U}) \\ &= \mathbb{P}(\pi_n \rightarrow \pi^* | \pi_1 \in \mathcal{U}, \mathcal{E}) \times \mathbb{P}(\mathcal{E} | \pi_1 \in \mathcal{U}) \geq 1 - \delta, \end{aligned}$$

and our proof is complete.  $\blacksquare$

*Proof of Corollary 1.* For [Models 1](#) and [2](#), taking  $\ell_b = \infty, \ell_\sigma = 0$ , we obtain  $p > 1/2$ . Since we have that  $\sum_{n=1}^{\infty} \gamma_n = \infty$ , we get that  $p \leq 1$ , i.e.,  $p \in (1/2, 1]$ .

For [Model 3](#), we have that  $B_n = \mathcal{O}(\varepsilon_n)$  and  $\sigma_n = \mathcal{O}(1/\sqrt{\varepsilon_n})$ , i.e.,  $\ell_b = r$  and  $\ell_\sigma = r/2$ . Now, since  $p \leq 1, p + \ell_b > 1$  and  $p - \ell_\sigma > 1/2$ , we obtain that  $p \in (2/3, 1]$  and  $(1 - p)/2 < r/2 < p - 1/2$ .  $\blacksquare$

## APPENDIX B. RATE OF CONVERGENCE TO SECOND-ORDER STATIONARY POLICIES

We now proceed with the proof of [Theorem 2](#), which we again restate below for convenience:

**Theorem 2.** *Let  $\pi^*$  be a second-order stationary policy, let  $\mathcal{B}$  be a neighborhood of  $\pi^*$  such that [\(3a\)](#) holds on  $\mathcal{B}$ , and let  $\pi_n$  be the sequence of play generated by [\(PG\)](#) with step-size  $\gamma_n = \gamma/(n + m)^p, p \in (1/2, 1]$ , and policy gradient estimates such that  $p + \ell_b > 1$  and  $p - \ell_\sigma > 1/2$  as per [\(8\)](#). Then:*

- (1) *There exists a neighborhood  $\mathcal{U}$  of  $\pi^*$  in  $\Pi$  such that, for any confidence level  $\delta > 0$ , the event*

$$\mathcal{E} = \{\pi_n \in \mathcal{B} \text{ for all } n = 1, 2, \dots\} \quad (17)$$

*occurs with probability  $\mathbb{P}(\mathcal{E} | \pi_1 \in \mathcal{U}) \geq 1 - \delta$  if  $m$  is large enough relative to  $\delta$ .*

- (2) *The sequence  $\pi_n$  converges to  $\pi^*$  with probability 1 on  $\mathcal{E}$ ; in particular, we have*

$$\mathbb{P}(\pi_n \text{ converges to } \pi^* | \pi_1 \in \mathcal{U}) \geq 1 - \delta \quad (18)$$

*if  $m$  is large relative to  $\delta$ . Moreover, conditioned on  $\mathcal{E}$  and taking  $q = \min\{\ell_b, p - 2\ell_\sigma\}$ , we have*

$$\mathbb{E}[\|\pi_n - \pi^*\|^2 | \mathcal{E}] = \begin{cases} \mathcal{O}(1/n^{2\mu\gamma}) & \text{if } p = 1 \text{ and } 2\mu\gamma < q, \\ \mathcal{O}(1/n^q) & \text{otherwise.} \end{cases} \quad (19)$$

*Proof.* We will follow an approach similar to [Theorem 1](#) for the first part of the theorem. More precisely, let  $\mathcal{B} = \{\pi \in \Pi : \|\pi - \pi^*\| \leq \varrho\}$  be a ball of radius  $\varrho$  centered at  $\pi^*$  in  $\Pi$  such that (SOS) holds for all  $\pi \in \mathcal{B}$ . Then, for all  $\pi \in \mathcal{B} \setminus \{\pi^*\}$ , we have  $\langle v(\pi), \pi - \pi^* \rangle \leq -\mu \|\pi - \pi^*\| < 0$  by [Proposition 1](#). Hence, defining the events  $\mathcal{E}_n$  and  $H_n$  as in [Eq. \(A.7\)](#), and assuming that  $\pi_1$  is initialized in a ball of radius  $\sqrt{2a}$  centered at  $\pi^*$ , viz.

$$\mathcal{U} = \{\pi \in \Pi : D(\pi) \leq a\} = \{\pi \in \Pi : \|\pi - \pi^*\|^2/2 \leq a\}. \quad (\text{B.1})$$

then, by [Lemma A.2](#) and [Proposition A.1](#), we readily obtain that

$$\mathbb{P}(H_n \mid \pi_1 \in \mathcal{U}) \geq 1 - \delta \quad \text{for all } n = 1, 2, \dots \quad (\text{B.2})$$

which implies that

$$\mathbb{P}(\mathcal{E} \mid \pi_1 \in \mathcal{U}) \geq 1 - \delta \quad (\text{B.3})$$

if  $m$  is large enough relative to  $\delta$ .

For the second part, constraining [Eq. \(A.2\)](#) on the event  $\mathcal{E}_n$ , we get:

$$\begin{aligned} D_{n+1} \mathbf{1}_{\mathcal{E}_n} &\leq D_n \mathbf{1}_{\mathcal{E}_n} + \gamma_n \langle v(\pi_n), \pi_n - \pi^* \rangle \mathbf{1}_{\mathcal{E}_n} + \mathbf{1}_{\mathcal{E}_n} (\gamma_n \xi_n + \gamma_n \chi_n + \gamma_n^2 \psi_n^2) \\ &\leq (1 - 2\mu\gamma_n) D_n \mathbf{1}_{\mathcal{E}_n} + \mathbf{1}_{\mathcal{E}_n} (\gamma_n \xi_n + \gamma_n \chi_n + \gamma_n^2 \psi_n^2) \end{aligned} \quad (\text{B.4})$$

where the last inequality comes from (SOS). Therefore, taking expectations, we obtain:

$$\begin{aligned} \mathbb{E}[D_{n+1} \mathbf{1}_{\mathcal{E}_n}] &\leq (1 - 2\mu\gamma_n) \mathbb{E}[D_n \mathbf{1}_{\mathcal{E}_n}] + \mathbb{E}[\mathbf{1}_{\mathcal{E}_n} (\gamma_n \xi_n + \gamma_n \chi_n + \gamma_n^2 \psi_n^2)] \\ &\leq (1 - 2\mu\gamma_n) \mathbb{E}[D_n \mathbf{1}_{\mathcal{E}_n}] + \gamma_n \mathbb{E}[\mathbf{1}_{\mathcal{E}_n} \xi_n] + \gamma_n \mathbb{E}[\mathbf{1}_{\mathcal{E}_n} \chi_n] + \gamma_n^2 \mathbb{E}[\mathbf{1}_{\mathcal{E}_n} \psi_n^2] \\ &= (1 - 2\mu\gamma_n) \mathbb{E}[D_n \mathbf{1}_{\mathcal{E}_n}] + \gamma_n \mathbb{E}[\mathbf{1}_{\mathcal{E}_n} \chi_n] + \gamma_n^2 \mathbb{E}[\mathbf{1}_{\mathcal{E}_n} \psi_n^2] \\ &\leq (1 - 2\mu\gamma_n) \mathbb{E}[D_n \mathbf{1}_{\mathcal{E}_n}] + \|\Pi\| \mathbb{P}(\mathcal{E}_n) \gamma_n B_n + \mathbb{P}(\mathcal{E}_n) (G\gamma_n^2 + 3\gamma_n^2 \sigma_n^2 + 3\gamma_n^2 B_n^2) \end{aligned} \quad (\text{B.5})$$

where the equality in the third line comes from the fact that

$$\mathbb{E}[\mathbf{1}_{\mathcal{E}_n} \xi_n] = \mathbb{E}[\mathbb{E}[\xi_n \mathbf{1}_{\mathcal{E}_n} \mid \mathcal{F}_n]] = \mathbb{E}[\mathbf{1}_{\mathcal{E}_n} \mathbb{E}[\xi_n \mid \mathcal{F}_n]] = 0. \quad (\text{B.6})$$

Now, since  $\mathbf{1}_{\mathcal{E}_{n+1}} \leq \mathbf{1}_{\mathcal{E}_n}$ , we further have

$$\mathbb{E}[D_{n+1} \mathbf{1}_{\mathcal{E}_{n+1}}] \leq \mathbb{E}[D_{n+1} \mathbf{1}_{\mathcal{E}_n}] \quad (\text{B.7})$$

and hence, setting  $\bar{D}_n := \mathbb{E}[D_n \mathbf{1}_{\mathcal{E}_n}]$ , we get

$$\begin{aligned} \bar{D}_{n+1} &\leq (1 - 2\mu\gamma_n) \bar{D}_n + \|\Pi\| \mathbb{P}(\mathcal{E}_n) \gamma_n B_n + \mathbb{P}(\mathcal{E}_n) (G\gamma_n^2 + 3\gamma_n^2 \sigma_n^2 + 3\gamma_n^2 B_n^2) \\ &\leq (1 - 2\mu\gamma_n) \bar{D}_n + \|\Pi\| \gamma_n B_n + G\gamma_n^2 + 3\gamma_n^2 \sigma_n^2 + 3\gamma_n^2 B_n^2. \end{aligned} \quad (\text{B.8})$$

Therefore, taking  $\gamma_n, B_n, \sigma_n$  as per the statement of the theorem and noting that the terms  $\gamma_n^2$  and  $\gamma_n^2 B_n^2$  are respectively dominated by the terms  $\gamma_n^2 \sigma_n^2$  and  $\gamma_n B_n$ , we obtain

$$\begin{aligned} \bar{D}_{n+1} &\leq \left(1 - \frac{2\mu\gamma}{(n+m)^p}\right) \bar{D}_n + \frac{C_1}{(n+m)^{p+\ell_b}} + \frac{C_2}{(n+m)^{2p-2\ell_\sigma}} \\ &\leq \left(1 - \frac{2\mu\gamma}{(n+m)^p}\right) \bar{D}_n + \frac{C_1 + C_2}{(n+m)^{p+q}} \end{aligned} \quad (\text{B.9})$$

for some  $C_1, C_2 > 0$ , where  $q = \min\{\ell_b, p - 2\ell_\sigma\}$ , as per the theorem's statement. Therefore, by a straightforward modification of Chung's lemma [[10](#), Lemmas 2&3], [[44](#), p. 45], we get

$$\bar{D}_n = \begin{cases} \mathcal{O}(1/n^{2\mu\gamma}) & \text{if } p = 1 \text{ and } 2\mu\gamma < q, \\ \mathcal{O}(1/n^q) & \text{otherwise.} \end{cases} \quad (\text{B.10})$$

Accordingly, letting  $n \rightarrow \infty$  and recalling that  $\mathbb{E}[D_n \mathbf{1}_{\mathcal{E}}] \leq \mathbb{E}[D_n \mathbf{1}_{\mathcal{E}_n}] = \bar{D}_n$

$$\lim_{n \rightarrow \infty} \mathbb{E}[D_n \mathbf{1}_{\mathcal{E}}] = 0. \quad (\text{B.11})$$

Then, by Fatou's lemma [18], we obtain

$$0 \leq \mathbb{E}[\liminf_{n \rightarrow \infty} D_n \mathbf{1}_{\mathcal{E}}] \leq \liminf_{n \rightarrow \infty} \mathbb{E}[D_n \mathbf{1}_{\mathcal{E}}] = 0, \quad (\text{B.12})$$

which readily shows that  $\mathbb{E}[\liminf_{n \rightarrow \infty} D_n \mathbf{1}_{\mathcal{E}}] = 0$ . Finally, since  $\liminf_{n \rightarrow \infty} D_n \mathbf{1}_{\mathcal{E}} \geq 0$  (a.s.) and  $\mathbb{E}[\liminf_{n \rightarrow \infty} D_n \mathbf{1}_{\mathcal{E}}] = 0$ , we get that

$$\liminf_{n \rightarrow \infty} D_n \mathbf{1}_{\mathcal{E}} = 0 \quad \text{with probability 1.} \quad (\text{B.13})$$

Therefore, there exists a subsequence  $D_{n_k}$  that converges to 0 with probability 1 on the event  $\mathcal{E}$ , i.e.,  $\pi_{n_k}$  converges to  $\pi^*$ . Hence, invoking [Proposition A.3](#), we further deduce that  $D_n$  converges to some  $D_\infty$  with probability 1 on  $\mathcal{E}$ , and thus, we obtain that  $\lim_{n \rightarrow \infty} D_n = 0$  on  $\mathcal{E}$ . We thus get

$$\begin{aligned} \mathbb{P}(\pi_n \rightarrow \pi^* \mid \pi_1 \in \mathcal{U}) &\geq \mathbb{P}(\mathcal{E} \cap \{\pi_n \rightarrow \pi^*\} \mid \pi_1 \in \mathcal{U}) \\ &= \mathbb{P}(\pi_n \rightarrow \pi^* \mid \pi_1 \in \mathcal{U}, \mathcal{E}) \times \mathbb{P}(\mathcal{E} \mid \pi_1 \in \mathcal{U}) \geq 1 - \delta, \end{aligned} \quad (\text{B.14})$$

as claimed.

For the last part of the theorem, note that

$$\begin{aligned} \bar{D}_n &= \mathbb{E}[D_n \mathbf{1}_{\mathcal{E}_n}] \geq \mathbb{E}[D_n \mathbf{1}_{\mathcal{E}}] = \mathbb{E}[\mathbb{E}[D_n \mid \sigma(\mathcal{E})] \mathbf{1}_{\mathcal{E}}] \\ &= \mathbb{E}[\mathbb{E}[D_n \mid \mathcal{E}] \mathbf{1}_{\mathcal{E}}] \\ &= \mathbb{E}[D_n \mid \mathcal{E}] \mathbb{P}(\mathcal{E}) \\ &= \mathbb{E}[D_n \mid \mathcal{E}] \mathbb{P}(\mathcal{E}) \end{aligned} \quad (\text{B.15})$$

where we used the fact that  $\mathbb{E}[D_n \mid \sigma(\mathcal{E})] \mathbf{1}_{\mathcal{E}} = \mathbb{E}[D_n \mid \mathcal{E}] \mathbf{1}_{\mathcal{E}}$ . We thus conclude that

$$\mathbb{E}[\|\pi_n - \pi^*\|^2 \mid \mathcal{E}] = 2 \mathbb{E}[D_n \mid \mathcal{E}] \leq \frac{2}{\mathbb{P}(\mathcal{E})} \bar{D}_n \leq \frac{2}{1 - \delta} \bar{D}_n \quad (\text{B.16})$$

and hence

$$\mathbb{E}[\|\pi_n - \pi^*\|^2 \mid \mathcal{E}] = \begin{cases} \mathcal{O}(1/n^{2\mu\gamma}) & \text{if } p = 1 \text{ and } 2\mu\gamma < q, \\ \mathcal{O}(1/n^q) & \text{otherwise.} \end{cases} \quad \blacksquare$$

*Proof of [Corollary 2](#).* For [Models 1](#) and [2](#), taking  $\ell_b = \infty, \ell_\sigma = 0$  we readily get that  $q = p$  and  $p > 1/2$ . Since we require that  $\sum_{n=1}^{\infty} \gamma_n = \infty$ , we obtain that  $p \in (1/2, 1]$ . Hence, for  $p = 1$  and  $2\mu\gamma > 1$  we obtain  $\mathcal{O}(1/n)$  rate of convergence.

For [Model 3](#), we have that  $B_n = \mathcal{O}(\varepsilon_n)$  and  $\sigma_n = \mathcal{O}(1/\sqrt{\varepsilon_n})$ , i.e.,  $\ell_b = p/2$  and  $\ell_\sigma = p/4$ , and, hence, we readily get that  $q = p/2$ . Now, since  $p \leq 1, p + \ell_b > 1$  and  $p - \ell_\sigma > 1/2$ , we obtain that  $p \in (2/3, 1]$ . Hence, for  $p = 1$  and  $\mu\gamma > 1$ , we obtain  $\mathcal{O}(1/\sqrt{n})$  rate of convergence.  $\blacksquare$

We conclude this appendix with a detailed statement and proof of the fact that, when run with perfect policy gradients (i.e., as per [Model 1](#)), the sequence of play generated by (PG) achieves a geometric convergence rate to Nash policies satisfying (SOS). The precise result is as follows:

**Proposition B.1.** *Let  $\pi^*$  be a second-order stationary policy, let  $\mathcal{B}$  be a neighborhood of  $\pi^*$  such that (3a) holds on  $\mathcal{B}$ , and let  $\pi_n$  be the sequence of play generated by (PG) with a sufficiently small constant step-size  $\gamma > 0$  and perfect policy gradients as per [Model 1](#). Then, there exists a neighborhood  $\mathcal{U}$  of  $\pi^*$  in  $\Pi$  and some  $\rho > 0$  such that*

$$\|\pi_n - \pi^*\| = \mathcal{O}(\exp(-\rho n)) \quad \text{whenever } \pi_1 \in \mathcal{U}. \quad (\text{B.17})$$

*Proof.* The crucial part of the proof is the observation that, in the case of [Model 1](#), the energy inequality [\(A.1\)](#) of [Lemma A.1](#) may be written in the sharper form:

$$D_{n+1} \leq D_n + \gamma \langle v(\pi_n) - v(\pi^*), \pi_n - \pi^* \rangle + \frac{1}{2} \gamma^2 \|v(\pi_n) - v(\pi^*)\|^2. \quad (\text{B.18})$$

To see this, consider the development

$$\begin{aligned} \|\pi_{n+1} - \pi^*\|^2 &= \|\pi_{n+1} - \pi_n + \pi_n - \pi^*\|^2 \\ &= \|\pi_n - \pi^*\|^2 + 2\langle \pi_{n+1} - \pi_n, \pi_n - \pi^* \rangle + \|\pi_{n+1} - \pi_n\|^2 \\ &= \|\pi_n - \pi^*\|^2 + 2\langle \pi_{n+1} - \pi_n, \pi_{n+1} - \pi^* \rangle - \|\pi_{n+1} - \pi_n\|^2 \\ &\leq \|\pi_n - \pi^*\|^2 + 2\gamma \langle v(\pi_n), \pi_{n+1} - \pi^* \rangle - \|\pi_{n+1} - \pi_n\|^2, \\ &\leq \|\pi_n - \pi^*\|^2 + 2\gamma \langle v(\pi_n) - v(\pi^*), \pi_{n+1} - \pi^* \rangle - \|\pi_{n+1} - \pi_n\|^2, \end{aligned} \quad (\text{B.19})$$

where the last line follows from [\(FOS\)](#) and, in the penultimate step, we used the fact that  $\pi^* \in \Pi$  so  $\langle \pi_{n+1} - (\pi_n + \gamma v(\pi_n)), \pi_{n+1} - \pi^* \rangle \leq 0$ . In addition, by Young's inequality, we have

$$\begin{aligned} \langle v(\pi_n) - v(\pi^*), \pi_{n+1} - \pi^* \rangle &= \langle v(\pi_n) - v(\pi^*), \pi_n - \pi^* \rangle + \langle v(\pi_n) - v(\pi^*), \pi_{n+1} - \pi_n \rangle \\ &\leq \langle v(\pi_n) - v(\pi^*), \pi_n - \pi^* \rangle + \frac{\gamma}{2} \|v(\pi_n) - v(\pi^*)\|^2 + \frac{1}{2\gamma} \|\pi_{n+1} - \pi_n\|^2 \end{aligned} \quad (\text{B.20})$$

so [\(B.18\)](#) follows by substituting [\(B.20\)](#) in [\(B.19\)](#) and simplifying.

Now, since  $v$  is  $G$ -Lipschitz (cf. [Lemma D.7](#) in [Appendix D](#)), we have  $\|v(\pi_n) - v(\pi^*)\| \leq G\|\pi_n - \pi^*\|$ , so the energy inequality [\(B.18\)](#) becomes

$$D_{n+1} \leq D_n + \gamma \langle v(\pi_n) - v(\pi^*), \pi_n - \pi^* \rangle + \gamma^2 G^2 D_n. \quad (\text{B.21})$$

However, if  $\pi_n \in \mathcal{B}$ , [Proposition 1](#) further yields

$$\langle v(\pi_n) - v(\pi^*), \pi_n - \pi^* \rangle \leq -\mu \|\pi_n - \pi^*\|^2 \quad (\text{B.22})$$

so

$$D_{n+1} \leq D_n - \mu\gamma \|\pi_n - \pi^*\|^2 + \gamma^2 G^2 D_n = (1 - 2\mu\gamma + \gamma^2 G^2) D_n. \quad (\text{B.23})$$

Thus, if  $\gamma < 2\mu/G^2$  and  $\pi_1$  is initialized in a ball centered at  $\pi^*$  and contained within  $\mathcal{B}$ , our assertion follows from a straightforward induction argument.  $\blacksquare$

## APPENDIX C. RATE OF CONVERGENCE TO STRICT NASH POLICIES

**C.1. Structural preliminaries.** To prove [Theorem 3](#), we will first require some notions describing the geometry of  $\Pi$  near  $\pi^*$ . Referring to [\[46\]](#) for a full treatment, we have:

**Definition 3.** Let  $\mathcal{C}$  be a convex set and let  $x \in \mathcal{C}$ . Then the tangent cone  $\text{TC}_{\mathcal{C}}(x)$  is defined as the set of all rays emanating from  $x$  and intersecting  $\mathcal{C}$  at at least one other point different from  $x$ . The *polar cone*  $\text{PC}_{\mathcal{C}}(x)$  to  $\mathcal{C}$  at  $x$  is then defined  $\text{PC}_{\mathcal{C}}(x) = \{y : \langle y, z \rangle \leq 0 \text{ for all } z \in \text{TC}_{\mathcal{C}}(x)\}$ , where  $y$  belong in the dual space of the vector space in which  $\mathcal{C}$  is defined.

With these general definitions in hand, we proceed to characterize some further projections of Euclidean projections on  $\Pi$  that will play an important role in the sequel. For notational simplicity, we suppress the player and state indices in the statement and proof of the next lemma.

**Lemma C.1.**  $x = \text{proj}(y)$  if and only if there exist  $\mu \in \mathbb{R}$  and  $\nu_\alpha \in \mathbb{R}_+$  such that, for all  $\alpha \in \mathcal{A}$ , we have  $y_\alpha = x_\alpha + \mu - \nu_\alpha$  with  $\nu_\alpha \geq 0$  and  $x_\alpha \nu_\alpha = 0$ .

*Proof.* Recall that  $\text{proj}(y) = \arg \min_{x \in \Delta(\mathcal{A})} \|y - x\|^2$ . Our result then follows by applying the KKT conditions to this optimization problem and noting that, since the constraints are affine, the KKT conditions are sufficient for optimality. Our Lagrangian is

$$\mathcal{L}(x, \mu, \nu) = \sum_{\alpha \in \mathcal{A}} \frac{1}{2} (y_\alpha - x_\alpha)^2 - \mu \left( \sum_{\alpha \in \mathcal{A}} x_\alpha - 1 \right) + \sum_{\alpha \in \mathcal{A}} \nu_\alpha x_\alpha \quad (\text{C.1})$$

where the set of constraints (i) of the statement of the lemma are the stationarity constraints, which in our case are  $\nabla \mathcal{L}(x, \mu, \nu) = 0 \Leftrightarrow \nabla \left( \sum_{\alpha \in \mathcal{A}} \frac{1}{2} (y_\alpha - x_\alpha)^2 \right) = \mu \nabla \left( \sum_{\alpha \in \mathcal{A}} x_\alpha - 1 \right) - \sum_{\alpha \in \mathcal{A}} \nu_\alpha \nabla x_\alpha$ , while the set of constraints (ii) of the statement of the lemmas are the complementary slackness constraints. Note that complementary slackness implies  $\nu_\alpha > 0$  whenever  $\alpha \notin \text{supp}(x)$ , so our proof is complete. ■

Our next result is a concrete consequence of [Proposition 1](#) which will be very useful in establishing the stability estimates required for the proof of [Theorem 3](#).

**Lemma C.2.** *Let  $\pi^* = (\alpha_{i,s}^*)_{i \in \mathcal{N}, s \in \mathcal{S}}$  be a strict Nash policy. Then there exists a neighborhood  $\mathcal{U}$  of  $\pi^*$  and constants  $c_{i,s}$  such that for each player  $i \in \mathcal{N}$  and state  $s \in \mathcal{S}$ , we have:*

$$v_{i\alpha_{i,s}^*}(\pi) - v_{i\alpha_{i,s}}(\pi) \geq c_{i,s} \quad \text{for all } \pi \in \mathcal{U} \text{ and } \alpha_i \neq \alpha_i^*, \alpha_i \in \mathcal{A}_i. \quad (\text{C.2})$$

*Proof.* Our claim is a consequence of the definition of strict Nash policies. Specifically, from [Proposition 1](#) we have

$$\langle v(\pi^*), z \rangle < 0 \quad \text{for all } z \in \text{TC}(\pi^*), z \neq 0 \quad (\text{C.3})$$

Let  $z = e_{i,\alpha_{i,s}} - e_{i,\alpha_{i,s}^*}$ , then we get that

$$v_{i\alpha_{i,s}^*}(\pi^*) - v_{i\alpha_{i,s}}(\pi^*) > 0 \quad (\text{C.4})$$

where  $e_{i,\alpha_{i,s}}$  is the vector that has one only in the index and zero anywhere else. By continuity there exists a neighborhood  $\mathcal{U} \subseteq \mathcal{X}$  and  $c_{i,s} > 0$  for each player  $i \in \mathcal{N}$  such that

$$v_{i\alpha_{i,s}^*}(\pi) - v_{i\alpha_{i,s}}(\pi) \geq c_{i,s} \quad \text{for all } \pi \in \mathcal{U} \quad \blacksquare$$

Our final result is intimately tied to the lazy projection step in [\(LPG\)](#), and quantifies the relation between initializations in  $\prod_i \mathbb{R}^{\mathcal{A}_i \times \mathcal{S}}$  and  $\Pi$ .

**Lemma C.3.** *Let  $\pi^* = (\alpha_{i,s}^*)_{i \in \mathcal{N}, s \in \mathcal{S}}$ , be a deterministic policy. For each agent  $i \in \mathcal{N}$  and each state  $s \in \mathcal{S}$ , let  $y_{i,\alpha_{i,s}} - y_{i,\alpha_{i,s}^*}$  be the difference of the aggregated gradients between the strategy of the equilibrium and any other strategy  $\alpha_i^* \neq \alpha_i \in \mathcal{A}_i$ . Then for any  $\varepsilon > 0$  such that  $\mathcal{U}_\varepsilon = \{\pi : \pi_{i,\alpha_{i,s}^*} \geq 1 - \varepsilon \text{ for all } i \in \mathcal{N} \text{ and } s \in \mathcal{S}\}$ , there exist  $M_{i,\varepsilon,s}$  such that if  $\mathcal{W}_{i,s} = \{y \in \mathbb{R}^{\mathcal{A}_i} : y_{i,\alpha_{i,s}} - y_{i,\alpha_{i,s}^*} < -M_{i,\varepsilon,s}\}$  then  $\prod_{i \in \mathcal{N}, s \in \mathcal{S}} \text{proj}_{\Pi_i}(\mathcal{W}_{i,s}) \subseteq \mathcal{U}_\varepsilon$ .*

*Proof.* Consider an arbitrary player  $i \in \mathcal{N}$ , a state  $s \in \mathcal{S}$ , and let  $\mathcal{W}_i(M_{i,\varepsilon,s})$  be an open set as defined in the statement of the lemma. For notational simplicity, we will drop the index  $s$ . We will show that any  $M_{i,\varepsilon} > 1 - \frac{\varepsilon}{|\mathcal{A}_i|} > 0$  satisfies our claim. By using [Lemma C.1](#) for a  $y_i \in \mathcal{W}_i(M_{i,\varepsilon})$  with  $\pi_i = \text{proj}(y_i)$  we have that

$$\begin{aligned} y_i \alpha_i^* - y_i \alpha_i &> M_{i,\varepsilon} \\ \pi_i \alpha_i^* - \pi_i \alpha_i - (\nu_{\alpha_i^*} - \nu_{\alpha_i}) &> M_{i,\varepsilon} \end{aligned} \quad (\text{C.5})$$

with  $\nu_{\alpha_i} \geq 0$  and  $\pi_i \alpha_i = 0$  whenever  $\nu_{\alpha_i} > 0$ . Notice that since  $M_{i,\varepsilon} > 1 - \frac{\varepsilon}{A_i}$  we have that  $\pi_i \alpha_i^* > \pi_i \alpha_i + 1 - \frac{\varepsilon}{A_i} + (\nu_{\alpha_i^*} - \nu_{\alpha_i})$  or

$$\pi_i \alpha_i < \pi_i \alpha_i^* - 1 + \frac{\varepsilon}{A_i} - (\nu_{\alpha_i^*} - \nu_{\alpha_i}) < \frac{\varepsilon}{A_i} \quad (\text{C.6})$$

Hence, by summing over all strategies of player  $i$  we get the desired result. ■

**C.2. Proof of the main theorem.** We are now in a position to prove our main result on the rate of convergence towards strict Nash policies. For ease of reference, we restate [Theorem 3](#) below.

**Theorem 3.** *Let  $\pi_n$  be the sequence of play under (LPG) with step-size and policy gradient estimates such that  $p + \ell_b > 1$  and  $p - \ell_\sigma > 1/2$  as per (8). If  $\pi^*$  is a deterministic Nash policy, there exists an unbounded open set  $\mathcal{W} \subseteq \prod_i \mathbb{R}^{\mathcal{A}_i \times \mathcal{S}}$  of initializations such that, for any  $\delta > 0$ , we have*

$$\mathbb{P}(\pi_n \text{ converges to } \pi^* \mid y_1 \in \mathcal{W}) \geq 1 - \delta, \quad (20)$$

provided that  $\gamma > 0$  is small enough. Moreover, conditioned on this event,  $\pi_n$  converges to  $\pi^*$  at a finite number of iterations, i.e., there exists some  $n_0$  such that  $\pi_n = \pi^*$  for all  $n \geq n_0$ .

*Proof of Theorem 3.* We start by fixing a confidence level  $\delta > 0$  and all the parameters of the algorithm, such that all the assumptions stated in the theorem are satisfied and. We will prove that for each agent  $i \in \mathcal{N}$ ,  $s \in \mathcal{S}$  there exist  $M_{1,i,s} > 0$ ,  $\mathcal{W}_{1,i,s} = \{y \in \mathbb{R}^{\mathcal{A}_i} : y_{i,\alpha_i} - y_{i,\alpha_i^*} < -M_{1,i,s} \text{ for all } \alpha_i \in \mathcal{A}_i, \alpha_i \neq \alpha_i^*\}$ , such that if  $y_1 \in \mathcal{W}_1 := \prod_{i \in \mathcal{N}, s \in \mathcal{S}} \mathcal{W}_{1,i,s}$  then the agents' sequence of play, converge to the deterministic Nash policy, in finite number of iterations.

To simplify the notation, we will drop the indices  $s$  and  $i$  referring to the states and agents, accordingly, and we will focus on a specific agent and a specific state. From [Lemma C.3](#), [Lemma C.2](#) we have that there exist constants  $c, M$ , neighborhood  $\mathcal{U}_c = \{\pi \in \Pi : \|\pi - \pi^*\| \leq \beta\}$  and open set  $\mathcal{W}_M$  such that

$$\begin{aligned} v_{\alpha^*}(\pi) - v_\alpha(\pi) &\geq c && \text{for all } \alpha \neq \alpha^*, \alpha \in \mathcal{A} \text{ and } \pi \in \mathcal{U}_c \\ y_{\alpha^*} - y_\alpha &> M_c && \text{for all } \alpha \neq \alpha^*, \alpha \in \mathcal{A} \text{ and } \pi = \text{proj}(y) \in \mathcal{U}_c \end{aligned} \quad (C.7)$$

The first step is to prove that for an appropriate initialization for  $y_1$ , we have  $y_n \in \mathcal{W}(M_c)$  for all  $n = 1, 2, \dots$ , with probability at least  $1 - \delta$ . Assume that  $y_k \in \mathcal{W}(M_c)$  for all  $k = 1, \dots, n$ ; then for the differences of the scores at a round  $n + 1$  between any  $\alpha \in \mathcal{A}$  and the equilibrium strategy  $\alpha^*$ , we have

$$\begin{aligned} y_{\alpha,n+1} - y_{\alpha^*,n+1} &= y_{\alpha,n} - y_{\alpha^*,n} + (\hat{v}_{\alpha,n} - \hat{v}_{\alpha^*,n}) \\ &= y_{\alpha,1} - y_{\alpha^*,1} + \sum_{k=1}^n \gamma_k [(v_{\alpha,k} - v_{\alpha^*,k}) + (U_{\alpha,k} - U_{\alpha^*,k}) + (b_{\alpha,k} - b_{\alpha^*,k})] \\ &\leq -M_1 + \sum_{k=1}^n \gamma_k [(v_{\alpha,k} - v_{\alpha^*,k}) + (U_{\alpha,k} - U_{\alpha^*,k}) + (b_{\alpha,k} - b_{\alpha^*,k})] \\ &\leq -M_1 - c \sum_{k=1}^n \gamma_k + \sum_{k=1}^n \gamma_k [(U_{\alpha,k} - U_{\alpha^*,k}) + (b_{\alpha,k} - b_{\alpha^*,k})] \\ &\leq -M_1 - c \sum_{k=1}^n \gamma_k + \sum_{k=1}^n \gamma_k [\xi_k + \chi_k] \end{aligned} \quad (C.8)$$

where  $\xi_k = (U_{\alpha,k} - U_{\alpha^*,k})$  and  $\chi_k = 2\|b_k\|$ . Now, similarly to the proofs of [Theorems 1](#) and [2](#) we will proceed to control the aggregate error terms

$$R_n = \sum_{k=1}^n \gamma_k \xi_k \quad \text{and} \quad S_n = \sum_{k=1}^n \gamma_k \chi_k. \quad (C.9)$$

Since  $\mathbb{E}[\xi_n \mid \mathcal{F}_n] = 0$ , we have  $\mathbb{E}[R_n \mid \mathcal{F}_n] = R_{n-1}$ , so  $R_n$  is a martingale; likewise,  $\mathbb{E}[S_n \mid \mathcal{F}_n] \geq S_{n-1}$ , so  $S_n$  is a sub-martingale. Furthermore from (7) we have:

$$\text{I. } \mathbb{E}[\xi_n^2] \leq \mathbb{E}[\|U_n\|^2] \leq \mathbb{E}[\mathbb{E}[\|U_n\|^2 \mid \mathcal{F}_n]] \leq \sigma_n^2$$

$$\text{II. } \mathbb{E}[\chi_n] = 2 \mathbb{E}[|b_n|] \leq \mathbb{E}[\mathbb{E}[|b_n| \mid \mathcal{F}_n]] \leq B_n$$

Moreover, for any  $\eta_1 > 0$ , we get by Doob's Maximal Inequality:

$$\mathbb{P}\left(\sup_{1 \leq k \leq n} R_k \geq \eta_1\right) \leq \frac{\mathbb{E}[R_n^2]}{\eta_1^2} \stackrel{(a)}{\leq} \frac{\sum_{k=1}^n \gamma_k^2 \mathbb{E}[\xi_k^2]}{\eta_1^2} \stackrel{(I)}{\leq} \frac{\sum_{k=1}^n \gamma_k^2 \sigma_k^2}{\eta_1^2} \quad (\text{C.10})$$

where (a) comes from the fact that  $\mathbb{E}[\xi_i \xi_j] = 0$  for  $i \neq j$ . Since  $\gamma_n = \gamma/n^p$ ,  $\sigma_n = \mathcal{O}(n^{\ell_\sigma})$  and  $p - \ell_\sigma > 1/2$ , there exists  $\gamma_1$  sufficiently small such that if  $\gamma \leq \gamma_1$  then

$$\sum_{k=1}^{\infty} \gamma_k^2 \sigma_k^2 < \frac{\delta \eta_1^2}{2} \quad (\text{C.11})$$

and so we automatically get that

$$\mathbb{P}\left(\sup_{1 \leq k \leq n} R_k \geq \eta_1\right) \leq \frac{\delta}{2} \quad (\text{C.12})$$

Furthermore, notice that the term  $\{S_n\}_{n \in \mathbb{N}}$  is a sub-martingale, since  $\mathbb{E}[|S_n| \mid \mathcal{F}_n] < \infty$  and  $\mathbb{E}[S_{n+1} \mid \mathcal{F}_n] > S_n$ , for all  $n$ . As before, using Doob's Maximal Inequality, we get for any  $\eta_2 > 0$ :

$$\mathbb{P}\left(\sup_{1 \leq k \leq n} S_k \geq \eta_2\right) \leq \frac{\mathbb{E}[S_n]}{\eta_2} = \frac{\sum_{k=1}^n \gamma_k \mathbb{E}[\chi_k]}{\eta_2} \leq \frac{2 \sum_{k=1}^n \gamma_k B_k}{\eta_2} \quad (\text{C.13})$$

So, since  $p + \ell_b > 1$  there exists  $\gamma_2$  sufficiently small such that if  $\gamma \leq \gamma_2$  then  $\sum_{k=1}^n \gamma_k B_k \leq \frac{\eta_2 \delta}{4}$  which immediately implies that

$$\mathbb{P}\left(\sup_{1 \leq k \leq n} S_k \geq \eta_2\right) \leq \frac{\delta}{2} \quad (\text{C.14})$$

By choosing  $\gamma \leq \min\{\gamma_1, \gamma_2\}$  we get that

$$\mathbb{P}\left(\sup_{1 \leq k \leq n} R_n + S_n \leq M_c\right) \geq 1 - \delta. \quad (\text{C.15})$$

Notice now that by choosing  $M_1 > M_c + \eta_1 + \eta_2$ , from (C.8) we have that with probability at least  $1 - \delta$ ,  $y_{\alpha, n+1} - y_{\alpha^*, n+1} < -M_c$ , which implies that  $\pi_{n+1} \in \mathcal{U}_c$ .

Defining the sequences of "good" events  $\{\mathcal{E}_n\}_{n \in \mathbb{N}}$  and  $\{\mathcal{E}'_n\}_{n \in \mathbb{N}}$  as  $\mathcal{E}_n := \{\pi_k \in \mathcal{U}_c \text{ for all } k = 1, 2, \dots, n\}$  and  $\mathcal{E}'_n := \{\sup_{1 \leq k \leq n} R_k + S_k \leq \eta_1 + \eta_2\}$ , accordingly, we get that  $\mathcal{E}'_n \subseteq \mathcal{E}_n$  for all  $n$ . Because  $\mathbb{P}(\mathcal{E}'_n) \geq 1 - \delta$ , we get that  $\mathbb{P}(\mathcal{E}_n) \geq 1 - \delta$  and since  $\{\mathcal{E}_n\}_{n \in \mathbb{N}}$  is a decreasing sequence converging to  $\mathcal{E} := \{\pi_n \in \mathcal{U}_c, \forall n \in \mathbb{N}\}$ , we obtain  $\mathbb{P}(\mathcal{E}) \geq 1 - \delta$ . i.e.,

$$\mathbb{P}(\pi_n \in \mathcal{U}_c, \forall n \mid y_1 \in \mathcal{W}_1) \geq 1 - \delta \quad (\text{C.16})$$

Notice that the above conclusions immediately imply convergence in finite time. More specifically, constrained to the event  $\mathcal{E}$  with probability at least  $1 - \delta$ , from Eq. (C.8) we have

$$y_{\alpha, n+1} - y_{\alpha^*, n+1} \leq -M_c - c \sum_{k=1}^n \gamma_k \quad (\text{C.17})$$

for all  $n = 1, 2, \dots$ . Assume ad absurdum that there exists at least one strategy  $\alpha \neq \alpha^*$ ,  $\alpha \in \mathcal{A}$  such that  $\limsup_{n \rightarrow \infty} \pi_{\alpha, n} \geq \varepsilon > 0$ . for all sufficiently large  $n$ . Recall also that for  $\pi \in \mathcal{U}_c$ , it holds that  $\pi_{\alpha^*} > 0$  by construction. Using Lemma C.1 we get

$$y_{\alpha, n+1} - y_{\alpha^*, n+1} = \pi_{\alpha, n+1} - \pi_{\alpha^*, n+1} \leq -M_c - c \sum_{k=1}^n \gamma_k \quad (\text{C.18})$$

Notice that the LHS of this inequality is bounded, while the RHS goes to  $-\infty$ , which is a contradiction. Thus, with probability at least  $1 - \delta$ ,  $\pi_n \rightarrow \pi^*$  as  $n \rightarrow \infty$ .

We can rewrite the previous inequality as

$$\pi_{\alpha, n+1} \leq 1 - M_c - c \sum_{k=1}^n \gamma_k \quad \text{for all } \alpha^* \neq \alpha \in \mathcal{A} \quad (\text{C.19})$$

Now aggregating over all strategies, on the previous inequality, we get that

$$\|\pi_{n+1} - \pi^*\|_1 = 2(1 - \pi_{\alpha^*, n+1}) \leq 2 \sum_{\alpha^* \neq \alpha \in \mathcal{A}} (1 - M_c - c \sum_{k=1}^n \gamma_k) \quad (\text{C.20})$$

Thus, once  $\sum_{k=1}^n \gamma_k$  becomes at least  $(1 - M_c)/c$ , which occurs in finite time, the convergence is implied.  $\blacksquare$

*Proof of Corollary 3.* For Models 1 and 2, taking  $\ell_b = \infty, \ell_\sigma = 0$  we readily get that  $p > 1/2$ . Since we require that  $\sum_{n=1}^{\infty} \gamma_n = \infty$ , we obtain that  $p \in (1/2, 1]$ .

For Model 3, we have that  $B_n = \mathcal{O}(\varepsilon_n)$  and  $\sigma_n = \mathcal{O}(1/\sqrt{\varepsilon_n})$ , i.e.,  $\ell_b = r$  and  $\ell_\sigma = r/2$ . Now, since  $p \leq 1, p + \ell_b > 1$  and  $p - \ell_\sigma > 1/2$ , we obtain that  $p \in (2/3, 1]$ .  $\blacksquare$

#### APPENDIX D. STRUCTURAL PROPERTIES OF POLICY GRADIENT METHODS

In this appendix we will establish the required properties for the gradient of the players' value function. More precisely, we prove the following intermediate results:

- In Lemma D.1 we prove that in the random stopping episodic framework visitation the notion of discounted state visitation distribution is well-defined.
- In Lemma 1, we prove the conversion lemma, a standard lemma that connects a sample by visitation distribution and a random trajectory.
- In Lemma D.4, we establish different versions of Policy Gradient theorem via  $Q$ -value function for the random stopping episodic framework.
- In Lemmas D.5 and D.7, we establish the boundedness and the Lipschitz smoothness of policy gradient vector field, i.e.,  $v(\pi) = (v_i(\pi))_{i \in \mathcal{N}}$  where  $v_i(\pi) = \nabla_{\pi_i} V_{i,s}(\pi)$

For a policy profile  $\pi \in \Pi$  and an arbitrary initial state distribution  $s_0 \sim \rho$ , let's recall the definition of discounted state visitation measure/distribution as

$$\tilde{d}_\rho^\pi(s) = \mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} \mathbf{1}\{s_t = s\} \mid s_0 \sim \rho \right], \quad d_\rho^\pi(s) := \tilde{d}_\rho^\pi(s) / Z_\rho^\pi$$

To begin with, we prove formally that the above definition is well-posed for the random stopping episodic framework described above, i.e.,  $\tilde{d}_\rho^\pi(s) < \infty$ , so  $Z_\rho^\pi := \sum_{s \in \mathcal{S}} \tilde{d}_\rho^\pi(s)$  is well-defined.

**Lemma D.1.** *For any  $s \in \mathcal{S}$ ,  $\tilde{d}_\rho^\pi(s) < \infty$  and  $Z_\rho^\pi \leq \frac{1}{\zeta}$ .*

*Proof.* For the sake of the proof, we define a new state  $s_f$ , indicating that the game has stopped. In other words, we have that  $\mathbb{P}(s_f \mid s, \alpha) = \zeta_{s,\alpha} \geq \zeta > 0$  for all  $\alpha \in \mathcal{A}, s \in \mathcal{S}$ . Hence, for  $s \in \mathcal{S}$  we obtain:

$$\begin{aligned} \tilde{d}_\rho^\pi(s) &= \mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} \mathbf{1}\{s_t = s\} \mid s_0 \sim \rho \right] \\ &= \mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{\infty} \mathbf{1}\{s_t = s, s_i \neq s_f, 1 \leq i \leq t\} \mid s_0 \sim \rho \right] \end{aligned}$$

$$\begin{aligned}
&\leq \sum_{s \in \mathcal{S}} \tilde{d}_\rho^\pi(s) = \mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{\infty} \mathbb{1}\{s_i \neq s_f, 1 \leq i \leq t\} | s_0 \sim \rho \right] \\
&= \sum_{t=0}^{\infty} \mathbb{P}(s_i \neq s_f, 1 \leq i \leq t | s_0 \sim \rho) \\
&= \sum_{t=0}^{\infty} \prod_{i=1}^t \mathbb{P}(s_i \neq s_f | s_0 \sim \rho, s_j \neq s_f, 1 \leq j \leq i-1) \\
&\leq \sum_{t=0}^{\infty} (1 - \zeta)^t \leq \frac{1}{\zeta} < \infty. \quad \blacksquare
\end{aligned}$$

**Lemma 1.** [Conversion Lemma] *For an arbitrary state-action function  $f: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , a policy profile  $\pi$  and an initial state distribution  $s_0 \sim \rho$ , we have*

$$\mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} f(s_t, \alpha_t) \right] = Z_\rho^\pi \mathbb{E}_{s \sim d_\rho^\pi} \mathbb{E}_{\alpha \sim \pi(\cdot|s)} [f(s, \alpha)] \quad (2)$$

*Proof.*

$$\begin{aligned}
\mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} f(s_t, \alpha_t) \right] &= \sum_{t=0}^{\infty} \sum_{s \in \mathcal{S}} \sum_{\alpha \in \mathcal{A}} \mathbb{E}_{\tau \sim \text{MDP}} [\mathbb{1}\{t \leq T(\tau), s_t = s, \alpha_t = \alpha\} f(s, \alpha)] \\
&= \sum_{s \in \mathcal{S}} \sum_{t=0}^{\infty} \sum_{\alpha \in \mathcal{A}} \mathbb{P}^\pi(s = s_t | s_0 \sim \rho) \pi(\alpha | s) f(s, \alpha) \\
&= \sum_{s \in \mathcal{S}} \sum_{t=0}^{\infty} \mathbb{P}^\pi(s = s_t | s_0 \sim \rho) \sum_{\alpha \in \mathcal{A}} \pi(\alpha | s) f(s, \alpha) \\
&= \sum_{s \in \mathcal{S}} \tilde{d}_\rho^\pi(s) \mathbb{E}_{\alpha \sim \pi(\cdot|s)} [f(s, \alpha)] \\
&= Z_\rho^\pi \mathbb{E}_{s \sim d_\rho^\pi} \mathbb{E}_{\alpha \sim \pi(\cdot|s)} [f(s, \alpha)] \quad (D.1)
\end{aligned}$$

where  $Z_\rho^\pi := \mathbb{E}_{s \sim \text{Unif}(\mathcal{S})} [\tilde{d}_\rho^\pi(s)] \cdot |\mathcal{S}|$  is well-defined by [Lemma D.1](#).  $\blacksquare$

A compact reformulation the aforementioned lemma is via the matrix representation of the discounted visitation distribution:

**Lemma D.2** (Conversion Lemma (Matrix form)). *For an arbitrary state-action function  $f: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  and a policy profile  $\pi$ , we have*

$$\mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} f(s_t, \alpha_t) | \alpha_0 = \alpha, s_0 = s \right] = e_{s,\alpha}^\top \mathcal{T}(\pi) f \quad (D.2)$$

where  $\mathcal{T}$  is a discounted visitation distribution (action-state)-matrix under policy profile  $\pi$  i.e.,

$$[\mathcal{T}(\pi)] \underbrace{(\alpha, s)}_{\text{Row Index}} \rightarrow \underbrace{(\alpha', s')}_{\text{Column Index}} = \sum_{t=0}^{\infty} \mathbb{P}^\pi(s_t = s', \alpha_t = \alpha' | s_0 = s, \alpha_0 = \alpha)$$

*Proof.* By definition we have

$$\begin{aligned}
e_{s,\alpha}^\top \mathcal{T}(\pi) f &= \langle e_{s,\alpha}^\top \mathcal{T}(\pi), f \rangle \\
&= \sum_{s' \in \mathcal{S}} \sum_{\alpha' \in \mathcal{A}} (e_{s,\alpha}^\top \mathcal{T}(\pi))_{(s', \alpha')} \cdot f(s', \alpha') \\
&= \sum_{s' \in \mathcal{S}} \sum_{\alpha' \in \mathcal{A}} e_{s,\alpha}^\top \mathcal{T}(\pi) e_{s', \alpha'} \cdot f(s', \alpha')
\end{aligned}$$

$$\begin{aligned}
&= \sum_{s' \in \mathcal{S}} \sum_{\alpha' \in \mathcal{A}} \sum_{t=0}^{\infty} \mathbb{P}^{\pi}(s_t = s', \alpha_t = \alpha' | s_0 = s, \alpha_0 = \alpha) \cdot f(s', \alpha') \\
&= \sum_{t=0}^{\infty} \sum_{s' \in \mathcal{S}} \sum_{\alpha' \in \mathcal{A}} \mathbb{E}_{\tau \sim \text{MDP}} [\mathbb{1}\{t \leq T(\tau), s'_t = s, \alpha'_t = \alpha, \} f(s, \alpha) | s_0 = s, \alpha_0 = \alpha] \\
&= \mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} f(s_t, \alpha_t) | \alpha_0 = \alpha, s_0 = s \right]. \quad \blacksquare
\end{aligned}$$

*Remark 4.* Notice that  $\mathcal{T}$  is a well-defined matrix. Indeed, let's us define  $\mathcal{P}(\pi)$  as the state-action one step transition matrix:

$$[\mathcal{P}(\pi)] \underbrace{(\alpha, s)}_{\text{Row Index}} \rightarrow \underbrace{(\alpha', s')}_{\text{Column Index}} = \mathbb{P}^{\pi}(s_1 = s', \alpha_1 = \alpha' | s_0 = s, \alpha_0 = \alpha) = \pi(\alpha' | s') P(s' | s, \alpha). \quad (\text{D.3})$$

Notice that  $\mathcal{P}(\pi)$  is a substochastic matrix and therefore  $\text{spectral}(\mathcal{P}(\pi)) < 1$  or equivalently  $(I - \mathcal{P}(\pi))^{-1}$  is invertible. Thus using Neumann series we have that  $(I - \mathcal{P}(\pi))^{-1} = \sum_{t=0}^{\infty} \mathcal{P}(\pi)^t$ . By induction, a folklore probabilistic-graph theoretic fact, we can show that  $\sum_{t=0}^{\infty} \mathcal{P}(\pi)^t = \mathcal{T}(\pi)$ .

In order to analyze the gradient of MARL policy gradient methods, we will introduce the notions  $Q, A$  and their per-player averages that are useful in the MDP analysis.

**Definition 4.** For a state  $s \in \mathcal{S}$ , a policy  $\pi$  and  $\alpha = (\alpha_1, \dots, \alpha_N) \in \mathcal{A}$ , we define:

(i) The  $Q$ -value function of player  $i$  as:

$$Q_i^{\pi}(s, \alpha) := \mathbb{E}_{\tau \sim \text{MDP}(\pi|s)} \left[ \sum_{t=0}^{T(\tau)} R_i(s_t(\tau), \alpha_t(\tau)) | s_0 = s, \alpha_0 = \alpha \right] \quad (\text{D.4})$$

(ii) The *advantage function* of player  $i$  as:

$$\text{adv}_i^{\pi}(s, \alpha) := Q_i^{\pi}(s, \alpha) - V_{i,s}(\pi) \quad (\text{D.5})$$

We also define  $\bar{Q}_i^{\pi}, \overline{\text{adv}}_i^{\pi}$  to be the averaged for  $i$ -th player single MDP  $Q$ -value and advantage functions:

(i) The averaged  $\bar{Q}_i^{\pi}$ -value function of player  $i$  as:

$$\bar{Q}_i^{\pi}(s, \alpha_i) := \mathbb{E}_{\alpha_{-i} \sim \pi_{-i}(\cdot|s)} [Q_i^{\pi}(s, (\alpha_i; \alpha_{-i}))] \quad (\text{D.6})$$

(ii) The *averaged advantage function*  $\overline{\text{adv}}_i^{\pi}$  of player  $i$  as:

$$\overline{\text{adv}}_i^{\pi}(s, \alpha_i) := \mathbb{E}_{\alpha_{-i} \sim \pi_{-i}(\cdot|s)} [\text{adv}_i^{\pi}(s, (\alpha_i; \alpha_{-i}))], \quad (\text{D.7})$$

By [Remark 4](#), we can rewrite the above notations using  $\mathcal{T}, \mathcal{P}$ .

**Lemma D.3.** For a policy profile  $\pi$ , we have:

$$(1) Q_i^{\pi}(s, \alpha) = e_{s,\alpha}^{\top} \mathcal{T}(\pi) r_i$$

$$(2) \tilde{d}_p^{\pi}(s) = [\sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} \pi(\alpha' | s') e_{s',\alpha'}]^{\top} \mathcal{T}(\pi) \sum_{\alpha \in \mathcal{A}} e_{s,\alpha}$$

*Proof.* We separately have using [Lemma D.3](#) and [Remark 4](#).

(1) For our first claim, a straightforward calculation gives:

$$Q_i^{\pi}(s, \alpha) = \mathbb{E}_{\tau \sim \text{MDP}(\pi|s)} \left[ \sum_{t=0}^{T(\tau)} R_i(s_t(\tau), \alpha_t(\tau)) | s_0 = s, \alpha_0 = \alpha \right] = e_{s,\alpha}^{\top} \mathcal{T}(\pi) R_i \quad (\text{D.8})$$

(2) As for the second part of the lemma, we have:

$$\begin{aligned}
\tilde{d}_\rho^\pi(s) &= \mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} \mathbb{1}\{s_t = s\} \middle| s_0 \sim \rho \right] \\
&= \mathbb{E}_{s' \sim \rho} \mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} \sum_{\alpha \in \mathcal{A}} \mathbb{1}\{s_t = s, \alpha_t = \alpha\} \middle| s_0 = s' \right] \\
&= \mathbb{E}_{s' \sim \rho} \mathbb{E}_{\alpha' \sim \pi(\cdot|s)} \mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} \sum_{\alpha \in \mathcal{A}} \mathbb{1}\{s_t = s, \alpha_t = \alpha\} \middle| s_0 = s', \alpha_0 = \alpha' \right] \\
&= \mathbb{E}_{s' \sim \rho} \mathbb{E}_{\alpha' \sim \pi(\cdot|s)} \left[ e_{s', \alpha'}^\top \mathcal{T}(\pi) \sum_{\alpha \in \mathcal{A}} e_{s, \alpha} \right] \\
&= \left[ \sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} \pi(\alpha' | s') e_{s', \alpha'} \right]^\top \mathcal{T}(\pi) \sum_{\alpha \in \mathcal{A}} e_{s, \alpha}. \quad \blacksquare
\end{aligned}$$

Having defined the above notions, we are ready to provide equivalent forms of  $v(\pi)$  that will permit us to prove its boundedness and smoothness. We start with the following versions of the policy gradient theorem for random stopping setting:

**Lemma D.4.** *For the independent gradient operator  $v(\pi)$  per player the following expressions are equal to  $v_i(\pi)$ :*

$$\begin{aligned}
(1) \quad v_i(\pi) &= \mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} \nabla_i (\log \pi_i(\alpha_{i,t}(\tau) | s_t(\tau))) \bar{Q}_i^\pi(s_t(\tau), \alpha_{i,t}(\tau)) \right] \\
(2) \quad v_i(\pi) &= Z_\rho^\pi \mathbb{E}_{s \sim d_\rho^\pi} \mathbb{E}_{\alpha_i \sim \pi_i(\cdot|s)} \left[ \nabla_i (\log \pi_i(\alpha_i | s)) \bar{Q}_i^\pi(s, \alpha_i) \right] \\
(3) \quad (v_i(\pi))_{\alpha_i^\circ, s^\circ} &= \frac{\partial V_{i,\rho}(\pi)}{\partial \pi_i(\alpha_i^\circ | s^\circ)} = \tilde{d}_\rho^\pi(s^\circ) \bar{Q}_i^\pi(s^\circ, \alpha_i^\circ) = Z_\rho^\pi d_\rho^\pi(s^\circ) \bar{Q}_i^\pi(s^\circ, \alpha_i^\circ)
\end{aligned}$$

*Proof.* Recall first that the independent gradient operator  $v(\pi)$  is given by  $v_i(\pi) = \nabla_i V_{i,\rho}(\pi)$ . Accordingly, we will begin by showing that:

$$\nabla_i (V_{i,\rho}(\pi)) = \mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} \nabla_i (\log \pi_i(\alpha_{i,t}(\tau) | s_t(\tau))) \bar{Q}_i^\pi(s_t(\tau), \alpha_{i,t}(\tau)) \right] \quad (\text{D.9})$$

To that end, we will start with an arbitrary  $s_0$ , and by linearity of  $\nabla_{\pi_i}(\cdot)$  and  $\mathbb{E}_{s_0 \sim \rho}[\cdot]$ , we will obtain the result. Indeed, we have:

$$\begin{aligned}
\nabla_i (V_{i,s_0}(\pi)) &= \nabla_i (\mathbb{E}_\tau [R_i(\tau)]) \\
&= \nabla_i (\mathbb{E}_{\alpha_i \sim \pi_i(\cdot|s_0)} [\bar{Q}_i^\pi(s_0, \alpha_i)]) \\
&= \nabla_i \left( \sum_{\alpha_i \in \mathcal{A}_i} \pi_i(\alpha_i | s_0) \bar{Q}_i^\pi(s_0, \alpha_i) \right) \\
&= \sum_{\alpha_i \in \mathcal{A}_i} \nabla_i (\pi_i(\alpha_i | s_0)) \bar{Q}_i^\pi(s_0, \alpha_i) + \pi_i(\alpha_i | s_0) \nabla_i (\bar{Q}_i^\pi(s_0, \alpha_i)) \\
&= \sum_{\alpha_i \in \mathcal{A}_i} \nabla_i (\log \pi_i(\alpha_i | s_0)) \pi_i(\alpha_i | s_0) \bar{Q}_i^\pi(s_0, \alpha_i) + \pi_i(\alpha_i | s_0) \nabla_i (\bar{Q}_i^\pi(s_0, \alpha_i)) \\
&= \mathbb{E}_{\alpha_i \sim \pi_i(\cdot|s_0)} \left[ \nabla_i (\log \pi_i(\alpha_i | s_0)) \bar{Q}_i^\pi(s_0, \alpha_i) \right] \\
&\quad + \sum_{\alpha_i \in \mathcal{A}_i} \pi_i(\alpha_i | s_0) \nabla_i \left( \mathbb{E}_{\alpha_{-i} \sim \pi_{-i}(\cdot|s_0)} \left[ R_i(s_0, \alpha) + \sum_{s_1 \in \mathcal{S}} P(s_1 | s_0, \alpha) V_{i,s_1}(\pi) \right] \right)
\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_{\alpha_i \sim \pi_i(\cdot|s_0)} \left[ \nabla_i (\log \pi_i(\alpha_i|s_0)) \bar{Q}_i^\pi(s_0, \alpha_i) \right] \\
&\quad + \sum_{\alpha_i \in \mathcal{A}_i} \pi_i(\alpha_i|s_0) \mathbb{E}_{\alpha_{-i} \sim \pi_{-i}(\cdot|s_0)} \left[ \sum_{s_1 \in \mathcal{S}} P(s_1|s_0, \alpha) \nabla_i (V_{i,s_1}(\pi)) \right] \\
&= \mathbb{E}_{\alpha_i \sim \pi_i(\cdot|s_0)} \left[ \nabla_i (\log \pi_i(\alpha_i|s_0)) \bar{Q}_i^\pi(s_0, \alpha_i) \right] \\
&\quad + \mathbb{E}_{\alpha \sim \pi(\cdot|s_0)} \left[ \sum_{s_1 \in \mathcal{S}} P(s_1|s_0, \alpha) \nabla_i (V_{i,s_1}(\pi)) \right] \tag{D.10}
\end{aligned}$$

Thus, we can rewrite it as:

$$\begin{aligned}
\nabla_i (V_{i,s_0}(\pi)) &= \mathbb{E}_{\alpha_i \sim \pi_i(\cdot|s_0)} \left[ \nabla_i (\log \pi_i(\alpha_i|s_0)) \bar{Q}_i^\pi(s_0, \alpha_i) \right] \\
&\quad + \mathbb{E}_{\alpha \sim \pi(\cdot|s_0)} \left[ \sum_{s_1 \in \mathcal{S}} P(s_1|s_0, \alpha) \nabla_i (V_{i,s_1}(\pi)) \right] \\
&= \mathbb{E}_{\tau \sim \text{MDP}(\pi|s_0)} \left[ \nabla_i (\log \pi_i(\alpha_{i,0}(\tau)|s_0)) \bar{Q}_i^\pi(s_0, \alpha_{i,0}(\tau)) \right] \\
&\quad + \mathbb{E}_{\tau \sim \text{MDP}(\pi|s_0)} \left[ \mathbf{1}\{T(\tau) \geq 1\} \nabla_i (V_{i,s_1}(\tau)(\pi)) \right] \\
&= \sum_{t=0}^{\infty} \mathbb{E}_{\tau \sim \text{MDP}(\pi|s_0)} \left[ \mathbf{1}\{t \leq T(\tau)\} \nabla_i (\log \pi_i(\alpha_{i,t}(\tau)|s_t(\tau))) \bar{Q}_i^\pi(s_t(\tau), \alpha_{i,t}(\tau)) \right] \\
&\quad + \mathbb{E}_{\tau \sim \text{MDP}(\pi|s_0)} \left[ \mathbf{1}\{T(\tau) = \infty\} A_\infty \right] \\
&\stackrel{(a)}{=} \mathbb{E}_{\tau \sim \text{MDP}(\pi|s_0)} \left[ \sum_{t=0}^{T(\tau)} \nabla_i (\log \pi_i(\alpha_{i,t}(\tau)|s_t(\tau))) \bar{Q}_i^\pi(s_t(\tau), \alpha_{i,t}(\tau)) \right] \tag{D.11}
\end{aligned}$$

where (a) holds because  $\mathbb{P}(T(\tau) = \infty) = 0$ , and  $A_\infty$  is some limiting quantity.

Hence, we readily obtain:

$$\nabla_i (V_{i,\rho}(\pi)) = \mathbb{E}_{s_0 \sim \rho} [\nabla_i (V_{i,s_0}(\pi))] \tag{D.12}$$

Now by [Lemma 1](#), we further have

$$\begin{aligned}
\nabla_i (V_{i,\rho}(\pi)) &= Z_\rho^\pi \mathbb{E}_{s \sim d_\rho^\pi} \mathbb{E}_{\alpha \sim \pi(\cdot|s)} \left[ \nabla_i (\log \pi_i(\alpha_i|s)) \bar{Q}_i^\pi(s, \alpha_i) \right] \\
&= Z_\rho^\pi \mathbb{E}_{s \sim d_\rho^\pi} \mathbb{E}_{\alpha_i \sim \pi_i(\cdot|s)} \left[ \nabla_i (\log \pi_i(\alpha_i|s)) \bar{Q}_i^\pi(s, \alpha_i) \right] \tag{D.13}
\end{aligned}$$

Decoupling  $\nabla_i$  per a state  $s^\circ$  and action  $\alpha_i^\circ$ , we get

$$\begin{aligned}
\frac{\partial V_{i,\rho}(\pi)}{\partial \pi_i(\alpha_i^\circ|s^\circ)} &= Z_\rho^\pi \mathbb{E}_{s \sim d_\rho^\pi} \mathbb{E}_{\alpha_i \sim \pi_i(\cdot|s)} \left[ \frac{\partial (\log \pi_i(\alpha_i|s))}{\partial \pi_i(\alpha_i^\circ|s^\circ)} \bar{Q}_i^\pi(s, \alpha_i) \right] \\
&= Z_\rho^\pi \mathbb{E}_{s \sim d_\rho^\pi} \mathbb{E}_{\alpha_i \sim \pi_i(\cdot|s)} \left[ \mathbf{1}\{\alpha_i^\circ = \alpha_i, s^\circ = s\} \frac{1}{\pi_i(\alpha_i^\circ|s^\circ)} \bar{Q}_i^\pi(s^\circ, \alpha_i^\circ) \right] \\
&= \sum_{s \in \mathcal{S}} \tilde{d}_\rho^\pi(s) \sum_{\alpha_i \in \mathcal{A}_i} \pi_i(\alpha_i|s) \mathbf{1}\{\alpha_i^\circ = \alpha_i, s^\circ = s\} \frac{1}{\pi_i(\alpha_i^\circ|s^\circ)} \bar{Q}_i^\pi(s^\circ, \alpha_i^\circ) \\
&= \tilde{d}_\rho^\pi(s^\circ) \bar{Q}_i^\pi(s^\circ, \alpha_i^\circ) = Z_\rho^\pi d_\rho^\pi(s^\circ) \bar{Q}_i^\pi(s^\circ, \alpha_i^\circ). \quad \blacksquare
\end{aligned}$$

We are ready to bound the amplitude of the independent player gradient operator:

**Lemma D.5.** *For a given initial state distribution  $\rho$ , the independent player policy gradient operator  $v(\pi)$  is bounded. More precisely,*

$$\|v_i(\pi)\| \leq \frac{\sqrt{A_i}}{\zeta^2} \quad \& \quad \|v(\pi)\| \leq \frac{\sum_{i \in \mathcal{N}} \sqrt{A_i}}{\zeta^2} \tag{D.14}$$

*Proof.* We start by analyzing  $\|v_i(\pi)\|^2$  using [Lemma D.4](#). Specifically, we have:

$$\begin{aligned}
\|v_i(\pi)\|^2 &= \sum_{\alpha_i^\circ, s^\circ \in \mathcal{A}_i, \mathcal{S}} (v_i(\pi)_{\alpha_i^\circ, s^\circ})^2 = \sum_{s^\circ \in \mathcal{S}} \sum_{\alpha_i^\circ \in \mathcal{A}_i} \left( \frac{\partial V_{i,\rho}(\pi)}{\partial \pi_i(\alpha_i^\circ | s^\circ)} \right)^2 \\
&= \sum_{s^\circ \in \mathcal{S}} \sum_{\alpha_i^\circ \in \mathcal{A}_i} (Z_\rho^\pi d_\rho^\pi(s^\circ) \bar{Q}_i^\pi(s^\circ, \alpha_i^\circ))^2 \\
&\leq (Z_\rho^\pi)^2 \max_{\alpha_i^\circ, s^\circ \in \mathcal{A}_i, \mathcal{S}} (\bar{Q}_i^\pi(s^\circ, \alpha_i^\circ))^2 \sum_{s^\circ \in \mathcal{S}} \sum_{\alpha_i^\circ \in \mathcal{A}_i} d_\rho^\pi(s^\circ)^2 \\
&\leq \frac{1}{\zeta^2} \max_{\alpha_i^\circ, s^\circ \in \mathcal{A}_i, \mathcal{S}} (\mathbb{E}_{\alpha_{-i} \sim \pi_{-i}(\cdot | s)} [Q_i^\pi(s^\circ, (\alpha_i^\circ; \alpha_{-i}))])^2 \sum_{s^\circ \in \mathcal{S}} \sum_{\alpha_i^\circ \in \mathcal{A}_i} d_\rho^\pi(s^\circ) \\
&\leq \frac{1}{\zeta^2} \max_{\alpha^\circ, s^\circ \in \mathcal{A}, \mathcal{S}} (Q_i^\pi(s^\circ, \alpha^\circ))^2 \sum_{\alpha_i^\circ \in \mathcal{A}_i} \sum_{s^\circ \in \mathcal{S}} d_\rho^\pi(s^\circ) \\
&\leq \frac{1}{\zeta^2} \max_{\alpha^\circ, s^\circ \in \mathcal{A}, \mathcal{S}} \left( \mathbb{E}_{\tau \sim \text{MDP}(\pi | s)} \left[ \sum_{t=0}^{T(\tau)} R_i(s_t(\tau), \alpha_t(\tau)) | s_0 = s^\circ, \alpha_0 = \alpha^\circ \right] \right)^2 |\mathcal{A}_i| \\
&\leq \frac{A_i}{\zeta^2} \left( \mathbb{E}_{\tau \sim \text{MDP}(\pi | s)} \left[ \sum_{t=0}^{T(\tau)} 1 | s_0 = s^\circ, \alpha_0 = \alpha^\circ \right] \right)^2 \leq \frac{A_i}{\zeta^4} \tag{D.15}
\end{aligned}$$

We thus conclude that

$$\|v_i(\pi)\| \leq \frac{\sqrt{A_i}}{\zeta^2} \quad \text{and} \quad \|v(\pi)\| \leq \frac{\sum_{i \in \mathcal{N}} \sqrt{A_i}}{\zeta^2}. \quad \blacksquare$$

To prove the smoothness of the policy gradient operator, we have first to establish the performance lemma for our setting. Respectively, we get

**Lemma D.6** (Performance lemma). *For any pair of policy profiles  $\pi = (\pi_i, \pi_{-i}), \pi' = (\pi'_i, \pi'_{-i})$ , it holds*

$$V_{i,\rho}(\pi_i, \pi_{-i}) - V_{i,\rho}(\pi'_i, \pi'_{-i}) = \mathbb{E}_{\tau \sim \text{MDP}(\pi | \rho)} \left[ \sum_{t=0}^{T(\tau)} \text{adv}_i^{\pi'_i, \pi'_{-i}}(s_t, \alpha_t) \right] \tag{D.16}$$

where  $\text{MDP}(\pi | \rho)$  signifies that players follow  $\pi$  as policy profile with  $\rho$  as the initial state distribution.

*Proof.* We will initial prove the aforementioned result for an arbitrary deterministic initial state  $s_0 = s$ :

$$\begin{aligned}
V_{i,s}(\pi) - V_{i,s}(\pi') &= \mathbb{E}_{\tau \sim \text{MDP}(\pi | \rho)} \left[ \sum_{t=0}^{T(\tau)} R_i(s_t, \alpha_t) \right] - V_{i,s}(\pi') \\
&= \mathbb{E}_{\tau \sim \text{MDP}(\pi | s)} \left[ \sum_{t=0}^{T(\tau)} (R_i(s_t, \alpha_t) + V_{i,s_t}(\pi') - V_{i,s_t}(\pi')) \right] - V_{i,s}(\pi') \\
&= \mathbb{E}_{\tau \sim \text{MDP}(\pi | s)} \left[ \sum_{t=0}^{T(\tau)} R_i(s_t, \alpha_t) + \sum_{t=0}^{T(\tau)} (V_{i,s_t}(\pi') - V_{i,s}(\pi') - V_{i,s_t}(\pi')) \right]
\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}_{\tau \sim \text{MDP}(\pi|s)} \left[ \sum_{t=0}^{T(\tau)} (R_i(s_t, \alpha_t) + \mathbb{1}\{T(\tau) \geq t+1\} V_{i, s_{t+1}}(\pi')) - V_{i, s_t}(\pi') \right] \\
&= \mathbb{E}_{\tau \sim \text{MDP}(\pi|s)} \left[ \sum_{t=0}^{T(\tau)} (Q_i^{\pi'}(s_t, \alpha_t) - V_{i, s_t}(\pi')) \right] \\
&= \mathbb{E}_{\tau \sim \text{MDP}(\pi|s)} \left[ \sum_{t=0}^{T(\tau)} \text{adv}_i^{\pi'}(s_t, \alpha_t) \right] \tag{D.17}
\end{aligned}$$

where in the last equation we recall the definition of the Advantage function and in the pre-last the equivalent definitions of  $Q_i^\pi(s, \alpha)$

$$\begin{aligned}
Q_i^\pi(s, \alpha) &= \mathbb{E}_{\tau \sim \text{MDP}(\pi|s)} \left[ \sum_{t=0}^{T(\tau)} R_i(s_t(\tau), \alpha_t(\tau)) | s_0 = s, \alpha_0 = \alpha \right] \\
&= R_i(s, \alpha) + \mathbb{E}_{\tau \sim \text{MDP}(\pi|s)} [\mathbb{1}\{T(\tau) \geq 1\} V_{i, s_1}(\pi) | s_0 = s, \alpha_0 = \alpha] \tag{D.18}
\end{aligned}$$

Applying the linearity of  $\mathbb{E}_{s \sim \rho}[\cdot]$ , we get the desired result:

$$V_{i, \rho}(\pi) - V_{i, \rho}(\pi') = \mathbb{E}_{\tau \sim \text{MDP}(\pi|s)} \left[ \sum_{t=0}^{T(\tau)} \text{adv}_i^{\pi'}(s_t, \alpha_t) \right] = Z_\rho^\pi \mathbb{E}_{s \sim d_\rho^\pi} \mathbb{E}_{\alpha \sim \pi(\cdot|s)} [\text{adv}_i^{\pi'}(s, \alpha)] \tag{D.19}$$

where the last expression comes from [Lemma 1](#).  $\blacksquare$

Before closing this section by proving the Lipschitz-smoothness of our operator, we describe a useful observation that will be helpful in the smoothness bounds.

**Proposition D.1.** *For any pair of policy profiles  $\pi = (\pi_i, \pi_{-i})$ ,  $\pi' = (\pi'_i, \pi'_{-i})$  and an arbitrary initial state distribution  $\rho$  and a subset  $\mathcal{M} \subseteq \mathcal{N}$ , it holds that:*

$$\sum_s d_\rho^\pi(s) \sum_{\alpha_{\mathcal{M}}} |(\pi_{\mathcal{M}} - \pi'_{\mathcal{M}})(\alpha_{\mathcal{M}}|s)| \leq \sum_{i \in \mathcal{M}} \sqrt{A_i} \|\pi_i - \pi'_i\| \tag{D.20}$$

where  $\pi_{\mathcal{M}} = (\pi_i)_{i \in \mathcal{M}}$  and  $\alpha_{\mathcal{M}} = (\alpha_i)_{i \in \mathcal{M}}$ , correspondingly.

*Proof.* A series of direct calculations gives:

$$\begin{aligned}
\sum_s d_\rho^\pi(s) \sum_{\alpha_{\mathcal{M}}} |(\pi_{\mathcal{M}} - \pi'_{\mathcal{M}})(\alpha_{\mathcal{M}}|s)| &= 2 \sum_s d_\rho^\pi(s) \frac{1}{2} \|(\pi_{\mathcal{M}} - \pi'_{\mathcal{M}})\|_1 \\
&= 2 \sum_s d_\rho^\pi(s) \frac{1}{2} d_{\text{TV}}(\pi_{\mathcal{M}}(\cdot|s), \pi'_{\mathcal{M}}(\cdot|s)) \\
&\leq 2 \sum_s d_\rho^\pi(s) \sum_{i \in \mathcal{M}} \frac{1}{2} d_{\text{TV}}(\pi_i(\cdot|s), \pi'_i(\cdot|s)) \\
&= \sum_s d_\rho^\pi(s) \sum_{i \in \mathcal{M}} \|(\pi_i(\cdot|s) - \pi'_i(\cdot|s))\|_1 \\
&= \sum_s d_\rho^\pi(s) \sum_{i \in \mathcal{M}} \sqrt{A_i} \|\pi_i - \pi'_i\|_2 \\
&= \sum_{i \in \mathcal{M}} \sqrt{A_i} \|\pi_i - \pi'_i\|_2 \left( \sum_s d_\rho^\pi(s) \right)
\end{aligned}$$

$$= \sum_{i \in \mathcal{M}} \sqrt{A_i} \|\pi_i - \pi'_i\|_2 \quad (\text{D.21})$$

where  $d_{\text{TV}}$  corresponds to the total variation distance, and the first inequality is a consequence of the triangle inequality for  $d_{\text{TV}}$ . ■

**Lemma D.7.** *For a given initial state distribution  $\rho$ , the independent player policy gradient operator  $v(\pi)$  is Lipschitz continuous. More precisely, for any pair of policy profiles  $\pi = (\pi_i, \pi_{-i})$ ,  $\pi' = (\pi'_i, \pi'_{-i})$ , it holds*

$$\|v_i(\pi) - v_i(\pi')\| = \|\nabla_i(V_{i,\rho}(\pi) - \nabla_i(V_{i,\rho}(\pi'))\| \leq \frac{3\sqrt{A_i}}{\zeta^3} \sum_{j=1}^N \sqrt{A_j} \|\pi_j - \pi'_j\| \quad \forall i \in \mathcal{N} \quad (\text{D.22})$$

and consequently,

$$\|v(\pi) - v(\pi')\| \leq \frac{3A}{\zeta^3} \|\pi - \pi'\| \quad (\text{D.23})$$

*Proof.* For the proof, we will follow the approach of Zhang et al. [61] and Agarwal et al. [1]. Our first task is to bound the directional derivative of the  $i$ -th player's value function. To that end, let  $\pi, \pi' \in \Pi$  and  $\text{pert} \in \mathcal{S} \times \mathcal{A}$  such that  $\|\text{pert}\| = 1$ , and consider  $\lambda$ -perturbed policies

$$\begin{aligned} \pi_\lambda^{\mathbb{A}}(\alpha|s) &= (\pi_i + \lambda \text{pert}, \pi_{-i}) \\ \pi_\lambda^{\mathbb{B}}(\alpha|s) &= (\pi'_i + \lambda \text{pert}, \pi'_{-i}) \end{aligned} \quad (\text{D.24})$$

We then have:

$$\begin{aligned} \left| \frac{\partial V_{i,\rho}(\pi_\lambda^{\mathbb{A}})}{\partial \lambda} - \frac{\partial V_{i,\rho}(\pi_\lambda^{\mathbb{B}})}{\partial \lambda} \right| &= \left| \frac{\partial V_{i,\rho}(\pi_\lambda^{\mathbb{A}}) - V_{i,\rho}(\pi_\lambda^{\mathbb{B}})}{\partial \lambda} \right| = \left| \frac{\partial (V_{i,\rho}(\pi_\lambda^{\mathbb{A}}) - V_{i,\rho}(\pi_\lambda^{\mathbb{B}}))}{\partial \lambda} \right| \\ &= \left| \frac{\partial \left( Z_\rho^{\pi_\lambda^{\mathbb{A}}} \mathbb{E}_{s \sim d_\rho^{\pi_\lambda^{\mathbb{A}}}} \mathbb{E}_{\alpha \sim \pi_\lambda^{\mathbb{A}}(\cdot|s)} \left[ \text{adv}_i^{\pi_\lambda^{\mathbb{B}}}(s, \alpha) \right] \right)}{\partial \lambda} \right| \end{aligned} \quad (\text{D.25a})$$

$$= \left| \frac{\partial \left( Z_\rho^{\pi_\lambda^{\mathbb{A}}} \sum_{s, \alpha} d_\rho^{\pi_\lambda^{\mathbb{A}}}(s) (\pi_\lambda^{\mathbb{A}} - \pi_\lambda^{\mathbb{B}})(\alpha|s) \text{adv}_i^{\pi_\lambda^{\mathbb{B}}}(s, \alpha) \right)}{\partial \lambda} \right| \quad (\text{D.25b})$$

$$= \left| \frac{\partial \left( Z_\rho^{\pi_\lambda^{\mathbb{A}}} \sum_{s, \alpha} d_\rho^{\pi_\lambda^{\mathbb{A}}}(s) (\pi_\lambda^{\mathbb{A}} - \pi_\lambda^{\mathbb{B}})(\alpha|s) Q_i^{\pi_\lambda^{\mathbb{B}}}(s, \alpha) \right)}{\partial \lambda} \right| \quad (\text{D.25c})$$

where (D.25a) follows from Lemma D.6 and (D.25b) uses the fact that  $\sum_{\alpha \in \mathcal{A}} \pi(\alpha|s) \text{adv}_i^\pi(s, \alpha) =$ , for all  $s \in \mathcal{S}$  and the last one is derived by the definition  $d_\rho^\pi(s) := \tilde{d}_\rho^\pi(s)/Z_\rho^\pi$ .

By triangular inequality, the linearity of  $\partial$  operator and Lemma D.1, we have:

$$\begin{aligned} \left| \frac{\partial (V_{i,\rho}(\pi_\lambda^{\mathbb{A}}) - V_{i,\rho}(\pi_\lambda^{\mathbb{B}}))}{\partial \lambda} \right|_{\lambda=0} &\leq \left| \sum_{s, \alpha} \frac{\partial \tilde{d}_\rho^{\pi_\lambda^{\mathbb{A}}}(s)}{\partial \lambda} \right|_{\lambda=0} (\pi - \pi')(\alpha|s) Q_i^{\pi'}(s, \alpha) \\ &\quad + Z_\rho^{\pi_\lambda^{\mathbb{A}}} \left| \sum_{s, \alpha} d_\rho^\pi(s) \frac{\partial (\pi_\lambda^{\mathbb{A}} - \pi_\lambda^{\mathbb{B}})(\alpha|s)}{\partial \lambda} \right|_{\lambda=0} Q_i^{\pi'}(s, \alpha) \end{aligned}$$

$$+ Z_\rho^{\pi^\mathbb{A}} \left| \sum_{s,\alpha} d_\rho^\pi(s) (\pi - \pi')(s|\alpha) \frac{\partial Q_i^{\pi^\mathbb{B}}(s,\alpha)}{\partial \lambda} \Big|_{\lambda=0} \right| \quad (\text{D.26})$$

We will bound the following three terms separately:

$$\begin{aligned} \text{Term}_A &= \left| \sum_{s,\alpha} \frac{\partial \tilde{d}_\rho^{\pi^\mathbb{A}}(s)}{\partial \lambda} \Big|_{\lambda=0} (\pi - \pi')(s|\alpha) Q_i^{\pi'}(s,\alpha) \right| \\ \text{Term}_B &= \left| \sum_{s,\alpha} d_\rho^\pi(s) \frac{\partial (\pi^\mathbb{A} - \pi^\mathbb{B})(s|\alpha)}{\partial \lambda} \Big|_{\lambda=0} Q_i^{\pi'}(s,\alpha) \right| \\ \text{Term}_C &= \left| \sum_{s,\alpha} d_\rho^\pi(s) (\pi - \pi')(s|\alpha) \frac{\partial Q_i^{\pi^\mathbb{B}}(s,\alpha)}{\partial \lambda} \Big|_{\lambda=0} \right| \end{aligned} \quad (\text{D.27})$$

For  $\text{Term}_A$ , we will use [Lemma D.3](#) in order to compute compactly the derivative:

$$\begin{aligned} \frac{\partial \tilde{d}_\rho^{\pi^\mathbb{A}}(s)}{\partial \lambda} &= \frac{\partial \left( \left[ \sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} \pi_\lambda^\mathbb{A}(\alpha'|s') e_{s',\alpha'} \right]^\top \mathcal{T}(\pi_\lambda^\mathbb{A}) \sum_{\alpha \in \mathcal{A}} e_{s,\alpha} \right)}{\partial \lambda} \\ &= \left( \left[ \sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} \frac{\partial \pi_\lambda^\mathbb{A}(\alpha'|s')}{\partial \lambda} e_{s',\alpha'} \right]^\top \mathcal{T}(\pi_\lambda^\mathbb{A}) \sum_{\alpha \in \mathcal{A}} e_{s,\alpha} \right) \\ &\quad + \left( \left[ \sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} \pi_\lambda^\mathbb{A}(\alpha'|s') e_{s',\alpha'} \right]^\top \frac{\partial \mathcal{T}(\pi_\lambda^\mathbb{A})}{\partial \lambda} \sum_{\alpha \in \mathcal{A}} e_{s,\alpha} \right) \\ &= \left( \left[ \sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} \text{pert}(\alpha'_i|s') \cdot \pi_{-i}(\alpha'_{-i}|s') e_{s',\alpha'} \right]^\top \mathcal{T}(\pi_\lambda^\mathbb{A}) \sum_{\alpha \in \mathcal{A}} e_{s,\alpha} \right) \\ &\quad + \left( \left[ \sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} \pi_\lambda^\mathbb{A}(\alpha'|s') e_{s',\alpha'} \right]^\top \frac{\partial (I - \mathcal{P}(\pi_\lambda^\mathbb{A}))^{-1}}{\partial \lambda} \sum_{\alpha \in \mathcal{A}} e_{s,\alpha} \right) \\ &= \left( \left[ \sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} \text{pert}(\alpha'_i|s') \cdot \pi_{-i}(\alpha'_{-i}|s') e_{s',\alpha'} \right]^\top \mathcal{T}(\pi_\lambda^\mathbb{A}) \sum_{\alpha \in \mathcal{A}} e_{s,\alpha} \right) \\ &\quad + \left( \left[ \sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} \pi_\lambda^\mathbb{A}(\alpha'|s') e_{s',\alpha'} \right]^\top \left( \mathcal{T}(\pi_\lambda^\mathbb{A}) \frac{\partial \mathcal{P}(\pi_\lambda^\mathbb{A})}{\partial \lambda} \mathcal{T}(\pi_\lambda^\mathbb{A}) \right) \sum_{\alpha \in \mathcal{A}} e_{s,\alpha} \right) \end{aligned} \quad (\text{D.28})$$

Thus for  $\lambda = 0$ , we get

$$\begin{aligned} \frac{\partial \tilde{d}_\rho^{\pi^\mathbb{A}}(s)}{\partial \lambda} \Big|_{\lambda=0} &= \left( \left[ \sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} \text{pert}(\alpha'_i|s') \cdot \pi_{-i}(\alpha'_{-i}|s') e_{s',\alpha'} \right]^\top \mathcal{T}(\pi) \sum_{\alpha \in \mathcal{A}} e_{s,\alpha} \right) \\ &\quad + \left( \left[ \sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} \pi(\alpha'|s') e_{s',\alpha'} \right]^\top \left( \mathcal{T}(\pi) \frac{\partial \mathcal{P}(\pi^\mathbb{A})}{\partial \lambda} \Big|_{\lambda=0} \mathcal{T}(\pi) \right) \sum_{\alpha \in \mathcal{A}} e_{s,\alpha} \right) \end{aligned} \quad (\text{D.29})$$

Notice that  $\left[ \frac{\partial \mathcal{P}(\pi^\mathbb{A})}{\partial \lambda} \Big|_{\lambda=0} \right]_{(s^\circ, \alpha^\circ) \rightarrow (s^*, \alpha^*)} = \text{pert}(\alpha_i^*|s^*) \cdot \pi_{-i}(\alpha_{-i}^*|s^*) P(s^*|s^\circ, \alpha^\circ)$ .

To simplify notation let us call  $\text{aux}_A := [\sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} \text{pert}(\alpha'_i | s') \cdot \pi_{-i}(\alpha'_{-i} | s') e_{s', \alpha'}]$ ,  $\text{aux}_B := [\sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} \pi(\alpha' | s') e_{s', \alpha'}]$  and  $\text{aux}_C(s) := \sum_{\alpha \in \mathcal{A}} e_{s, \alpha}$ .

We thus get:

$$\begin{aligned}
\text{Term}_A &= \left| \sum_{s, \alpha} \frac{\partial \tilde{d}_{\rho}^{\pi_{\lambda}^{\hat{\Lambda}}}(s)}{\partial \lambda} \Big|_{\lambda=0} (\pi' - \pi')(\alpha | s) Q_i^{\pi' s'}(s, \alpha) \right| \\
&= \left| \sum_{s, \alpha} \left( \text{aux}_A^{\top} \mathcal{T}(\pi) \text{aux}_C(s) + \text{aux}_B^{\top} (\mathcal{T}(\pi) \frac{\partial \mathcal{P}(\pi_{\lambda}^{\hat{\Lambda}})}{\partial \lambda} \Big|_{\lambda=0} \mathcal{T}(\pi)) \text{aux}_C(s) \right) (\pi - \pi')(\alpha | s) Q_i^{\pi'}(s, \alpha) \right| \\
&= \left| \left( \text{aux}_A^{\top} \mathcal{T}(\pi) + \text{aux}_B^{\top} (\mathcal{T}(\pi) \frac{\partial \mathcal{P}(\pi_{\lambda}^{\hat{\Lambda}})}{\partial \lambda} \Big|_{\lambda=0} \mathcal{T}(\pi)) \right) \underbrace{\sum_{s, \alpha} (\pi - \pi')(\alpha | s) Q_i^{\pi'}(s, \alpha) \text{aux}_C(s)}_{\text{aux}_D} \right| \\
&\leq \|\text{aux}_A\|_1 \|\mathcal{T}(\pi) \text{aux}_D\|_{\infty} + \|\text{aux}_B\|_1 \|\mathcal{T}(\pi) \frac{\partial \mathcal{P}(\pi_{\lambda}^{\hat{\Lambda}})}{\partial \lambda} \Big|_{\lambda=0} \mathcal{T}(\pi) \text{aux}_D\|_{\infty} \quad (\text{D.30})
\end{aligned}$$

It is easy to see that  $\|\text{aux}_A\|_1 \leq \sqrt{A_i}$ ,  $\|\text{aux}_B\|_1 = 1$ . Indeed,

$$\begin{aligned}
\|\text{aux}_A\|_1 &= \sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} |\text{pert}(\alpha'_i | s')| \cdot \pi_{-i}(\alpha'_{-i} | s') = \sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha'_i \in \mathcal{A}_i} |\text{pert}(\alpha'_i | s')| \\
&= \sum_{s' \in \mathcal{S}} \rho(s') \|\text{pert}_{i|s'}\|_1 \leq \sum_{s' \in \mathcal{S}} \rho(s') \sqrt{A_i} \|\text{pert}_{i|s'}\|_2 \leq \sqrt{A_i} \\
\|\text{aux}_B\|_1 &= \sum_{s' \in \mathcal{S}} \rho(s') \sum_{\alpha' \in \mathcal{A}} \pi(\alpha' | s') = 1 \quad (\text{D.31})
\end{aligned}$$

Additionally by Conversion Lemma in Matrix form (See [Lemma D.2](#)), we have that:

$$\|\mathcal{T}(\pi)x\|_{\infty} = \max_{s, \alpha} |e_{s, \alpha}^{\top} \mathcal{T}(\pi)x| = \max_{s, \alpha} |\mathbb{E}_{\tau \sim \text{MDP}} \left[ \sum_{t=0}^{T(\tau)} x(s_t, \alpha_t) | \alpha_0 = \alpha, s_0 = s \right]| \leq \frac{1}{\zeta} \|x\|_{\infty} \quad (\text{D.32})$$

Similarly, for the matrix  $\frac{\partial \mathcal{P}(\pi_{\lambda}^{\hat{\Lambda}})}{\partial \lambda} \Big|_{\lambda=0}$ , we have that

$$\begin{aligned}
\left\| \frac{\partial \mathcal{P}(\pi_{\lambda}^{\hat{\Lambda}})}{\partial \lambda} \Big|_{\lambda=0} x \right\|_{\infty} &= \max_{s, \alpha} \left| e_{s, \alpha}^{\top} \frac{\partial \mathcal{P}(\pi_{\lambda}^{\hat{\Lambda}})}{\partial \lambda} \Big|_{\lambda=0} x \right| \\
&= \max_{s, \alpha} \left| \sum_{s', \alpha'} \text{pert}(\alpha'_i | s') \cdot \pi_{-i}(\alpha'_{-i} | s') P(s' | s, \alpha) x_{s', \alpha'} \right| \\
&\leq \sum_{s', \alpha'} |\text{pert}(\alpha'_i | s')| \cdot \pi_{-i}(\alpha'_{-i} | s') P(s' | s, \alpha) \\
&\leq \sqrt{A_i} \|\text{pert}_{i|s'}\|_2 \|x\|_{\infty} \leq \sqrt{A_i} \|x\|_{\infty} \quad (\text{D.33})
\end{aligned}$$

since  $\|\text{pert}\|_2 = 1$ . Then, using [\(D.33\)](#) and [\(D.32\)](#) in [\(D.30\)](#) we get that :

$$\begin{aligned}
\text{Term}_A &\leq \frac{\sqrt{A_i}}{\zeta} \|\text{aux}_D\|_{\infty} + \frac{\sqrt{A_i}}{\zeta^2} \|\text{aux}_D\|_{\infty} \\
&\leq \frac{\sqrt{A_i}}{\zeta} \left( 1 + \frac{1}{\zeta} \right) \left\| \sum_{s, \alpha} (\pi - \pi')(\alpha | s) Q_i^{\pi'}(s, \alpha) \text{aux}_C(s) \right\|_{\infty} \\
&\leq \frac{\sqrt{A_i}}{\zeta^2} \left( 1 + \frac{1}{\zeta} \right) \max_s \left| \sum_{\alpha} (\pi - \pi')(\alpha | s) \right| \|\text{aux}_C(s)\|_{\infty}
\end{aligned}$$

$$\begin{aligned}
&\leq \frac{\sqrt{A_i}}{\zeta^2} \left(1 + \frac{1}{\zeta}\right) \sum_{j=1}^N \sqrt{A_i} \|\pi_j - \pi'_j\| \\
&\leq \frac{\sqrt{A_i}}{\zeta^3} \sum_{j=1}^N \sqrt{A_i} \cdot \|\pi_j - \pi'_j\|
\end{aligned} \tag{D.34}$$

where we used above the fact that  $Q$  function is bounded by  $1/\zeta$ ,  $\|\text{pert}\| = 1$  and [Proposition D.1](#) to bound the difference of the policy profiles.

Moving forward, for  $\text{Term}_B$ , we have:

$$\begin{aligned}
\text{Term}_B &= \left| \sum_{s,\alpha} d_\rho^\pi(s) \frac{\partial(\pi_\lambda^\mathbb{A} - \pi_\lambda^\mathbb{B})(\alpha|s)}{\partial\lambda} \Big|_{\lambda=0} Q_i^{\pi'}(s, \alpha) \right| \\
&= \left| \sum_{s,\alpha} d_\rho^\pi(s) \text{pert}(\alpha_i|s) (\pi_{-i} - \pi'_{-i})(\alpha|s) Q_i^{\pi'}(s, \alpha) \right| \\
&\leq \frac{1}{\zeta} \left| \sum_s d_\rho^\pi(s) \sum_{\alpha_i} \text{pert}(\alpha_i|s) \sum_{\alpha_{-i}} (\pi_{-i} - \pi'_{-i})(\alpha|s) \right| \\
&\leq \frac{1}{\zeta} \sum_s \left| d_\rho^\pi(s) \max_{\alpha_i} \sum_{\alpha_i} |\text{pert}(\alpha_i|s)| \sum_{\alpha_{-i}} (\pi_{-i} - \pi'_{-i})(\alpha|s) \right| \\
&\leq \frac{1}{\zeta} \max_s \|\text{pert}_{i|s}\|_1 \sum_s d_\rho^\pi(s) \sum_{\alpha_{-i}} |(\pi_{-i} - \pi'_{-i})(\alpha|s)| \\
&\leq \frac{\sqrt{A_i}}{\zeta} \max_s \|\text{pert}_{i|s}\|_2 \left( \sum_s d_\rho^\pi(s) \sum_{\alpha_{-i}} |(\pi_{-i} - \pi'_{-i})(\alpha|s)| \right) \\
&\leq \frac{\sqrt{A_i}}{\zeta} \sum_{j \in \mathcal{N} \setminus \{i\}} \sqrt{A_i} \|\pi_j - \pi'_j\| \leq \frac{\sqrt{A_i}}{\zeta} \sum_{j=1}^N \sqrt{A_i} \|\pi_j - \pi'_j\|
\end{aligned} \tag{D.35}$$

where we used again the fact that  $Q$  function is bounded by  $1/\zeta$  and [Proposition D.1](#) to bound the difference of the policy profiles.

Finally, for  $\text{Term}_C$ , we have:

$$\begin{aligned}
\text{Term}_C &= \left| \sum_{s,\alpha} d_\rho^\pi(s) (\pi - \pi')(\alpha|s) \frac{\partial Q_i^{\pi^\mathbb{B}}(s, \alpha)}{\partial\lambda} \Big|_{\lambda=0} \right| \\
&\leq \max_{s,\alpha} \left| \frac{\partial Q_i^{\pi^\mathbb{B}}(s, \alpha)}{\partial\lambda} \Big|_{\lambda=0} \right| \left| \sum_{s,\alpha} d_\rho^\pi(s) (\pi - \pi')(\alpha|s) \right| \\
&\leq \max_{s,\alpha} \left| \frac{\partial Q_i^{\pi^\mathbb{B}}(s, \alpha)}{\partial\lambda} \Big|_{\lambda=0} \right| \sum_{j=1}^N \sqrt{A_j} \|\pi_j - \pi'_j\| \\
&\leq \max_{s,\alpha} \left| e_{s,\alpha}^\top \frac{\partial \mathcal{T}(\pi_\lambda^\mathbb{B})}{\partial\lambda} \Big|_{\lambda=0} r_i \right| \sum_{j=1}^N \sqrt{A_j} \|\pi_j - \pi'_j\| \\
&\leq \max_{s,\alpha} \left| e_{s,\alpha}^\top \frac{\partial (I - \mathcal{P}(\pi_\lambda^\mathbb{A}))^{-1}}{\partial\lambda} \Big|_{\lambda=0} r_i \right| \sum_{j=1}^N \sqrt{A_j} \|\pi_j - \pi'_j\|
\end{aligned}$$

$$\begin{aligned}
&\leq \max_{s,\alpha} \left| e_{s,\alpha}^\top (\mathcal{T}(\pi) \frac{\partial \mathcal{P}(\pi_\lambda^\mathbb{A})}{\partial \lambda} \Big|_{\lambda=0} \mathcal{T}(\pi)) r_i \right| \sum_{j=1}^N \sqrt{A_j} \|\pi_j - \pi'_j\| \\
&\leq \frac{\sqrt{A_j}}{\zeta^2} \sum_{j=1}^N \sqrt{A_j} \|\pi_j - \pi'_j\|
\end{aligned} \tag{D.36}$$

using again (D.33) and (D.32) and Proposition D.1. Thus, we are ready now to bound the gradient per player:

$$\left| \frac{\partial (V_{i,\rho}(\pi_\lambda^\mathbb{A}) - V_{i,\rho}(\pi_\lambda^\mathbb{B}))}{\partial \lambda} \Big|_{\lambda=0} \right| \leq \text{Term}_A + Z_\rho^{\pi^\mathbb{A}} (\text{Term}_B + \text{Term}_C) \leq \frac{3\sqrt{A_i}}{\zeta^3} \sum_{j=1}^N \sqrt{A_j} \|\pi_j - \pi'_j\| \tag{D.37}$$

where we recall that  $Z_\rho^{\pi^\mathbb{A}} \leq \frac{1}{\zeta}$ . Since we prove it for an arbitrary perturbation vector  $\text{pert}$  for the directional derivative, for the independent player's policy gradient it holds also that:

$$\|v_i(\pi) - v_i(\pi')\| = \|\nabla_i(V_{i,\rho}(\pi) - V_{i,\rho}(\pi'))\| \leq \frac{3\sqrt{A_i}}{\zeta^3} \sum_{j=1}^N \sqrt{A_j} \|\pi_j - \pi'_j\| \quad \forall i \in \mathcal{N} \tag{D.38}$$

Finally for the concatenated gradient operator we get:

$$\begin{aligned}
\|v(\pi) - v(\pi')\| &= \sqrt{\sum_{i \in \mathcal{N}} \|v_i(\pi) - v_i(\pi')\|^2} = \sqrt{\sum_{i \in \mathcal{N}} \|\nabla_i(V_{i,\rho}(\pi) - V_{i,\rho}(\pi'))\|^2} \\
&\leq \sqrt{\sum_{i \in \mathcal{N}} \frac{9A_i}{\zeta^6} \left( \sum_{j \in \mathcal{N}} \sqrt{A_j} \|\pi_j - \pi'_j\| \right)^2} \leq \sqrt{\sum_{i \in \mathcal{N}} \frac{9A_i}{\zeta^6} \sum_{j \in \mathcal{N}} A_j \sum_{j \in \mathcal{N}} \|\pi_j - \pi'_j\|^2} \\
&\leq \frac{3}{\zeta^3} \sqrt{\left( \sum_{i \in \mathcal{N}} A_i \right)^2 \|\pi - \pi'\|^2} \leq \frac{3A}{\zeta^3} \|\pi - \pi'\|
\end{aligned} \tag{D.39}$$

which completes our proof.  $\blacksquare$

*Remark 5.* In the proof above, we considered perturbations that may formally lie outside the game's policy space. However, it is not difficult to see that for sufficiently small  $\lambda$  both  $V_i(\pi(\lambda))$ ,  $\nabla V_i(\pi(\lambda))$ ,  $Z_\rho^{\pi(\lambda)}$  are well-defined and bounded, for  $\pi_\lambda^\mathbb{A}(\alpha|s) = (\pi_i + \lambda \text{pert}, \pi_{-i}, \pi_{-i})$ . In view of this, we may harmlessly assume that all functions considered above are defined in an open neighborhood of the players' policy space.

## APPENDIX E. STATISTICS OF REINFORCE

In this appendix, we prove the two fundamental properties of the REINFORCE Policy Gradient estimator that we stated in Lemma 4, namely:

- a) REINFORCE is an unbiased estimator of  $v(\pi)$ .
- b) The variance of REINFORCE is bounded from above by  $\mathcal{O}(1/\min_{s \in \mathcal{S}, \alpha_i \in \mathcal{A}_i} \pi_i(\alpha_i|s))$  for each  $i \in \mathcal{N}$ .

We recall here that  $\nabla_i$  denotes the gradient of the quantity in question with respect to  $\pi_i$ , i.e., when  $\pi_{-i}$  is kept fixed and only  $\pi_i$  is varied. For concision, we will write  $v_i(\pi) = \nabla_i V_{i,\rho}(\pi)$  for the individual gradient of player  $i$ 's value function, and  $v(\pi) = (v_i(\pi))_{i \in \mathcal{N}}$  for the ensemble thereof.

With all this said and done, we begin by restating Lemma 4 for convenience:

**Lemma 4.** *Suppose that each agent  $i \in \mathcal{N}$  follows a stationary policy  $\pi_i \in \Pi_i$ . Then:*

$$a) \quad \mathbb{E}_{\tau \sim \text{MDP}}[\text{REINFORCE}(\pi)] = v(\pi) \quad (12a)$$

$$b) \quad \mathbb{E}_{\tau \sim \text{MDP}}[\|\text{REINFORCE}_i(\pi) - v_i(\pi)\|^2] \leq \frac{24A_i}{\kappa_i \zeta^4} \quad (12b)$$

where  $\kappa_i = \min_{s \in \mathcal{S}, \alpha_i \in \mathcal{A}_i} \pi_i(\alpha_i | s)$ .

*Proof.* Without loss of generality let's assume that  $\text{MDP} \equiv \text{MDP}(\pi | \rho)$  for some initial state distribution  $\rho$ . Additionally, we denote  $\mathbb{P}^\pi(\tau)$  the induced probability of a random trajectory  $\tau = (s_t, \alpha_t, r_t)_{t \leq T(\tau)}$ .

$$\begin{aligned} \mathbb{E}_{\tau \sim \text{MDP}}[\hat{v}_i] &= \mathbb{E}_{\tau \sim \text{MDP}}[R_i(\tau) \cdot \Lambda_i(\tau)] = \sum_{\tau \in \mathcal{T}} \mathbb{P}^\pi(\tau) R_i(\tau) \cdot \Lambda_i(\tau) \\ &= \sum_{\tau \in \mathcal{T}} \mathbb{P}^\pi(\tau) R_i(\tau) \cdot \left[ \sum_{t=0}^{T(\tau)} \nabla_i(\log \pi_i(\alpha_{i,t} | s_t)) \right] \\ &= \sum_{\tau \in \mathcal{T}} \mathbb{P}^\pi(\tau) R_i(\tau) \cdot \nabla_i \left[ \sum_{t=0}^{T(\tau)} \log \pi_i(\alpha_{i,t} | s_t) \right] \\ &= \sum_{\tau \in \mathcal{T}} \mathbb{P}^\pi(\tau) R_i(\tau) \nabla_i \sum_{t=0}^{T(\tau)} \log \pi_i(\alpha_{i,t} | s_t) \\ &\quad + \sum_{\tau \in \mathcal{T}} \mathbb{P}^\pi(\tau) R_i(\tau) \left( \nabla_i \sum_{j \neq i} \sum_{t=0}^{T(\tau)} \log \pi_j(\alpha_{j,t} | s_t) + \nabla_i \sum_{t=0}^{T(\tau)} \log \mathbb{P}(s_t | s_{t-1}, \alpha_{t-1}) \right) \\ &\quad + \sum_{\tau \in \mathcal{T}} \mathbb{P}^\pi(\tau) R_i(\tau) \nabla_i \log \rho(s_0) \\ &= \sum_{\tau \in \mathcal{T}} \mathbb{P}^\pi(\tau) R_i(\tau) \nabla_i(\log \mathbb{P}^\pi(\tau)) \\ &= \sum_{\tau \in \mathcal{T}} (\nabla_i \mathbb{P}^\pi(\tau)) R_i(\tau) = \nabla_i \left( \sum_{\tau \in \mathcal{T}} \mathbb{P}^\pi(\tau) R_i(\tau) \right) = \nabla_i V_{i,\rho}(\pi) \end{aligned} \quad (E.1)$$

where in the penultimate inequality we used the definition for the derivative of the logarithm. We also note here that

$$\mathbb{E}_{\tau \sim \text{MDP}}[\hat{v}_i] = \mathbb{E}_{\tau \sim \text{MDP}}[R_i(\tau) \nabla_i(\log \mathbb{P}^\pi(\tau))] \quad (E.2)$$

For the variance of REINFORCE estimator we have that

$$\begin{aligned} \mathbb{E}_{\tau \sim \text{MDP}}[\|\text{REINFORCE}_i(\pi) - v_i(\pi)\|^2] &= \mathbb{E}_{\tau \sim \text{MDP}}[\|\text{REINFORCE}_i(\pi)\|^2] \\ &\quad - 2 \mathbb{E}_{\tau \sim \text{MDP}}[\langle \text{REINFORCE}_i(\pi), v_i(\pi) \rangle] \\ &\quad + \mathbb{E}_{\tau \sim \text{MDP}}[\|v_i(\pi)\|^2] \end{aligned}$$

or equivalently  $\mathbb{E}_{\tau \sim \text{MDP}}[\|\text{REINFORCE}_i(\pi) - v_i(\pi)\|^2] = \mathbb{E}_{\tau \sim \text{MDP}}[\|\text{REINFORCE}_i(\pi)\|^2] - \mathbb{E}_{\tau \sim \text{MDP}}[\|v_i(\pi)\|^2]$ . Therefore, we have that

$$\mathbb{E}_{\tau \sim \text{MDP}}[\|\text{REINFORCE}_i(\pi) - v_i(\pi)\|^2] \leq \mathbb{E}_{\tau \sim \text{MDP}}[\|\text{REINFORCE}_i(\pi)\|^2] = \mathbb{E}[\|\hat{v}_i\|^2] \quad (E.3)$$

and, after a series of – tedious but otherwise straightforward – calculations, we get:

$$\mathbb{E}[\|\hat{v}_i\|^2] = \mathbb{E}_{\tau \sim \text{MDP}}[\|R_i(\tau) \Lambda_i(\tau)\|^2] \leq \mathbb{E}_{\tau \sim \text{MDP}}[\|R_i(\tau)\|^2 \|\Lambda_i(\tau)\|^2]$$

$$\begin{aligned}
&\leq \mathbb{E}_{\tau \sim \text{MDP}}[(T(\tau) + 1)^2 \|\sum_{t=0}^{T(\tau)} \nabla_i \log \pi_i(\alpha_{i,t}, s_t)\|^2] \\
&\leq \mathbb{E}_{\tau \sim \text{MDP}}[(T(\tau) + 1)^3 \sum_{t=0}^{\infty} \sum_{s, \alpha \in \mathcal{S} \times \mathcal{A}_i} \mathbb{1}\{t \leq T\} \mathbb{1}\{s_t = s, \alpha_{i,t} = \alpha\} \|\nabla_i \log \pi_i(\alpha, s)\|^2] \\
&= \sum_{t=0}^{\infty} \sum_{s, \alpha \in \mathcal{S} \times \mathcal{A}_i} \mathbb{E}_{\tau \sim \text{MDP}}[(T(\tau) + 1)^3 \mathbb{1}\{t \leq T\} \mathbb{1}\{s_t = s, \alpha_{i,t} = \alpha\} \frac{1}{(\pi_i(\alpha, s))^2}] \\
&\leq \sum_{t=0}^{\infty} \sum_{s, \alpha \in \mathcal{S} \times \mathcal{A}_i} \frac{1}{(\pi_i(\alpha, s))^2} \mathbb{E}_{\tau \sim \text{MDP}}[(T(\tau) + 1)^3 \mathbb{1}\{t \leq T\} \mathbb{1}\{s_t = s, \alpha_{i,t} = \alpha\}] \\
&\leq \sum_{t=0}^{\infty} \sum_{s, \alpha \in \mathcal{S} \times \mathcal{A}_i} \frac{1}{\pi_i(\alpha, s)} \mathbb{E}_{\tau \sim \text{MDP}}[(T(\tau) + 1)^3 \mathbb{1}\{t \leq T\} \mathbb{1}\{s_t = s\}] \\
&\leq \sum_{t=0}^{\infty} \sum_{s, \alpha \in \mathcal{S} \times \mathcal{A}_i} \frac{1}{\kappa_i} \{(T(\tau) + 1)^3 \mathbb{1}\{t \leq T\} \mathbb{1}\{s_t = s\}\} \\
&= \sum_{t=0}^{\infty} \sum_{s \in \mathcal{S}} \frac{|A_i|}{\kappa_i} \mathbb{E}_{\tau \sim \text{MDP}}[(T(\tau) + 1)^3 \mathbb{1}\{t \leq T\} \mathbb{1}\{s_t = s\}] \\
&= \frac{|A_i|}{\kappa_i} \mathbb{E}_{\tau \sim \text{MDP}}[(T(\tau) + 1)^3 \sum_{t=0}^T \mathbb{1}\{t \leq T\}] \\
&\leq \frac{|A_i|}{\kappa_i} \mathbb{E}_{\tau \sim \text{MDP}}[(T(\tau) + 1)^4] \\
&\leq \frac{|A_i|}{\kappa_i} \sum_{t=0}^{\infty} (1 - \zeta)^t \zeta (t + 1)^4 \leq \frac{24}{\zeta^4} \frac{|A_i|}{\kappa_i} \tag{E.4}
\end{aligned}$$

where, to go from the first to the second inequality we used the boundness by one of the rewards, while from the second to the third, we used Jensen's inequality.  $\blacksquare$

#### APPENDIX F. SOLUTION CONCEPTS

In this last appendix, we proceed to establish three important facts regarding the gradient characterization of stationary Nash policies. More precisely, we prove the following:

- In [Lemma 2](#), we prove the crucial property of Gradient Dominance for the multi-agent random stopping setting.
- In [Lemma 3](#), we establish that any stationary point corresponds to Nash Equilibria.
- In [Proposition 1](#), we prove the “drift” inequalities for all the different types of stationary points.

We begin with the gradient dominance property of the game, which we restate below for convenience:

**Lemma 2** (Gradient dominance property). *For any policy profile  $\pi = (\pi_i)_{i \in \mathcal{N}} \in \Pi$ , we have that*

$$V_{i,\rho}(\pi'_i; \pi_{-i}) - V_{i,\rho}(\pi_i; \pi_{-i}) \leq \mathcal{C}_G \max_{\bar{\pi}_i \in \Pi_i} \langle \nabla_i V_{i,\rho}(\pi), \bar{\pi}_i - \pi_i \rangle \tag{GDP}$$

for any unilateral deviation  $\pi'_i \in \Pi_i$  of player  $i \in \mathcal{N}$ .

*Proof.* We start by rewriting the LHS of (GDP) using Lemmas 1 and D.6 for  $\pi^{\mathbb{A}} = (\pi'_i; \pi_{-i})$  and  $\pi^{\mathbb{B}} = (\pi_i; \pi_{-i})$ :

$$\begin{aligned}
V_{i,\rho}(\pi^{\mathbb{A}}) - V_{i,\rho}(\pi^{\mathbb{B}}) &= \sum_{s \in \mathcal{S}} \tilde{d}_\rho^{\pi^{\mathbb{A}}}(s) \mathbb{E}_{\alpha \sim \pi^{\mathbb{A}}(\cdot|s)} \left[ A_i^{\pi^{\mathbb{B}}}(s, \alpha) \right] \\
&= \sum_{s \in \mathcal{S}} \tilde{d}_\rho^{\pi^{\mathbb{A}}}(s) \sum_{\alpha_i \in \mathcal{A}_i} \pi'_i(\alpha_i|s) \sum_{\alpha_{-i} \in \mathcal{A}_{-i}} \pi_{-i}(\alpha_{-i}|s) A_i^{\pi^{\mathbb{B}}}(s, \alpha) \\
&= \sum_{s \in \mathcal{S}} \tilde{d}_\rho^{\pi^{\mathbb{A}}}(s) \sum_{\alpha_i \in \mathcal{A}_i} \pi'_i(\alpha_i|s) \bar{A}_i^{\pi^{\mathbb{B}}}(s, \alpha_i) \\
&\leq \sum_{s \in \mathcal{S}} \tilde{d}_\rho^{\pi^{\mathbb{A}}}(s) \sum_{\alpha_i \in \mathcal{A}_i} \pi'_i(\alpha_i|s) \max_{\alpha_i \in \mathcal{A}_i} \bar{A}_i^{\pi^{\mathbb{B}}}(s, \alpha_i) \tag{F.1}
\end{aligned}$$

Thus, by a series of direct calculations, we obtain:

$$\begin{aligned}
V_{i,\rho}(\pi^{\mathbb{A}}) - V_{i,\rho}(\pi^{\mathbb{B}}) &\leq \max_{\tilde{\pi}_i \in \Delta(\mathcal{A})^{\mathcal{S}}} \sum_{s \in \mathcal{S}} \tilde{d}_\rho^{\pi^{\mathbb{A}}}(s) \sum_{\alpha_i \in \mathcal{A}_i} \tilde{\pi}_i(\alpha_i|s) \bar{A}_i^{\pi^{\mathbb{B}}}(s, \alpha_i) \\
&\leq \max_{\tilde{\pi}_i \in \Delta(\mathcal{A})^{\mathcal{S}}} \sum_{s \in \mathcal{S}} \tilde{d}_\rho^{\pi^{\mathbb{A}}}(s) \sum_{\alpha_i \in \mathcal{A}_i} (\tilde{\pi}_i(\alpha_i|s) - \pi_i(\alpha_i|s)) \bar{A}_i^{\pi^{\mathbb{B}}}(s, \alpha_i) \\
&\leq \max_{\tilde{\pi}_i \in \Delta(\mathcal{A})^{\mathcal{S}}} \sum_{s \in \mathcal{S}} \frac{\tilde{d}_\rho^{\pi^{\mathbb{A}}}(s)}{\tilde{d}_\rho^{\pi^{\mathbb{B}}}(s)} \tilde{d}_\rho^{\pi^{\mathbb{B}}}(s) \sum_{\alpha_i \in \mathcal{A}_i} (\tilde{\pi}_i(\alpha_i|s) - \pi_i(\alpha_i|s)) \bar{A}_i^{\pi^{\mathbb{B}}}(s, \alpha_i) \\
&\leq \left\| \frac{\tilde{d}_\rho^{\pi^{\mathbb{A}}}(s)}{\tilde{d}_\rho^{\pi^{\mathbb{B}}}(s)} \right\|_{\infty} \max_{\tilde{\pi}_i \in \Delta(\mathcal{A})^{\mathcal{S}}} \sum_{s \in \mathcal{S}} \sum_{\alpha_i \in \mathcal{A}_i} \tilde{d}_\rho^{\pi^{\mathbb{B}}}(s) (\tilde{\pi}_i(\alpha_i|s) - \pi_i(\alpha_i|s)) \bar{Q}_i^{\pi^{\mathbb{B}}}(s, \alpha_i) \\
&\leq \left\| \frac{\tilde{d}_\rho^{\pi^{\mathbb{A}}}(s)}{\tilde{d}_\rho^{\pi^{\mathbb{B}}}(s)} \right\|_{\infty} \max_{\tilde{\pi}_i \in \Delta(\mathcal{A})^{\mathcal{S}}} \sum_{s \in \mathcal{S}, \alpha_i \in \mathcal{A}_i} (\tilde{\pi}_i(\alpha_i|s) - \pi_i(\alpha_i|s)) \tilde{d}_\rho^{\pi^{\mathbb{B}}}(s) \bar{Q}_i^{\pi^{\mathbb{B}}}(s, \alpha_i) \\
&\leq \left\| \frac{\tilde{d}_\rho^{\pi^{\mathbb{A}}}(s)}{\tilde{d}_\rho^{\pi^{\mathbb{B}}}(s)} \right\|_{\infty} \max_{\tilde{\pi}_i \in \Delta(\mathcal{A})^{\mathcal{S}}} \sum_{s \in \mathcal{S}, \alpha_i \in \mathcal{A}_i} (\tilde{\pi}_i(\alpha_i|s) - \pi_i(\alpha_i|s)) \frac{\partial V_{i,\rho}(\pi)}{\partial \pi_i(\alpha_i|s)} \tag{F.2}
\end{aligned}$$

so

$$V_{i,\rho}(\pi'_i; \pi_{-i}) - V_{i,\rho}(\pi_i; \pi_{-i}) \leq \mathcal{C}_G \max_{\tilde{\pi}_i \in \Pi_i} \langle \nabla_i V_{i,\rho}(\pi), \tilde{\pi}_i - \pi_i \rangle \tag{F.3}$$

and our proof is complete.  $\blacksquare$

*Remark.* Notice that we have assumed that  $\tilde{d}_\rho^{\pi^{\mathbb{B}}} > 0$ . If this wasn't the case we could take a trivial bound of  $\infty$ .

We now proceed to establish the link between (FOS) and (NE):

**Lemma 3** (First-order stationary policies are Nash). *A policy  $\pi^* = (\pi_i^*)_{i \in \mathcal{N}} \in \Pi$  is Nash if and only if it satisfies the first-order stationary condition*

$$\langle v(\pi^*), \pi - \pi^* \rangle \leq 0 \quad \text{for all } \pi \in \Pi. \tag{FOS}$$

*Proof.* By the definition of first-order stationarity applied to the policies  $\pi^*$  and  $\pi$ , it is straightforward to check that  $\langle v(\pi^*), \pi^* - \pi \rangle \geq 0$  if and only if  $\max_{\tilde{\pi}_i \in \Pi_i} \langle \nabla_i V_{i,\rho}(\pi^*), \pi_i - \tilde{\pi}_i^* \rangle \leq 0$ . However, by the gradient dominance property established in Lemma 2, we readily get

$$V_{i,\rho}(\pi_i; \pi_{-i}^*) - V_{i,\rho}(\pi_i^*; \pi_{-i}^*) \leq \mathcal{C}_G \max_{\tilde{\pi}_i \in \Pi_i} \langle \nabla_i V_{i,\rho}(\pi^*), \tilde{\pi}_i - \pi_i^* \rangle \leq 0 \tag{F.4}$$

and hence

$$V_{i,\rho}(\pi_i; \pi_{-i}^*) \leq V_{i,\rho}(\pi_i^*; \pi_{-i}^*) \quad \text{for all } \pi_i \in \Pi_i \quad (\text{F.5})$$

and our claim follows.  $\blacksquare$

With all this in place, we are finally in a position to prove the characterization of second-order stationary and strict Nash policies that of [Proposition 1](#). For ease of reference, we restate the relevant claims below.

**Proposition 1.** *Let  $\pi^* = (\pi_i^*)_{i \in \mathcal{N}} \in \Pi$  be a Nash policy. Then:*

a) *If  $\pi^*$  is second-order stationary, there exists some  $\mu > 0$  such that*

$$\langle v(\pi), \pi - \pi^* \rangle \leq -\mu \|\pi - \pi^*\|^2 \quad \text{for all } \pi \text{ sufficiently close to } \pi^*. \quad (3a)$$

b) *If  $\pi^*$  is strict, there exists some  $\mu > 0$  such that*

$$\langle v(\pi), \pi - \pi^* \rangle \leq -\mu \|\pi - \pi^*\| \quad \text{for all } \pi \text{ sufficiently close to } \pi^*. \quad (3b)$$

*Proof.* We begin with the characterization of second-order stationary policies. To that end, let  $d = |\mathcal{S}| \sum_i |\mathcal{A}_i|$  denote the ambient dimension of  $\prod_i \mathbb{R}^{\mathcal{A}_i \times \mathcal{S}}$  and consider the mapping  $\varphi: \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$  mapping  $H \mapsto \max\{z^\top H z : z \in \text{TC}(\pi^*), \|z\| = 1\}$ . Clearly,  $\varphi$  is convex as the pointwise maximum of a set of linear – and hence convex – functions. This in turn implies the continuity of  $\varphi$  as every convex function is continuous on the interior of its effective domain. Since  $\pi^*$  satisfies (SOS) by assumption, we have  $\varphi(\text{Jac}_v(\pi^*)) < 0$ , so, by continuity and the convexity of  $\Pi$ , there exists some  $\mu > 0$  and a convex neighborhood  $\mathcal{U}$  of  $\pi^*$  in  $\Pi$  such that  $\varphi(\text{Jac}_v(\pi)) \leq -\mu$  for all  $\pi \in \mathcal{U}$ .

With this in mind, letting  $z = \pi - \pi^* \in \text{TC}(\pi^*)$  for some  $\pi \in \mathcal{U}$ , a straightforward Taylor expansion with integral remainder yields

$$v(\pi) - v(\pi^*) = \int_0^1 \text{Jac}_v(\pi^* + \tau z) z \, d\tau \quad (\text{F.6})$$

and hence, setting  $\pi_\tau = \pi^* + \tau z$ , we get

$$\begin{aligned} \langle v(\pi) - v(\pi^*), \pi - \pi^* \rangle &= \int_0^1 z^\top \text{Jac}_v(\pi_\tau) z \, d\tau \\ &\leq \|z\|^2 \int_0^1 \varphi(\text{Jac}_v(\pi_\tau)) \, d\tau \leq -\mu \|z\|^2 = -\mu \|\pi - \pi^*\|^2 \end{aligned} \quad (\text{F.7})$$

However, by (FOS), we have  $\langle v(\pi^*), \pi - \pi^* \rangle \leq 0$  which, combined with the above, yields  $\langle v(\pi), \pi - \pi^* \rangle \leq -\mu \|\pi - \pi^*\|^2$ , as claimed.

For the second part of our lemma, pick some  $\pi \neq \pi^*$  and let  $z = (\pi - \pi^*)/\|\pi - \pi^*\|$ , so  $z \in \text{TC}(\pi^*)$  and  $\|z\| = 1$ . Then, given that (FOS) is satisfied as a strict inequality for all  $\pi \neq \pi^*$ , we readily get  $\langle v(\pi^*), z \rangle < 0$  for all  $z \in \text{TC}(\pi^*)$  with  $\|z\| = 1$ . Thus, by the joint continuity of the function  $\langle v(\pi), z \rangle$  in  $\pi$  and  $z$ , there exists a compact convex neighborhood  $\mathcal{K}$  of  $\pi^*$  in  $\Pi$  such that  $\mu := \min\{\langle v(\pi), z \rangle : \pi \in \mathcal{K}, z \in \text{TC}(\pi^*), \|z\| = 1\} < 0$ . Thus, letting  $z = (\pi - \pi^*)/\|\pi - \pi^*\|$  as above, we conclude that  $\langle v(\pi), \pi - \pi^* \rangle \leq -\mu \|\pi - \pi^*\|$ , as claimed.  $\blacksquare$

#### ACKNOWLEDGMENTS

Part of this work was done while the authors were visiting the Simons Institute for the Theory of Computing. P. Mertikopoulos gratefully acknowledges financial support by the French National Research Agency (ANR) in the framework of the ‘‘Investissements d’avenir’’ program (ANR-15-IDEX-02), the LabEx PERSYVAL (ANR-11-LABX-0025-01), MIAI@Grenoble Alpes (ANR-19-P3IA-0003), and the bilateral ANR-NRF grant ALIAS (ANR-19-CE48-0018-01). K. Lotidis and E. Vlatakis are

grateful for financial support by the Onassis Foundation (F ZR 033-1/2021-2022, 010-1/2018-2019). A. Giannou is grateful for financial support by ONR: “A Theoretically Principled Framework for Learning by Pruning”. E. V. Vlatakis-Gkaragkounis is grateful for financial support by the Google-Simons Fellowship, Pancretan Association of America and Simons Collaboration on Algorithms and Geometry. This project was completed while he was a visiting research fellow at the Simons Institute for the Theory of Computing. Additionally, he would like to acknowledge the following series of NSF-CCF grants under the numbers 1763970/2107187/1563155/1814873.

## REFERENCES

- [1] Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. Optimality and approximation with policy gradient methods in markov decision processes. In Abernethy, J. D. and Agarwal, S. (eds.), *Conference on Learning Theory, COLT 2020, 9-12 July 2020, Virtual Event [Graz, Austria]*, volume 125 of *Proceedings of Machine Learning Research*, pp. 64–66. PMLR, 2020. URL <http://proceedings.mlr.press/v125/agarwal20a.html>.
- [2] Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98):1–76, 2021.
- [3] Azizian, W., Iutzeler, F., Malick, J., and Mertikopoulos, P. The last-iterate convergence rate of optimistic mirror descent in stochastic variational inequalities. In *COLT '21: Proceedings of the 34th Annual Conference on Learning Theory*, 2021.
- [4] Baxter, J. and Bartlett, P. L. Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research*, 15:319–350, 2001.
- [5] Bhandari, J. and Russo, D. Global optimality guarantees for policy gradient methods. *arXiv preprint arXiv:1906.01786*, 2019.
- [6] Bubeck, S. Convex optimization: Algorithms and complexity. *Foundations and Trends in Machine Learning*, 8(3-4):231–358, 2015.
- [7] Cen, S., Wei, Y., and Chi, Y. Fast policy extragradient methods for competitive games with entropy regularization. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021*, pp. 27952–27964, 2021.
- [8] Chasnov, B., Ratliff, L., Mazumdar, E., and Burden, S. Convergence analysis of gradient-based learning in continuous games. In *UAI '20: Proceedings of the 35th Annual Conference on Uncertainty in Artificial Intelligence*, 2020.
- [9] Chen, J. and Jiang, N. Information-theoretic considerations in batch reinforcement learning. In *International Conference on Machine Learning*, pp. 1042–1051. PMLR, 2019.
- [10] Chung, K.-L. On a stochastic approximation method. *The Annals of Mathematical Statistics*, 25(3): 463–483, 1954.
- [11] Daskalakis, C., Foster, D. J., and Golowich, N. Independent policy gradient methods for competitive reinforcement learning. *Advances in neural information processing systems*, 33:5527–5540, 2020.
- [12] Daskalakis, C., Golowich, N., and Zhang, K. The complexity of markov equilibrium in stochastic games. *arXiv preprint arXiv:2204.03991*, 2022.
- [13] Ding, D., Wei, C.-Y., Zhang, K., and Jovanović, M. R. Independent policy gradient for large-scale markov potential games: Sharper rates, function approximation, and game-agnostic convergence. *arXiv preprint arXiv:2202.04129*, 2022.
- [14] Fan, J., Wang, Z., Xie, Y., and Yang, Z. A theoretical analysis of deep q-learning. In *Learning for Dynamics and Control*, pp. 486–489. PMLR, 2020.
- [15] Filar, J. and Vrieze, K. *Competitive Markov Decision Processes*. Springer, 1997.
- [16] Fink, A. M. Equilibrium in a stochastic  $n$ -person game. *Journal of science of the hiroshima university, series ai (mathematics)*, 28(1):89–93, 1964.
- [17] Flokas, L., Vlatakis-Gkaragkounis, E. V., Lianas, T., Mertikopoulos, P., and Piliouras, G. No-regret learning and mixed Nash equilibria: They do not mix. In *NeurIPS '20: Proceedings of the 34th International Conference on Neural Information Processing Systems*, 2020.
- [18] Folland, G. B. *Real Analysis*. Wiley-Interscience, 2 edition, 1999.

- [19] Giannou, A., Vlatakis-Gkaragkounis, E. V., and Mertikopoulos, P. Survival of the strictest: Stable and unstable equilibria under regularized learning with partial information. In *COLT '21: Proceedings of the 34th Annual Conference on Learning Theory*, 2021.
- [20] Hall, P. and Heyde, C. C. *Martingale Limit Theory and Its Application*. Probability and Mathematical Statistics. Academic Press, New York, 1980.
- [21] Hart, S. and Mas-Colell, A. A simple adaptive procedure leading to correlated equilibrium. *Econometrica*, 68(5):1127–1150, September 2000.
- [22] Hart, S. and Mas-Colell, A. Uncoupled dynamics do not lead to Nash equilibrium. *American Economic Review*, 93(5):1830–1836, 2003.
- [23] Hsieh, Y.-G., Iutzeler, F., Malick, J., and Mertikopoulos, P. On the convergence of single-call stochastic extra-gradient methods. In *NeurIPS '19: Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pp. 6936–6946, 2019.
- [24] Hsieh, Y.-G., Iutzeler, F., Malick, J., and Mertikopoulos, P. Explore aggressively, update conservatively: Stochastic extragradient methods with variable stepsize scaling. In *NeurIPS '20: Proceedings of the 34th International Conference on Neural Information Processing Systems*, 2020.
- [25] Hsieh, Y.-P., Mertikopoulos, P., and Cevher, V. The limits of min-max optimization algorithms: Convergence to spurious non-critical sets. In *ICML '21: Proceedings of the 38th International Conference on Machine Learning*, 2021.
- [26] Jin, C., Liu, Q., Wang, Y., and Yu, T. V-Learning: A simple, efficient, decentralized algorithm for multiagent RL. <https://arxiv.org/abs/2110.14555>, 2021.
- [27] Jin, Y., Muthukumar, V., and Sidford, A. The complexity of infinite-horizon general-sum stochastic games, 2022.
- [28] Kakade, S. M. A natural policy gradient. *Advances in neural information processing systems*, 14, 2001.
- [29] Konda, V. and Tsitsiklis, J. Actor-critic algorithms. *Advances in neural information processing systems*, 12, 1999.
- [30] Laraki, R., Renault, J., and Sorin, S. *Mathematical Foundations of Game Theory*. Universitext. Springer, 2019.
- [31] Leonardos, S., Overman, W., Panageas, I., and Piliouras, G. Global convergence of multi-agent policy gradient in markov potential games. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=gfwON7rAm4>.
- [32] Leslie, D. S., Perkins, S., and Xu, Z. Best-response dynamics in zero-sum stochastic games. *Journal of Economic Theory*, 189:105095, 2020.
- [33] Mertikopoulos, P. and Zhou, Z. Learning in games with continuous action sets and unknown payoff functions. *Mathematical Programming*, 173(1-2):465–507, January 2019.
- [34] Mertikopoulos, P., Papadimitriou, C. H., and Piliouras, G. Cycles in adversarial regularized learning. In *SODA '18: Proceedings of the 29th annual ACM-SIAM Symposium on Discrete Algorithms*, 2018.
- [35] Mertikopoulos, P., Lecouat, B., Zenati, H., Foo, C.-S., Chandrasekhar, V., and Piliouras, G. Optimistic mirror descent in saddle-point problems: Going the extra (gradient) mile. In *ICLR '19: Proceedings of the 2019 International Conference on Learning Representations*, 2019.
- [36] Mertikopoulos, P., Hsieh, Y.-P., and Cevher, V. Learning in games from a stochastic approximation viewpoint. <https://arxiv.org/abs/2206.03922>, 2022.
- [37] Moravčík, M., Schmid, M., Burch, N., Lisý, V., Morrill, D., Bard, N., Davis, T., Waugh, K., Johanson, M., and Bowling, M. Deepstack: Expert-level artificial intelligence in heads-up no-limit poker. *Science*, 356(6337):508–513, 2017. doi: 10.1126/science.aam6960.
- [38] Munos, R. Error bounds for approximate policy iteration. In *ICML*, volume 3, pp. 560–567, 2003.
- [39] Munos, R. Error bounds for approximate value iteration. In *Proceedings of the National Conference on Artificial Intelligence*, volume 20, pp. 1006. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 2005.
- [40] Nesterov, Y. *Introductory Lectures on Convex Optimization: A Basic Course*. Number 87 in Applied Optimization. Kluwer Academic Publishers, 2004.
- [41] Neyman, A. and Sorin, S. (eds.). *Stochastic Games and Applications*. NATO ASI. Kluwer Academic Publishers, 2003.
- [42] Perkins, S. *Advanced stochastic approximation frameworks and their applications*. PhD thesis, University of Bristol, 2013.

- [43] Polyak, B. Gradient methods for the minimisation of functionals. *USSR Computational Mathematics and Mathematical Physics*, 3(4):864–878, 1963. ISSN 0041-5553.
- [44] Polyak, B. T. *Introduction to Optimization*. Optimization Software, New York, NY, USA, 1987.
- [45] Ratliff, L. J., Burden, S. A., and Sastry, S. S. On the characterization of local Nash equilibria in continuous games. *IEEE Trans. Autom. Control*, 61(8):2301–2307, August 2016.
- [46] Rockafellar, R. T. and Wets, R. J. B. *Variational Analysis*, volume 317 of *A Series of Comprehensive Studies in Mathematics*. Springer-Verlag, Berlin, 1998.
- [47] Sayin, M., Zhang, K., Leslie, D., Basar, T., and Ozdaglar, A. Decentralized q-learning in zero-sum markov games. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 18320–18334. Curran Associates, Inc., 2021.
- [48] Sayin, M. O., Parise, F., and Ozdaglar, A. Fictitious play in zero-sum stochastic games. *arXiv preprint arXiv:2010.04223*, 2020.
- [49] Shalev-Shwartz, S., Shammah, S., and Shashua, A. Safe, multi-agent, reinforcement learning for autonomous driving. *arXiv preprint arXiv:1610.03295*, 2016.
- [50] Shapley, L. S. Stochastic games. *Proceedings of the National Academy of Sciences of the USA*, 39: 1095–1100, 1953.
- [51] Shoham, Y. and Leyton-Brown, K. *Multiagent systems: Algorithmic, game-theoretic, and logical foundations*. Cambridge University Press, 2008.
- [52] Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., Chen, Y., Lillicrap, T. P., Hui, F., Sifre, L., van den Driessche, G., Graepel, T., and Hassabis, D. Mastering the game of go without human knowledge. *Nat.*, 550(7676):354–359, 2017. doi: 10.1038/nature24270.
- [53] Solan, E. and Vieille, N. Stochastic games. *Proceedings of the National Academy of Sciences*, 112(45): 13743–13746, 2015.
- [54] Sutton, R. S. and Barto, A. G. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, 1998.
- [55] Sutton, R. S., McAllester, D., Singh, S., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.
- [56] Takahashi, M. Stochastic games with infinitely many strategies. *Journal of Science of the Hiroshima University, Series AI (Mathematics)*, 26(2):123–134, 1962.
- [57] Vinyals, O., Babuschkin, I., Czarnecki, W. M., Mathieu, M., Dudzik, A., Chung, J., Choi, D. H., Powell, R., Ewalds, T., Georgiev, P., Oh, J., Horgan, D., Kroiss, M., Danihelka, I., Huang, A., Sifre, L., Cai, T., Agapiou, J. P., Jaderberg, M., Vezhnevets, A. S., Leblond, R., Pohlen, T., Dalibard, V., Budden, D., Sulsky, Y., Molloy, J., Paine, T. L., Gülçehre, Ç., Wang, Z., Pfaff, T., Wu, Y., Ring, R., Yogatama, D., Wünsch, D., McKinney, K., Smith, O., Schaul, T., Lillicrap, T. P., Kavukcuoglu, K., Hassabis, D., Apps, C., and Silver, D. Grandmaster level in starcraft II using multi-agent reinforcement learning. *Nat.*, 575(7782):350–354, 2019. doi: 10.1038/s41586-019-1724-z.
- [58] Wei, C.-Y., Lee, C.-W., Zhang, M., and Luo, H. Last-iterate convergence of decentralized optimistic gradient descent/ascent in infinite-horizon competitive markov games. In Belkin, M. and Kpotufe, S. (eds.), *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134 of *Proceedings of Machine Learning Research*, pp. 4259–4299. PMLR, 15–19 Aug 2021.
- [59] Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
- [60] Zhang, K., Koppel, A., Zhu, H., and Basar, T. Global convergence of policy gradient methods to (almost) locally optimal policies. *SIAM Journal on Control and Optimization*, 58(6):3586–3612, 2020.
- [61] Zhang, R., Ren, Z., and Li, N. Gradient play in multi-agent markov stochastic games: Stationary points and convergence. *arXiv e-prints*, pp. arXiv–2106, 2021.
- [62] Zinkevich, M. Online convex programming and generalized infinitesimal gradient ascent. In *ICML '03: Proceedings of the 20th International Conference on Machine Learning*, pp. 928–936, 2003.