

# Robustness and sample complexity of model-based MARL for general-sum Markov games

Jayakumar Subramanian · Amit Sinha ·  
Aditya Mahajan

the date of receipt and acceptance should be inserted later

**Abstract** Multi-agent reinforcement learning (MARL) is often modeled using the framework of Markov games (also called stochastic games or dynamic games). Most of the existing literature on MARL concentrates on zero-sum Markov games but is not applicable to general-sum Markov games. It is known that the best-response dynamics in general-sum Markov games are not a contraction. Therefore, different equilibria in general-sum Markov games can have different values. Moreover, the Q-function is not sufficient to completely characterize the equilibrium. Given these challenges, model based learning is an attractive approach for MARL in general-sum Markov games. In this paper, we investigate the fundamental question of *sample complexity* for model-based MARL algorithms in general-sum Markov games. We show two results. We first use Hoeffding inequality based bounds to show that  $\tilde{O}((1-\gamma)^{-2}\alpha^{-2})$  samples per state-action pair are sufficient to obtain a  $\alpha$ -approximate Markov perfect equilibrium with high probability, where  $\gamma$  is the discount factor, and the  $\tilde{O}(\cdot)$  notation hides logarithmic terms. We then use Bernstein inequality based bounds to show that  $\tilde{O}((1-\gamma)^{-1}\alpha^{-2})$  samples are sufficient. We characterize the exact constants for these two bounds. To obtain these results, we study the robustness of Markov perfect equilibrium to model approximations. We show

---

The work of Amit Sinha and Aditya Mahajan was supported in part by the Innovation for Defence Excellence and Security (IDEaS) Program of the Canadian Department of National Defence through grant CFPMN2-30.

A preliminary version of this work appeared in the 2021 Indian Control Conference.

---

J. Subramanian  
Media and Data Science Research Lab, Digital Experience Cloud,  
Adobe Inc., Noida, Uttar Pradesh, India.  
E-mail: jasubram@adobe.com

A. Sinha and A. Mahajan  
Department of Electrical and Computer Engineering,  
McGill University, Montreal, Canada.  
E-mail: amit.sinha@mail.mcgill.ca, aditya.mahajan@mcgill.ca

that the Markov perfect equilibrium of an approximate (or perturbed) game is always an approximate Markov perfect equilibrium of the original game and provide explicit bounds on the approximation error. We illustrate the results via a numerical example.

## 1 Introduction

Markov games (also called stochastic games or dynamic games) are a commonly used framework to model strategic interaction between multiple players interacting in a dynamic environment. Examples include applications in cyber-security (Sengupta et al. 2019), industrial organization (Ericson and Pakes 1995; Fershtiman and Pakes 2000), political economics (Acemoglu and Robinson 2001), advertisement and pricing Albright and Winston (1979), and many others (Başar and Zaccour 2018). Starting with the seminal work of Shapley (1953), several variations of Markov games have been considered in the literature. We refer the reader to Filar and Vrieze (1996) for an overview.

**Overview of Markov games.** In the basic setup of a dynamic game, the payoffs of players at any time not only depend on their current joint action profile but also on the current “state of the system”. Furthermore, the state of the system evolves in a controlled Markov manner conditioned on the current action profile of the players. It is typically assumed that the state of the system and the action profile of all players are publicly monitored by all players. Although Markov games may be viewed as a special case of extensive form games with perfect information, rather than using the standard solution concept of sub-game perfect equilibrium, attention is often restricted to a refinement of sub-game perfect equilibrium called Markov perfect equilibrium (MPE) where all players play Markov strategies (i.e., choose their actions as a (possibly randomized) function of the current state) (Maskin and Tirole 1988a,b). MPE is an attractive refinement of sub-game perfect equilibrium, both from a computational as well as conceptual point of view, but has some limitations because it excludes some history dependent strategies (such as tit-for-tat and grim trigger) commonly used in the repeated games setup. See Maskin and Tirole (2001); Mailath and Samuelson (2006) for a discussion.

Games can also be classified based on the sum of per-step payoffs of players as zero-sum or general-sum games. The nature of results in these two cases are different as are the tools used to prove them. The differences stem from the fact that the best response mapping (called the Shapley operator) for two-player zero-sum games is a contraction (Shapley 1953). Therefore, zero-sum games have a unique value (i.e., all equilibria in zero-sum games have the same value). Moreover, the MPE (also called minimax equilibrium for the zero-sum case) can be computed via recursive operations of the Shapley operator (Shapley 1953; Hoffman and Karp 1966). In contrast, the best response mapping for general-sum games is not a contraction. Therefore, the existence of MPE needs to be proved using variations of Kakutani’s fixed point theorem (Fink 1964; Takahashi 1964; Rogers 1969; Vrieze 1987). A consequence of this

is that, in general, different MPEs do not have the same value, which makes it difficult to compute MPEs. Various algorithms have been proposed to compute MPE, including non-linear programming (Breton 1991; Filar et al. 1991) and homotopy methods (Herings et al. 2004; Herings and Peeters 2010). It was recently established by Deng et al. (2021) that the computational complexity of computing MPE is PPAD-complete.

In spite of these challenges, computing MPE of general-sum games is an important research direction because several real-world problems are not zero-sum. The applications of network security (Sengupta et al. 2019), industrial organization (Ericson and Pakes 1995; Fershtiman and Pakes 2000), and political economics (Acemoglu and Robinson 2001) mentioned above are all general-sum games.

**Multi-agent reinforcement learning.** In recent years, there has been significant interest in understanding interaction between strategic agents operating in unknown environments. Such multi-player problems are studied under the heading of multi-agent reinforcement learning (MARL) and often modeled as Markov games (Littman 1994; Busoniu et al. 2008; Zhang et al. 2021a). Although there have been significant recent successes in single agent RL, these do not directly translate into the multi-agent setting. Part of the difficulty is that when multiple agents are learning simultaneously, the “environment” as viewed by any single agent is non-stationary (Busoniu et al. 2008); so it is not possible to use the theoretical guarantees of single agent RL algorithms, which are derived for a stationary or time-homogeneous environment.

Nonetheless, MARL for two player zero-sum games is well understood due to two properties. First, if two strategies  $(\pi^1, \pi^2)$  and  $(\mu^1, \mu^2)$  are minimax equilibrium, then so are strategies  $(\pi^1, \mu^2)$  and  $(\mu^1, \pi^2)$ . Therefore, to identify equilibrium strategies, it is sufficient to learn the action-value function (i.e., the  $Q$ -function). Second, the action-value function can be learnt using variants of  $Q$ -learning (called minimax  $Q$ -learning) because the Shapley operator is a contraction (Littman 1994, 2001). We refer the reader to Shoham et al. (2003) for an overview of MARL for zero-sum games.

However, the situation is different for general-sum MARL, where fewer convergence guarantees are available. Part of the difficulty is that the action-value function (or  $Q$ -function) is insufficient to characterize MPE (Zinkevich et al. 2006, Theorem 1).<sup>1</sup> For this reason, algorithms developed for two-player zero-sum games fail to converge to an MPE in general-sum games (Pérolat et al. 2017). There are some partial results, e.g., minimizing Bellman residual error to identify  $\varepsilon$ -MPE (Pérolat et al. 2017), using two-time scale stochastic approximation algorithms (Prasad et al. 2015), and using replicator dynamics

<sup>1</sup> Zinkevich et al. (2006) construct two player general-sum games with the following properties. The game has two states: in state 1, player 1 has two actions and player 2 has one action; in state 2, player 1 has one action and player 2 has two actions. The transition probabilities are chosen such that there is a unique Markov perfect equilibrium in mixed strategies. This means that in state 1, both actions of player 1 maximize the  $Q$ -function; in state 2, both actions of player 2 minimize the  $Q$ -function. However, the  $Q$ -function in itself is insufficient to determine the randomizing probabilities for the mixed strategy MPE.

based algorithms (Akchurina 2010). However, in general, developing MARL algorithms with convergence guarantees remains a challenging research direction.

There is some recent work on learning in general-sum games. Leonardos et al. (2021) show global convergence of a policy gradient algorithm in a special class of Markov games called Markov potential games, which are a generalization of normal form potential games and assume the existence of a common potential function for all players. Zhang et al. (2021b) also consider Markov potential games and present convergence analysis for a sample-based RL method in such games. While both these results are interesting, they do not apply to Markov games which do not have a potential function. Another recent result is presented by Song et al. (2021), who present and analyse algorithms for computing correlated and coarse correlated equilibria for general-sum games. The solution concepts of correlated equilibrium and its variations are different from MPE. In correlated equilibrium, players agree on a joint randomization strategies before the system starts running; such pre-game agreement is not allowed in MPE.

**Model based MARL, sample complexity, and robustness of equilibria.** One potential approach to alleviate the difficulties in MARL for general-sum games is to use model based algorithms, which explicitly learn (or estimate) the system model and then use a “planning algorithm” to find the solution of the estimated model (Sutton 1990). There has been significant recent interest in model based RL for single agent systems (see Wang et al. (2019) and references therein) and some interest in model-based approaches for MARL for zero-sum games (Krupnik et al. 2019; Sidford et al. 2020; Zhang et al. 2021c, 2020). However, as far as we are aware, there are no model based MARL algorithms for general-sum Markov games.

An important consideration in model-based RL is to determine how many samples are needed to identify an  $\alpha$ -approximate solution (for a pre-specified accuracy level  $\alpha$ ). This is known as *sample complexity* of learning and is typically analyzed under the assumption that the learning agent has access to a generative model, i.e., a black box simulator that takes the current state and action profile as input and generates samples of the next state as output.

Starting with the work of Kearns and Singh (1999); Kakade (2003), there is an extensive literature on the sample complexity of Markov decision processes (MDPs) (Azar et al. 2013; Sidford et al. 2018; Agarwal et al. 2020; Li et al. 2020). The simplest approach in this setting is to use a plug-in estimator,<sup>2</sup> i.e., estimating the transition matrix using the generated samples and using the optimal policy corresponding to the estimated model in the true system. Recent results of Agarwal et al. (2020) show that the sample complexity of the plug-in estimator matches the lower bounds on sample complexity (Azar et al. 2013) modulo logarithmic factors. Recently, Zhang et al. (2020), build on this line of work to establish sample complexity bounds for zero-sum games. As

<sup>2</sup> The plug-in estimator is also known as a certainty equivalent controller in the stochastic control literature.

far as we are aware, sample complexity of generative models for general-sum games hasn't been investigated before.

The analyses of model-based RL algorithms rely on the *robustness* of the “planning solution” to model approximations, i.e., *if the estimated model is close to the true model in some sense, does that imply that the strategy generated from the estimated model is approximately appropriate in some sense (optimality, equilibrium, etc.)?* This question is well understood for Markov decision processes (see Müller (1997) and follow-up work) and zero-sum Markov games (Tidball and Altman 1996; Tidball et al. 1997). In this paper, we address the question of robustness for general-sum Markov games. In particular, we show that if a dynamic game is approximated by another game such that the reward functions and transitions of the approximate game are close to those of the original game (in an appropriate sense), then a MPE of the approximate game is an approximate MPE of the original game. We quantify the exact relationship between the degree of approximation of the games and the approximation error in the MPE. We then build up on these results to establish sample complexity bounds for learning with a generative model for general-sum Markov games.

The notion of robustness is also useful in its own right. In many applications, the model of a dynamic game is estimated using modern econometric techniques Aguirregabiria and Mira (2007); Bajari et al. (2007); Pakes et al. (2007); Pesendorfer and Schmidt-Dengler (2008). In such situations, robustness characterizes the approximation error in using a MPE of an approximate game, in terms of the approximation errors in estimating the reward function and transition dynamics of the game.

**Other notions of robustness.** Our notion of robustness is different from that of robust control (Başar and Bernhard 2008) and robust Markov perfect equilibrium (Jaśkiewicz and Nowak 2014), both of which are Markov decision processes with uncertain dynamics and are treated as zero-sum games where nature acts as an adversary and picks the worst-case realization of the transition dynamics. Our notion of robustness is also different from uniformly  $\varepsilon$ -equilibrium (Solan 2021), which captures robustness with respect to time-horizon and discount factor.

Our notion of robustness is similar in spirit to robust MPEs considered in Maskin and Tirole (2001), who defined a MPE to be robust if for any small perturbation of the payoffs, there exists a nearby MPE. Maskin and Tirole (2001) showed that almost all finite horizon general-sum games have a finite number of MPEs, all of which are robust. Our results are of a different nature and it is difficult to compare the two results because Maskin and Tirole (2001) considered an atypical model where the states are not specified exogenously but are rather determined as the payoff relevant component of the history. Consequently, perturbing the payoffs changes the state *space*, which is not the case for our model.

The notion of strong stability considered in Doraszelski and Escobar (2010) is related to the work in Maskin and Tirole (2001). It is shown in Doraszelski

and Escobar (2010) that almost all Markov games have finite number of MPEs and these equilibria can be approximated by equilibria of nearby games. The dynamics in Doraszelski and Escobar (2010) are exogenous and, therefore, their result does not have the same limitations as that of Maskin and Tirole (2001). The result of Doraszelski and Escobar (2010) is stronger than ours because we only show that equilibria of nearby games are approximate equilibria of the original game but we do not establish that they are also close to the equilibria of the original game. However, the results of Doraszelski and Escobar (2010) rely on continuity arguments and do not explicitly characterize bounds on the size of the neighborhood. In contrast, for any  $\varepsilon$  perturbation in payoffs and  $\delta$  perturbation in dynamics, we explicitly characterize an  $\alpha$  such that the MPE of the perturbed game is an  $\alpha$ -MPE of the original game.

Perhaps the result most similar to ours is Whitt (1980), who consider a more general model and allow the approximate game to have a different state and action space than the original game. Their main result is to show that any  $\alpha_{\text{opt}}$ -MPE of the approximate game is an  $\alpha$ -MPE of the original game and an explicit relationship between  $\alpha_{\text{opt}}$  and  $\alpha$  is established. Our results are similar in spirit but the specific details are different.

**Organization.** The rest of the paper is organized as follows. In Sec. 2, we present our notion of approximation of a dynamic game and state our main results. In Sec. 3, we present background results on approximation of Markov decision processes. In Sec. 4, we provide the proof of our main results. In Sec. 5, we present numerical examples to validate our theoretical results. We conclude in Sec. 6.

**Notation.** We use  $\mathbb{R}$  to denote the set of real numbers,  $\mathbb{P}(\cdot)$  to denote the probability of an event,  $\mathbb{E}[\cdot]$  to denote the expectation of a random variable, and  $\mathcal{P}(\cdot)$  denotes the set of probability measures on a set.

We use calligraphic letters (e.g.,  $\mathcal{S}$ ,  $\mathcal{A}$ , etc.) to denote sets, uppercase letters (e.g.,  $S$ ,  $A$ , etc.) to denote random variables and lowercase letters (e.g.,  $s$ ,  $a$ , etc.) to denote their realization. Superscripts index players and subscripts index time. For example,  $a_t^i$  denotes the action of player  $i$  at time  $t$ . For sequence of variables  $\{s_t\}_{t \geq 1}$ , we use the short hand notation  $s_{1:t}$  to denote the sequence  $(s_1, \dots, s_t)$ . We use  $\mathbf{1}$  to denote a vector of ones of an appropriate size which is determined by context.

Given a function  $f: \mathcal{S} \rightarrow \mathbb{R}$ , we use  $\text{span}(f)$  to denote the span seminorm of  $f$ , i.e.,  $\text{span}(f) := \sup_{s \in \mathcal{S}} f(s) - \inf_{s \in \mathcal{S}} f(s)$ . Given a metric space  $(\mathcal{S}, d)$  and a function  $f: \mathcal{S} \rightarrow \mathbb{R}$ , we use  $\text{Lip}(f)$  to denote the Lipschitz constant of  $f$ , i.e.,

$$\text{Lip}(f) := \sup_{s, s' \in \mathcal{S}} \frac{|f(s) - f(s')|}{d(s, s')}.$$

## 2 System model, robustness, and sample complexity

We restrict the discussion in this paper to models with finite state and action spaces. The robustness results can be extended to models with continuous state

and action spaces under standard technical assumptions on the existence of equilibria in that setting.

## 2.1 Dynamic games

An infinite horizon dynamic game (also called stochastic game or Markov game) is a tuple  $\langle \mathcal{N}, \mathcal{S}, (\mathcal{A}^i)_{i \in \mathcal{N}}, \mathbf{P}, (r^i)_{i \in \mathcal{N}}, \gamma \rangle$  where:

- $\mathcal{N}$  is the (finite) set of players.
- $\mathcal{S}$  is the (finite) set of possible states of the game. We use  $S_t \in \mathcal{S}$  to denote the state of the game at time  $t$ .
- $(\mathcal{A}^i)_{i \in \mathcal{N}}$  is the (finite) set of actions available to player  $i$  at each time. We also use  $\mathcal{A} = \prod_{i \in \mathcal{N}} \mathcal{A}^i$  to denote the set of actions of all players. We use  $A_t = (A_t^i)_{i \in \mathcal{N}}$  to denote the action profile of all players at time  $t$ . Given an action profile  $A_t = (A_t^i)_{i \in \mathcal{N}}$  and a player  $j \in \mathcal{N}$ , we use the notation  $A_t^{-j} = (A_t^i)_{i \in \mathcal{N} \setminus \{j\}}$  to denote the action profile of all players except  $j$ .
- $\mathbf{P} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$  is the controlled transition probability of the state of the game. In particular, at any time  $t$ , given a realization  $s_{1:t+1}$  of  $S_{1:t+1}$  and choice of action profile  $a_{1:t}$  of  $A_{1:t}$ , we have

$$\mathbf{P}(s_{t+1} \mid s_t, a_t) := \mathbf{P}(S_{t+1} = s_{t+1} \mid S_{1:t} = s_{1:t}, A_{1:t} = a_{1:t}).$$

- $r^i : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  denotes the per-step reward of player  $i$ .
- $\gamma \in (0, 1)$  is the discount factor.

We assume that all players have perfect monitoring. At time  $t$ , all players observe the current state  $S_t$  and simultaneously choose their respective actions. At the end of time period  $t$ , all players observe all the actions, and the state of the game evolves according to the transition kernel  $\mathbf{P}$ .

Following Shapley (1953); Maskin and Tirole (1988a,b), we assume that each player chooses its action according to a time homogeneous Markov strategy. Let

$$\Pi^i := \{\pi^i : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A}^i)\}$$

denote the set of all Markov strategies for player  $i$ .

Given a strategy profile  $\pi = (\pi^i)_{i \in \mathcal{N}}$ , where  $\pi^i \in \Pi^i$ , and an initial state  $s_0$ , the expected discounted total reward of player  $i$  is given by:

$$V_{(\pi^i, \pi^{-i})}^i(s_0) := (1 - \gamma) \mathbb{E}_{(\pi^i, \pi^{-i})} \left[ \sum_{t=0}^{\infty} \gamma^t r^i(S_t, A_t) \mid S_0 = s_0 \right], \quad (1)$$

where the expectation is with respect to the joint measure on all the system variables induced by the choice of the strategy profile of all players.

Although the above model is formulated for an infinite horizon, it can capture interactions for a finite horizon by considering time as part of the state space and by assuming that, at the end of the horizon, the game moves to an absorbing state with zero rewards for all players. In the special case when the game has a single state, a dynamic game is equivalent to an infinitely

repeated matrix game. In the special case when the game has only one player, a dynamic game is equivalent to a Markov decision process.

*Remark 1* We follow the standard game theoretic convention of normalizing the expected total reward by pre-multiplying by  $(1 - \gamma)$ . An immediate implication of this is that for any strategy  $\pi$ ,  $\|V_\pi^i\|_\infty \leq \|r\|_\infty$ . In some of the AI literature, the expected reward is not normalized. In such cases  $\|V_\pi^i\|_\infty \leq \|r\|_\infty / (1 - \gamma)$ .

There are two solution concepts commonly used for Markov games, which we state below.

**Definition 1 (Markov perfect equilibrium)** A Markov strategy profile  $\pi = (\pi^i)_{i \in \mathcal{N}}$ , where  $\pi^i \in \Pi^i$ , is called a Markov perfect equilibrium (MPE) if for every initial state  $s \in \mathcal{S}$ , and every player  $i \in \mathcal{N}$ ,

$$V_{(\pi^i, \pi^{-i})}^i(s) \geq V_{(\tilde{\pi}^i, \pi^{-i})}^i(s), \quad \forall \tilde{\pi}^i \in \Pi^i. \quad (2)$$

A Markov perfect equilibrium can be viewed as a refinement of subgame perfect equilibrium where all players play Markov strategies (Maskin and Tirole 1988a,b). For games with finite state and action spaces, a Markov perfect equilibrium always exists (Fink 1964; Rogers 1969; Vrieze 1987; Filar and Vrieze 1996). For general state and action spaces, see Takahashi (1964).

We can also describe MPE equilibrium in terms of best responses. Given a player  $i \in \mathcal{N}$  and a strategy profile  $\pi^{-i}$  of players other than  $i$ , a strategy  $\pi^i$  for player  $i$  is called a best response of  $\pi^{-i}$  if it satisfies (2). We denote this relationship by  $\pi^i = \text{BR}^i(\pi^{-i})$ . Then, Definition 1 is equivalent to stating that a strategy profile  $\pi$  is an MPE if, for each player  $i \in \mathcal{N}$ ,  $\pi^i$  is a best response of  $\pi^{-i}$ .

**Definition 2 (Approximate Markov perfect equilibrium)** Given approximation level  $\alpha = (\alpha^i)_{i \in \mathcal{N}}$ , where  $\alpha^i$  are positive constants, a strategy profile  $\pi = (\pi^i)_{i \in \mathcal{N}}$ , where  $\pi^i \in \Pi^i$ , is called an  $\alpha$ -approximate Markov perfect equilibrium ( $\alpha$ -MPE) if for every initial state  $s \in \mathcal{S}$ , and every player  $i \in \mathcal{N}$ ,

$$V_{(\pi^i, \pi^{-i})}^i(s) \geq V_{(\tilde{\pi}^i, \pi^{-i})}^i(s) - \alpha^i, \quad \forall \tilde{\pi}^i \in \Pi^i. \quad (3)$$

When all  $\alpha^i$  are identical and equal to say  $\alpha'$ , we simply call the approximate MPE an  $\alpha'$ -MPE rather than  $(\alpha', \dots, \alpha')$ -MPE.

## 2.2 Preliminaries on integral probability metrics

Our results rely on a class of metrics on probability spaces known as integral probability metrics (IPMs) (Müller 1997).



**Definition 3** Let  $(\mathcal{X}, \mathcal{G})$  be a measurable space and  $\mathfrak{F}$  denote a class of uniformly bounded measurable functions on  $(\mathcal{X}, \mathcal{G})$ . The integral probability metric (IPM) between two probability distributions  $\mu, \nu \in \mathcal{P}(\mathcal{X})$  with respect to the function class  $\mathfrak{F}$  is defined as

$$d_{\mathfrak{F}}(\mu, \nu) := \sup_{f \in \mathfrak{F}} \left| \int_{\mathcal{X}} f d\mu - \int_{\mathcal{X}} f d\nu \right|.$$

Two specific forms of IPMs are used in this paper:

1. **Total variation distance:** If  $\mathfrak{F}$  is chosen as  $\mathfrak{F}^{\text{TV}} := \{f : \text{span}(f) \leq 1\}$ , then  $d_{\mathfrak{F}}$  is the total variation distance<sup>3</sup>.
2. **Wasserstein distance:** If  $\mathcal{X}$  is a metric space and  $\mathfrak{F}$  is chosen as  $\mathfrak{F}^{\text{W}} := \{f : \text{Lip}(f) \leq 1\}$  (where the Lipschitz constant is computed with respect to the metric on  $\mathcal{X}$ ), then  $d_{\mathfrak{F}}$  is the Wasserstein distance.

See Subramanian et al. (2022) for a discussion of other IPMs and their role in approximate planning for single agent problems. Our approximation results are stated in terms of the Minkowski functional of a function  $f$  (not necessarily in  $\mathfrak{F}$ ) with respect to a function class  $\mathfrak{F}$ , which is defined as follows:

$$\rho_{\mathfrak{F}}(f) := \inf\{\rho \in \mathbb{R}_{>0} : \rho^{-1}f \in \mathfrak{F}\}. \quad (4)$$

A key implication of this definition is that for any function  $f$ ,

$$\left| \int_{\mathcal{X}} f d\mu - \int_{\mathcal{X}} f d\nu \right| \leq \rho_{\mathfrak{F}}(f) \cdot d_{\mathfrak{F}}(\mu, \nu). \quad (5)$$

The Minkowski functional of the two IPMs considered in this paper are as follows:

1. **Total variation distance:** If  $\mathfrak{F}$  is chosen as  $\mathfrak{F}^{\text{TV}}$ ,  $|\int_{\mathcal{X}} f d\mu - \int_{\mathcal{X}} f d\nu| \leq \text{span}(f) d_{\mathfrak{F}}(\mu, \nu)$ . Thus, for total variation,  $\rho_{\mathfrak{F}^{\text{TV}}}(f) = \text{span}(f)$ .
2. **Wasserstein distance:** If  $\mathfrak{F}$  is chosen as  $\mathfrak{F}^{\text{W}}$ ,  $|\int_{\mathcal{X}} f d\mu - \int_{\mathcal{X}} f d\nu| \leq \text{Lip}(f) \cdot d_{\mathfrak{F}}(\mu, \nu)$ . Thus, for the Wasserstein distance,  $\rho_{\mathfrak{F}^{\text{W}}}(f) = \text{Lip}(f)$ .

<sup>3</sup> If  $\mu$  and  $\nu$  are absolutely continuous with respect to some measure  $\lambda$  and let  $p = d\mu/d\lambda$  and  $q = d\nu/d\lambda$ , then total variation is typically defined as  $\frac{1}{2} \int_{\mathcal{X}} |p(x) - q(x)| \lambda(dx)$ . This is consistent with our definition. Let  $\bar{f} = (\sup f + \inf f)/2$ . Then

$$\begin{aligned} \left| \int_{\mathcal{X}} f d\mu - \int_{\mathcal{X}} f d\nu \right| &= \left| \int_{\mathcal{X}} f(x)p(x)\lambda(dx) - \int_{\mathcal{X}} f(x)q(x)\lambda(dx) \right| \\ &= \left| \int_{\mathcal{X}} [f(x) - \bar{f}] [p(x) - q(x)] \lambda(dx) \right| \leq \|f - \bar{f}\|_{\infty} \int_{\mathcal{X}} |p(x) - q(x)| \lambda(dx) \\ &\leq \frac{1}{2} \text{span}(f) \int_{\mathcal{X}} |p(x) - q(x)| \lambda(dx). \end{aligned}$$

### 2.3 Robustness of Markov games to model approximation

**Definition 4** Given a function class  $\mathfrak{F}$  and positive constants  $(\varepsilon, \delta)$ , a game  $\mathcal{G} := \langle \mathcal{N}, \mathcal{S}, (\mathcal{A}^i)_{i \in \mathcal{N}}, \hat{\mathbf{P}}, (\hat{r}^i)_{i \in \mathcal{N}}, \gamma \rangle$  is an  $(\varepsilon, \delta)$ -approximation of the game  $\mathcal{G} := \langle \mathcal{N}, \mathcal{S}, (\mathcal{A}^i)_{i \in \mathcal{N}}, \mathbf{P}, (r^i)_{i \in \mathcal{N}}, \gamma \rangle$  if the following conditions are satisfied:

1. **Reward approximation:** For all  $i \in \mathcal{N}$ ,  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$ ,

$$|r^i(s, a) - \hat{r}^i(s, a)| \leq \varepsilon. \quad (6)$$

2. **Transition approximation:** For all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$ ,

$$d_{\mathfrak{F}}(\mathbf{P}(\cdot | s, a), \hat{\mathbf{P}}(\cdot | s, a)) \leq \delta. \quad (7)$$

Our main result is the following.

**Theorem 1** If game  $\hat{\mathcal{G}}$  is an  $(\varepsilon, \delta)$ -approximation of game  $\mathcal{G}$  and  $\hat{\pi}$  is an  $\alpha_{\text{opt}}$ -MPE of  $\hat{\mathcal{G}}$ , where  $\alpha_{\text{opt}} := (\alpha_{\text{opt}}^i)_{i \in \mathcal{N}}$ , then  $\hat{\pi}$  is also an  $\alpha$ -MPE of  $\mathcal{G}$ , where  $\alpha := (\alpha^i)_{i \in \mathcal{N}}$  can be bounded as

$$\alpha^i \leq 2\varepsilon + \frac{\gamma}{1-\gamma} [\Delta_{(\hat{\pi}^i, \hat{\pi}^{-i})}^i + \Delta_{(\tilde{\pi}_*^i, \hat{\pi}^{-i})}^i] + \alpha_{\text{opt}}^i, \quad \forall i \in \mathcal{N}, \quad (8)$$

where  $\tilde{\pi}_*^i = \text{BR}^i(\hat{\pi}^{-i})$  and

$$\Delta_{(\hat{\pi}^i, \hat{\pi}^{-i})}^i := \max_{s \in \mathcal{S}, a \in \mathcal{A}} \left| \sum_{s' \in \mathcal{S}} \left[ \mathbf{P}(s' | s, a) \hat{V}_{(\hat{\pi}^i, \hat{\pi}^{-i})}^i(s') - \hat{\mathbf{P}}(s' | s, a) \hat{V}_{(\hat{\pi}^i, \hat{\pi}^{-i})}^i(s') \right] \right|$$

and  $\Delta_{(\tilde{\pi}_*^i, \hat{\pi}^{-i})}^i$  is defined similarly. Furthermore, we have  $\Delta_{(\hat{\pi}^i, \hat{\pi}^{-i})}^i \leq \delta \rho_{\mathfrak{F}}(\hat{V}_{(\hat{\pi}^i, \hat{\pi}^{-i})}^i)$ . and  $\Delta_{(\tilde{\pi}_*^i, \hat{\pi}^{-i})}^i \leq \delta \rho_{\mathfrak{F}}(\hat{V}_{(\tilde{\pi}_*^i, \hat{\pi}^{-i})}^i)$ . Therefore, a looser upper bound on  $\alpha^i$  is given by

$$\alpha^i \leq 2\varepsilon + \frac{\delta\gamma}{1-\gamma} [\rho_{\mathfrak{F}}(\hat{V}_{(\hat{\pi}^i, \hat{\pi}^{-i})}^i) + \rho_{\mathfrak{F}}(\hat{V}_{(\tilde{\pi}_*^i, \hat{\pi}^{-i})}^i)] + \alpha_{\text{opt}}^i, \quad \forall i \in \mathcal{N}. \quad (9)$$

See Sec. 4 for the proof.

*Remark 2* Both upper bounds (8) and (9) are instance dependent. The bound (8) depends on the transition probability  $\mathbf{P}$  of the original game, while the only information about the original game  $\mathcal{G}$  needed in (9) are the modeling errors  $(\varepsilon, \delta)$ . The approximation bound (9) also depends on the choice of IPM used to measure the approximation in the dynamics. Instance independent bounds on the approximation error are presented below.

### 2.3.1 Instance independent bounds

It is possible to obtain instance-independent bounds by using worst case upper bounds on  $\rho_{\mathfrak{F}}(\hat{V}_{\hat{\pi}}^i)$ . In particular, the Minkowski functional  $\rho_{\mathfrak{F}}(\hat{V}_{\hat{\pi}}^i)$  is  $\text{span}(\hat{V}_{\hat{\pi}}^i)$  when using the total variation distance and is  $\text{Lip}(\hat{V}_{\hat{\pi}}^i)$  when using the Wasserstein distance. Using worst-case upper bounds on  $\text{span}(\hat{V}_{\hat{\pi}}^i)$  and  $\text{Lip}(\hat{V}_{\hat{\pi}}^i)$  gives us the following bounds.

**Corollary 1** *When  $\mathfrak{F} = \mathfrak{F}^{\text{TV}}$ , then*

$$\alpha^i \leq 2 \left( \varepsilon + \frac{\gamma \delta \text{span}(\hat{r}^i)}{(1 - \gamma)} \right) + \alpha_{\text{opt}}^i, \quad \forall i \in \mathcal{N}. \quad (10)$$

The next bound holds for games where the transition matrix and reward function are Lipschitz.

**Definition 5** Suppose the state space  $\mathcal{S}$  is a metric space with metric  $d$ . Then, a game  $\mathcal{G}$  is said to be  $(L_r, L_P)$ -Lipschitz if for any  $i \in \mathcal{N}$ ,  $s_1, s_2 \in \mathcal{S}$  and  $a \in \mathcal{A}$ ,

$$|r^i(s_1, a) - r^i(s_2, a)| \leq L_r d(s_1, s_2),$$

and

$$d_{\mathfrak{F}^{\text{W}}}(\mathbf{P}(\cdot|s_1, a), \mathbf{P}(\cdot|s_2, a)) \leq L_P d(s_1, s_2),$$

where  $d_{\mathfrak{F}^{\text{W}}}$  denotes the Wasserstein distance.

Furthermore, define

$$\text{Lip}(\hat{\pi}^i) = \sup_{\substack{s, s' \in \mathcal{S} \\ s \neq s'}} \frac{d_{\mathfrak{F}^{\text{W}}}(\hat{\pi}^i(\cdot|s), \hat{\pi}^i(\cdot|s'))}{d(s, s')}.$$

**Corollary 2** *When  $\mathfrak{F} = \mathfrak{F}^{\text{W}}$  and  $\hat{\mathcal{G}}$  is  $(L_r, L_P)$ -Lipschitz with  $\gamma L_P < 1$  and  $\gamma(1 + \text{Lip}(\hat{\pi}^i))L_P < 1$  for all  $i \in \mathcal{N}$ , then*

$$\alpha^i \leq 2\varepsilon + \frac{\gamma L_r \delta}{(1 - \gamma L_P)} + \frac{\gamma L_r \delta}{(1 - \gamma(1 + \text{Lip}(\hat{\pi}^i))L_P)} + \alpha_{\text{opt}}^i, \quad \forall i \in \mathcal{N}. \quad (11)$$

Furthermore, when  $\alpha_{\text{opt}} = 0$ , we have

$$\alpha^i \leq 2 \left( \varepsilon + \frac{\gamma L_r \delta}{(1 - \gamma L_P)} \right) \quad \forall i \in \mathcal{N}. \quad (12)$$

*Remark 3* The result of Theorem 1 and Corollaries 1 and 2 imply that MPE is continuous in model approximation for total variation and Wasserstein distances. To show this, we first start with the case when  $\mathfrak{F} = \mathfrak{F}^{\text{TV}}$ . Consider a sequence  $\{\mathcal{G}_n\}_{n \geq 0}$  of approximations such that  $\mathcal{G}_n$  is an  $(\varepsilon_n, \delta_n)$  approximation of  $\mathcal{G}$ , where  $\varepsilon_n \rightarrow 0$  and  $\delta_n \rightarrow 0$  as  $n \rightarrow \infty$ . Let  $\hat{\pi}_n$  be an MPE of  $\mathcal{G}_n$ . Then, Theorem 1 shows that  $\hat{\pi}_n$  is an  $\alpha_n$  MPE of  $\mathcal{G}$ , where  $\alpha_n$  is upper bounded by  $2\varepsilon_n + 2\gamma\delta_n \text{span}(\hat{r}_n^i)/(1 - \gamma)$ . Since  $\text{span}(\hat{r}_n^i)$  is bounded,  $\alpha_n \rightarrow 0$  as  $n \rightarrow \infty$ . Thus, any limit point  $\hat{\pi}_\infty$  of  $\{\hat{\pi}_n\}_{n \geq 0}$  is an MPE of game  $\mathcal{G}$ . The same argument applies for the case when  $\mathfrak{F} = \mathfrak{F}^{\text{W}}$ , provided that each approximate game  $\mathcal{G}_n$  is  $(L_r, L_p)$  Lipschitz with  $\gamma L_p < 1$ .

*Remark 4* Although we have only elaborated on two specific choices of IPMs (total variation and Wasserstein distances), the result of Theorem 1 is applicable for *any* IPM. Many other IPMs have been considered in the literature including Kolmogorov distance, bounded Lipschitz metric, and maximum mean discrepancy. See, for example, Müller (1997); Subramanian et al. (2022). The choice of the metric often depends on the specific properties of the model.

## 2.4 Model based RL for Markov games

Now, we consider the setting where the components  $\langle \mathcal{S}, (\mathcal{A}^i)_{i \in \mathcal{N}}, (r^i)_{i \in \mathcal{N}}, \gamma \rangle$  of a game  $\mathcal{G}$  are known but the transition probability matrix  $\mathbf{P}$  is not known. Suppose we have access to a *generative model*, i.e., a black-box simulator which provides samples  $S_+ \sim \mathbf{P}(\cdot|s, a)$  of the next state  $S_+$  for a given state-action pair  $(s, a)$  as input. Suppose we call the simulator  $n$  times at each state-action pair and estimate an empirical model  $\hat{\mathbf{P}}_n$  as  $\hat{\mathbf{P}}_n(s'|s, a) := \text{count}(s'|s, a)/n$ , where  $\text{count}(s'|s, a)$  is the number of times  $s'$  is sampled with the input is  $(s, a)$ . The game  $\hat{\mathcal{G}}_n := \langle \mathcal{S}, (\mathcal{A}^i)_{i \in \mathcal{N}}, \hat{\mathbf{P}}_n, (r^i)_{i \in \mathcal{N}}, \gamma \rangle$  may be viewed as an approximation of game  $\mathcal{G}$ . We further assume that there is a planning oracle, which takes the approximate game  $\hat{\mathcal{G}}_n$  as input and generates an  $\alpha_{\text{opt}}$ -MPE  $\hat{\pi}_n$ , where  $\alpha_{\text{opt}} \in \mathbb{R}_{>0}$  is a property of the planning oracle.

Note that we are assuming that there is a system planner which generates samples from the generative model and computes an  $\alpha_{\text{opt}}$ -MPE of  $\hat{\mathcal{G}}_n$ . A more interesting setting is where each player generates independent samples from the generative model and computes a different MPE. This setting is challenging due to the multiplicity of MPE and is not discussed in this paper.

One fundamental question in this setting is the following. Given an  $\alpha > 0$ , how many samples  $n$  are needed from the generative model to ensure that  $\hat{\pi}_n$  is an  $\alpha$ -MPE for game  $\mathcal{G}$ . This is called the *sample complexity* of learning. Below, we obtain two bounds on sample complexity using the approximation bounds of Theorem 1.

**Theorem 2 (Hoeffding-Type Bound)** For any  $\alpha_{\text{opt}} > 0$ ,  $\alpha > \alpha_{\text{opt}}$  and  $p > 0$ , let

$$n \geq n_H(\gamma) := \left\lceil \left( \frac{\gamma}{1-\gamma} \text{span}(r) \right)^2 \frac{2 \log(4|\mathcal{S}|(\prod_{i \in \mathcal{N}} |\mathcal{A}^i|)|\mathcal{N}|p^{-1})}{(\alpha - \alpha_{\text{opt}})^2} \right\rceil.$$

Then, with probability  $1-p$ , any  $\alpha_{\text{opt}}$ -MPE  $\hat{\pi}_n$  of game  $\hat{\mathcal{G}}_n$  is  $\alpha$ -MPE for game  $\mathcal{G}$ .

**Theorem 3 (Bernstein-Type Bound)** For any  $\alpha_{\text{opt}} > 0$ ,  $\alpha \in (5\alpha_{\text{opt}}, \|r\|_\infty \sqrt{1-\gamma} + 5\alpha_{\text{opt}})$  and  $p > 0$ , let

$$n \geq n_B(\gamma) := \left\lceil \frac{c\gamma\|r\|_\infty^2 \log(8|\mathcal{S}|(\prod_{i \in \mathcal{N}} |\mathcal{A}^i|)|\mathcal{N}|(1-\gamma)^{-2}p^{-1})}{(1-\gamma)(\alpha - 5\alpha_{\text{opt}})^2} \right\rceil,$$

where  $c$  is an absolute constant. Then, with probability  $1-p$ , any  $\alpha_{\text{opt}}$ -MPE  $\hat{\pi}_n$  of game  $\hat{\mathcal{G}}_n$  is an  $\alpha$ -MPE for game  $\mathcal{G}$ .

See Sec. 4 for the proofs.

*Remark 5* In general, game  $\hat{\mathcal{G}}_n$  may have multiple MPE. The results of Theorems 2 and 3 are true for *every* MPE of game  $\hat{\mathcal{G}}_n$ .

*Remark 6* We have assumed that we generate  $n$  samples for every state-action pair. Therefore, the total number of samples required in Theorems 2 and 3 are  $n|\mathcal{S}| \prod_{i=1}^n |\mathcal{A}^i|$ .

## 2.5 Discussion of the results

### 2.5.1 Bounds on the values of the Markov perfect equilibrium

The robustness and sample complexity bounds presented above only show that an  $\alpha_{\text{opt}}$ -MPE  $\hat{\pi}$  of approximate game  $\hat{\mathcal{G}}$  is an  $\alpha$ -MPE of the original game  $\mathcal{G}$ . A natural question is whether the value of policy  $\hat{\pi}$  in game  $\mathcal{G}$  is close to the value of *some* MPE of  $\mathcal{G}$ . We do not know an answer to this question in general, but we provide an answer for the special case of two-player zero sum games (ZSG).

**Definition 6** A two player Markov game is called *zero-sum* if

$$r^1(s, a) + r^2(s, a) = 0, \quad \forall s \in \mathcal{S}, a \in \mathcal{A}.$$

Two important properties of zero-sum games are the following (Filar and Vrieze 1996):

(P1) For any Markov policy  $\pi$ ,

$$V_\pi^1(s) + V_\pi^2(s) = 0, \quad \forall s \in \mathcal{S}.$$

(P2) All MPE of ZSG have the same value, i.e., if  $\pi_*$  and  $\tilde{\pi}_*$  are MPEs of a ZSG, then

$$V_{\pi_*}^i(s) = V_{\tilde{\pi}_*}^i(s), \quad i \in \{1, 2\}, s \in \mathcal{S}.$$

We use  $V_*^i$  to denote the *value of the game*.

We now establish an important property of approximate MPEs in ZSGs.

**Proposition 1** *If  $\pi$  is an  $\alpha$ -MPE of a two-player ZSG with value  $V_* = (V_*^1, V_*^2)$ , where  $\alpha = (\alpha^1, \alpha^2)$ , then*

$$V_\pi^i(s) \geq V_*^i(s) - \alpha^i, \quad \forall s \in \mathcal{S}. \quad (13)$$

*Proof* Let  $\pi_* = (\pi_*^1, \pi_*^2)$  be any MPE of  $\mathcal{G}$ . For any  $s \in \mathcal{S}$ , we have the following:

$$\begin{aligned} V_\pi^1(s) &\stackrel{(a)}{\geq} V_{(\pi^1, \pi_*^2)}^1(s) - \alpha^1 \\ &\stackrel{(b)}{=} -V_{(\pi^1, \pi_*^2)}^2(s) - \alpha^1 \\ &\stackrel{(c)}{\geq} -V_{(\pi_*^1, \pi_*^2)}^2(s) - \alpha^1 \\ &\stackrel{(d)}{=} V_{(\pi_*^1, \pi_*^2)}^1(s) - \alpha^1 \\ &\stackrel{(e)}{=} V_*^1(s) - \alpha^1 \end{aligned}$$

where (a) follows from the definition of  $\alpha$ -MPE, (b) and (d) follow from (P1), and (c) follows from the fact that  $\pi_*$  is an MPE, and (e) follows from (P2). This proves (13) for  $i = 1$ . The result for  $i = 2$  follows from a similar argument.

Combining Proposition 1 with the robustness results of Theorem 1, we get the following.

**Corollary 3** *If the game  $\mathcal{G}$  in Theorem 1 is a two-player ZSG with value  $V_* = (V_*^1, V_*^2)$ , then under the conditions of Theorem 1, we have*

$$V_{\tilde{\pi}}^i(s) \geq V_*^i(s) - \alpha^i, \quad i \in \{1, 2\}, s \in \mathcal{S}$$

where  $\alpha^i$  is given by (8).

Similarly, combining Proposition 1 with Theorems 2 and 3, we get the following.

**Corollary 4** *If the game  $\mathcal{G}$  is a two-player ZSG with value  $V_* = (V_*^1, V_*^2)$ , then under the conditions of Theorem 2 or Theorem 3, we have*

$$V_{\tilde{\pi}_n}^i(s) \geq V_*^i(s) - \alpha, \quad i \in \{1, 2\}, s \in \mathcal{S}.$$

Note that Corollary 4 under the conditions of Theorem 3 is identical to Theorem 3.2 of Zhang et al. (2020). Thus, the results of our paper provide an alternative proof of the sample complexity results for zero-sum games.

### 2.5.2 Discussion on the sample complexity bounds

The bounds of Theorem 2 and 3 provide an upper bound on sample complexity for learning MPE. As argued in Zhang et al. (2020), the lower bound on the sample complexity for MDPs obtained in Azar et al. (2013) (see Remark 10 later) can be directly translated to games to provide a lower bound of

$$\Omega\left(|\mathcal{S}|\left(\prod_{i=1}^n |\mathcal{A}^i|\right) \frac{\log(|\mathcal{S}|(\sum_{i \in \mathcal{N}} |\mathcal{A}^i|)p^{-1})}{(1-\gamma)\alpha^2}\right).$$

Thus, upper bounds of Theorems 2 and 3 are tight in  $|\mathcal{S}|$ , but loose in the logarithmic factor in  $\{|\mathcal{A}^i|\}$  ( $\sum_{i \in \mathcal{N}} |\mathcal{A}^i|$  vs  $\prod_{i \in \mathcal{N}} |\mathcal{A}^i|$ ). The sample complexity bounds for zero-sum games obtained in Zhang et al. (2020) were also loose in the action space. Zhang et al. (2020) also highlight the difficulty in tightening the lower bound to  $\Omega(\log(\prod_{i \in \mathcal{N}} |\mathcal{A}^i|))$ . Since, zero-sum games are a special case of general-sum games, the same difficulties hold for the general-sum games as well.

In addition, the bound of Theorem 3 is tight in the discount factor while the bound of Theorem 2 is loose by a factor of  $1/(1-\gamma)$ . This means that as the discount factor  $\gamma \rightarrow 1$ , the bound of Theorem 3 is tighter. Therefore, we must have the following.

**Proposition 2** *There exists a critical discount factor  $\gamma^\circ$  (perhaps depending on  $|\mathcal{S}|$  and  $\{|\mathcal{A}^i|\}_{i \in \mathcal{N}}$ ) such that for all  $\gamma \geq \gamma^\circ$ ,  $n_H(\gamma) \geq n_B(\gamma)$ , where  $n_H(\gamma)$  and  $n_B(\gamma)$  are as defined in Theorems 2 and 3, respectively.*

It is not possible to derive a closed form expression on  $\gamma^\circ$ , but we can obtain a lower bound as follows.

**Lemma 1** *When  $\alpha_{\text{opt}} = 0$ , for all  $\gamma \leq c/(c+8)$ , we have  $n_H(\gamma) \leq n_B(\gamma)$ . Thus,  $c/(c+8)$  is a lower bound on  $\gamma^\circ$ .*

*Proof* Consider the following sequence of inequalities (for  $\alpha_{\text{opt}} = 0$ )

$$\begin{aligned} n_H(\gamma) &\stackrel{(a)}{\leq} \frac{8\gamma^2 \|r\|_\infty^2 \log(4|\mathcal{S}|(\prod_{i \in \mathcal{N}} |\mathcal{A}^i|)|\mathcal{N}|p^{-1})}{(1-\gamma)^2 \alpha^2} \\ &\stackrel{(b)}{\leq} \frac{c\gamma \|r\|_\infty^2 \log(4|\mathcal{S}|(\prod_{i \in \mathcal{N}} |\mathcal{A}^i|)|\mathcal{N}|p^{-1})}{(1-\gamma)\alpha^2}, \\ &\stackrel{(c)}{\leq} n_B(\gamma) \end{aligned}$$

where (a) uses the fact that  $\text{span}(r) \leq 2\|r\|_\infty$ , (b) uses the fact that for any  $\gamma \leq c/(c+8)$ , we have  $\gamma/(1-\gamma) \leq c/8$  and (c) relies on the fact that  $2(1-\gamma)^{-2} \geq 1$ .  $\square$

We now present a lower bound on  $c$ , which implies a lower bound on  $\gamma^\circ$ .

**Lemma 2** *A lower bound on the constant  $c$  in Theorem 3 is given by  $c \geq 64$ . Therefore, for all  $\gamma \leq 64/(64+8) = 0.8889$ ,  $n_H(\gamma) \leq n_B(\gamma)$ .*

*Proof* The exact calculation of  $c$  is slightly nuanced, but a simple upper bound can be obtained as follows. The bound in Theorem 3 relies on (Agarwal et al. 2020, Theorem 1) (stated as Theorem 6 later in this paper). In the proof of (Agarwal et al. 2020, Theorem 1), the number of samples  $n$  are chosen such that  $(\gamma/(1-\gamma))\sqrt{8L'/n} \leq 1/2$ , where  $L' = \log(8|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-2}p^{-1})$  (also see the proof of Theorem 6). This implies that  $n \geq 64\gamma^2L'/(1-\gamma)^2$ . Comparing this with the bound on  $n$  in Theorem 3 (and, more precisely, the bound in Theorem 6 later), we get that  $c \geq 64$ . The bound on  $\gamma^\circ$  now follows from observing that  $c/(c+8)$  is an increasing function of  $c$ .  $\square$

*Remark 7* The result of Lemma 2 shows that for all  $\gamma \leq 0.8889$ , the sample complexity bound of Theorem 2 is tighter than the sample complexity bound of Theorem 3. Note that the actual value of the critical discount factor  $\gamma^\circ$  is higher, and depends on the values of  $|\mathcal{S}|$ ,  $|\mathcal{A}^i|$  and  $|\mathcal{N}|$ . See Sec. 5 for an example.

### 3 Background on MDPs

The main idea for proving the results of Sec. 2 is to pick a player, say  $i$ , fix the strategy profile of all players other than  $i$ , and then look at the best response of player  $i$ . The problem of finding the best response at player  $i$  is a single agent Markov decision process. So, we start by reviewing the pertinent results from Markov decision theory. All the results in this section are either standard or variations of existing results.

#### 3.1 MDP, Bellman Operators, and Dynamic Programming

A Markov Decision Process (MDP) is a tuple  $\langle \mathcal{S}, \mathcal{A}, \mathbf{P}, r, \gamma \rangle$  where

- $\mathcal{S}$  is the (finite) set of states of the environment. The state at time  $t$  is denoted by  $S_t$ .
- $\mathcal{A}$  is the (finite) set of actions available to the agent. The action at time  $t$  is denoted by  $A_t$ .
- $\mathbf{P} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$  is the controlled transition probability. For any realization  $s_{1:t+1}$  of  $S_{1:t+1}$  and choice  $a_{1:t}$  of  $A_{1:t}$ , we have

$$\mathbf{P}(s_{t+1}|s_t, a_t) := \mathbf{P}(S_{t+1} = s_{t+1} | S_{1:t} = s_{1:t}, A_{1:t} = a_{1:t}). \quad (14)$$

- $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the per-step reward function.
- $\gamma \in (0, 1)$  is the discount factor.

It is assumed that the agent observes the state  $S_t$  and chooses the action  $A_t$  according to a Markov strategy  $\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$ . The performance of a Markov



strategy  $\pi$  starting from initial state  $s_0 \in \mathcal{S}$  is given by:<sup>4</sup>

$$V_\pi(s_0) := (1 - \gamma) \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t r(S_t, A_t) \mid S_0 = s_0 \right], \quad (15)$$

where the expectation is with respect to the joint measure on the system variables induced by the choice of strategy  $\pi$ . A strategy  $\pi$  is called optimal if for any other Markov strategy  $\tilde{\pi}$ , we have

$$V_\pi(s) \geq V_{\tilde{\pi}}(s), \quad \forall s \in \mathcal{S}. \quad (16)$$

In addition, given a positive constant  $\alpha$ , a strategy  $\pi$  is called  $\alpha$ -optimal if

$$V_\pi(s) \geq V_{\tilde{\pi}}(s) - \alpha, \quad \forall s \in \mathcal{S}. \quad (17)$$

Given an MDP  $\mathcal{M} := (\mathcal{S}, \mathcal{A}, \mathbf{P}, r, \gamma)$  and a Markov strategy  $\pi$ , define the Bellman operators  $\mathcal{B}_\pi : \mathbb{R}^{|\mathcal{S}|} \rightarrow \mathbb{R}^{|\mathcal{S}|}$  and  $\mathcal{B}_* : \mathbb{R}^{|\mathcal{S}|} \rightarrow \mathbb{R}^{|\mathcal{S}|}$  as follows: for any  $v \in \mathbb{R}^{|\mathcal{S}|}$  and  $s \in \mathcal{S}$

$$[\mathcal{B}_\pi v](s) := \sum_{a \in \mathcal{A}} \pi(a|s) \left[ (1 - \gamma)r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathbf{P}(s'|s, a)v(s') \right], \quad (18)$$

$$[\mathcal{B}_* v](s) := \max_{a \in \mathcal{A}} \left[ (1 - \gamma)r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathbf{P}(s'|s, a)v(s') \right]. \quad (19)$$

Then, optimal and approximately optimal strategies can be characterized using the Bellman operators as shown below. These are standard results. See Bertsekas (2017), for example.

An immediate implication of Proposition 2 and the definition of an  $\alpha$ -optimal strategy is the following.

**Proposition 3** *A Markov strategy  $\pi$  is optimal if and only if there exists a value function  $V \in \mathbb{R}^{|\mathcal{S}|}$  such that*

$$V = \mathcal{B}_\pi V \quad \text{and} \quad V = \mathcal{B}_* V. \quad (20)$$

*Remark 8* Note that an MDP can have more than one optimal strategy but all optimal strategies have the same performance and hence the same value function.

**Proposition 4** *Given a Markov strategy  $\pi$ , let  $V_\pi$  be the unique fixed point of  $V_\pi = \mathcal{B}_\pi V_\pi$  and let  $V_*$  be the unique fixed point of  $V_* = \mathcal{B}_* V_*$ . Then, the strategy  $\pi$  is  $\alpha$ -optimal if and only if*

$$V_\pi \geq V_* - \alpha \mathbf{1}. \quad (21)$$

---

<sup>4</sup> For consistency with the normalized rewards considered in the game formulation (see Remark 1), we use normalized rewards for MDPs as well. Although most of the literature on MDPs uses unnormalized rewards, normalized rewards are commonly used in the literature on constrained MDPs (Altman 1999).

*Proof* This is an immediate implication of Proposition 3 and the definition of  $\alpha$ -optimal strategy. We provide a proof for completeness.

- ( $\Rightarrow$ ) Let  $\pi^*$  be an optimal policy and  $\pi$  be an  $\alpha$ -optimal strategy. Then, by definition of an  $\alpha$ -optimal strategy, we have  $V_\pi(s) \geq V_{\pi^*}(s) - \alpha$ , for all  $s \in \mathcal{S}$ . Moreover, since  $\pi^*$  is optimal, we have  $V_{\pi^*}(s) = V_*(s)$ . Combining the two, we have  $V_\pi(s) \geq V_*(s) - \alpha$ , which is same as (21).
- ( $\Leftarrow$ ) Let  $\pi$  be a policy which satisfies (21) and  $\tilde{\pi}$  be any policy. By definition of optimal policy, we have  $V_*(s) \geq V_{\tilde{\pi}}(s)$  for all  $s \in \mathcal{S}$ . Therefore, from (21), we have that  $V_\pi(s) \geq V_{\pi^*}(s) - \alpha \geq V_{\tilde{\pi}}(s) - \alpha$ . Since this inequality holds for every  $\tilde{\pi}$ , the policy  $\pi$  is  $\alpha$ -optimal.  $\square$

*Remark 9* A sufficient condition to verify (21) in Proposition 4 is that

$$V_\pi \geq \mathcal{B}_* V_\pi - \alpha \mathbf{1}. \quad (22)$$

We now present some basic properties of the value function which are used later.

**Lemma 3** *If  $V$  is the optimal value function of MDP  $\mathcal{M}$ , then*

$$\text{span}(V) \leq \text{span}(r).$$

*Proof* This result follows immediately by observing that the per-step reward  $r(S_t, A_t) \in [\min(r), \max(r)]$ . Therefore,  $\max(V) \leq \max(r)$  and  $\min(V) \geq \min(r)$ .  $\square$

We now define the notion of a Lipschitz MDP.

**Definition 7** Let  $d$  be a metric on the state space  $\mathcal{S}$ . The MDP  $\mathcal{M}$  is said to be  $(L_r, L_P)$ -Lipschitz if for any  $s_1, s_2 \in \mathcal{S}$  and  $a \in A$ , the reward function  $r$  and transition kernel  $P$  of  $\mathcal{M}$  satisfy the following

$$|r(s_1, a) - r(s_2, a)| \leq L_r d(s_1, s_2),$$

and

$$d_{\mathfrak{F}^w}(P(\cdot|s_1, a), P(\cdot|s_2, a)) \leq L_P d(s_1, s_2),$$

where  $d_{\mathfrak{F}^w}$  denotes the Wasserstein distance.

**Lemma 4** *If an MDP  $\mathcal{M}$  is  $(L_r, L_P)$ -Lipschitz, then for any policy  $\pi$ , the corresponding value function  $V_\pi$  is Lipschitz with*

$$\text{Lip}(V_\pi) \leq \frac{(1 - \gamma)L_r}{1 - \gamma(1 + \text{Lip}(\pi))L_P},$$

*provided  $\gamma(1 + \text{Lip}(\pi))L_P < 1$ , where  $\text{Lip}(\pi)$  is the Lipschitz-constant of the strategy  $\pi$ , i.e.,*

$$\text{Lip}(\pi) = \sup_{\substack{s, s' \in \mathcal{S} \\ s \neq s'}} \frac{d_{\mathfrak{F}^w}(\pi(\cdot|s), \pi(\cdot|s'))}{d(s, s')}.$$

*Proof* The result follows from (Hinderer 2005, Theorem 4.1).  $\square$

The above result can be strengthened when the policy  $\pi$  is the optimal policy.

**Lemma 5** *If an MDP  $\mathcal{M}$  is  $(L_r, L_P)$ -Lipschitz and  $\gamma L_P < 1$ , and  $V_*$  is the optimal value function of  $\mathcal{M}$ , then*

$$\text{Lip}(V_*) \leq \frac{(1 - \gamma)L_r}{1 - \gamma L_P}.$$

*Proof* The result follows from (Hinderer 2005, Theorem 4.2).  $\square$

### 3.2 Robustness of MDPs to model approximation

**Definition 8** Given a function class  $\mathfrak{F}$  and positive constants  $(\varepsilon, \delta)$ , we say that an MDP  $\widehat{\mathcal{M}} := \langle \mathcal{S}, \mathcal{A}, \widehat{P}, \widehat{r}, \gamma \rangle$  is an  $(\varepsilon, \delta)$ -approximation of the MDP  $\mathcal{M} := \langle \mathcal{S}, \mathcal{A}, P, r, \gamma \rangle$  if it satisfies the following properties:

1. **Reward approximation:** For all  $s \in \mathcal{S}$ , and  $a \in \mathcal{A}$ ,

$$|r(s, a) - \widehat{r}(s, a)| \leq \varepsilon. \quad (23)$$

2. **Transition approximation:** For all  $s \in \mathcal{S}$ , and  $a \in \mathcal{A}$ ,

$$d_{\mathfrak{F}}(P(\cdot|s, a), \widehat{P}(\cdot|s, a)) \leq \delta. \quad (24)$$

The main approximation result for MDPs relevant for our analysis is the following.

**Theorem 4** *Given a function class  $\mathfrak{F}$  and an MDP  $\mathcal{M} := \langle \mathcal{S}, \mathcal{A}, P, r, \gamma \rangle$ , suppose  $\widehat{\mathcal{M}} := \langle \mathcal{S}, \mathcal{A}, \widehat{P}, \widehat{r}, \gamma \rangle$  is an  $(\varepsilon, \delta)$ -approximation of  $\mathcal{M}$ . Let  $\hat{\pi}_*$  and  $\hat{\pi}$  be an optimal and an  $\alpha_{\text{opt}}$ -optimal strategy of  $\widehat{\mathcal{M}}$ . Let  $\hat{V}_{\hat{\pi}_*}$  and  $\hat{V}_{\hat{\pi}}$  be the corresponding value functions. Then  $\hat{\pi}$  is an  $\alpha$ -optimal strategy of  $\mathcal{M}$  with*

$$\alpha \leq 2\varepsilon + \frac{\gamma}{1 - \gamma} [\Delta_{\hat{\pi}_*} + \Delta_{\hat{\pi}}] + \alpha_{\text{opt}}, \quad (25)$$

where

$$\Delta_{\hat{\pi}} := \max_{s \in \mathcal{S}, a \in \mathcal{A}} \left| \sum_{s' \in \mathcal{S}} \left[ P(s'|s, a) \hat{V}_{\hat{\pi}}(s') - \widehat{P}(s'|s, a) \hat{V}_{\hat{\pi}}(s') \right] \right|$$

and  $\Delta_{\hat{\pi}_*}$  is defined similarly. Furthermore, we can show that

$$\Delta_{\hat{\pi}} \leq \delta \rho_{\mathfrak{F}}(\hat{V}_{\hat{\pi}})$$

and a similar bound holds for  $\Delta_{\hat{\pi}_*}$ .

*Proof* The bound for (25) is established in Appendix A. The upper bound on  $\Delta_{\hat{\pi}}$  and  $\Delta_{\hat{\pi}_*}$  follows from (5) and definition of  $\delta$ .  $\square$

Note that the results of Theorem 4 are instance dependent. We present instance independent bounds on  $\alpha$  by upper bounding  $\rho_{\mathfrak{F}}(V)$ .

### 3.2.1 Instance independent bounds

**Corollary 5** *If the function class  $\mathfrak{F}$  in Theorem 4 is  $\mathfrak{F}^{\text{TV}}$ , then*

$$\alpha \leq 2\left(\varepsilon + \frac{\gamma\delta \text{span}(\hat{r})}{(1-\gamma)}\right) + \alpha_{\text{opt}}.$$

*Proof* The result follows from the observation that  $\rho_{d_{\mathfrak{F}^{\text{TV}}}}(\hat{V}) = \text{span}(\hat{V})$  and then using Lemma 3 in Theorem 4.  $\square$

**Corollary 6** *If the function class  $\mathfrak{F}$  in Theorem 4 is  $\mathfrak{F}^{\text{W}}$ , and the approximate MDP  $\widehat{\mathcal{M}}$  is  $(L_r, L_P)$ -Lipschitz with  $\gamma L_P < 1$ , then*

$$\alpha \leq 2\varepsilon + \frac{\gamma\delta L_r}{(1-\gamma L_P)} + \frac{\gamma\delta L_r}{(1-\gamma(1+\text{Lip}(\hat{\pi}))L_P)} + \alpha_{\text{opt}}.$$

*Furthermore, if  $\alpha_{\text{opt}} = 0$ , the above expression simplifies to*

$$\alpha \leq 2\left(\varepsilon + \frac{\gamma\delta L_r}{(1-\gamma L_P)}\right).$$

*Proof* The result follows from the observation that  $\rho_{\mathfrak{F}^{\text{W}}}(\hat{V}) = \text{Lip}(\hat{V})$  and then using Lemmas 4 and 5 in Theorem 4.  $\square$

### 3.3 Model based RL for MDPs

In this section, we consider a setting similar to Sec. 2.4, but for MDPs. Suppose the components  $\langle \mathcal{S}, \mathcal{A}, r, \gamma \rangle$  of an MDP  $\mathcal{M}$  are known but the transition probability matrix  $P$  is not known. Similar to the game setting, we assume that we have access to a *generative model*, i.e., a black-box simulator which provides samples  $S_+ \sim P(\cdot|s, a)$  of the next state  $S_+$  for a given state-action pair  $(s, a)$  as input. Suppose we call the simulator  $n$  times at each state-action pair and estimate an empirical model  $\widehat{P}_n$  as  $\widehat{P}_n(s'|s, a) := \text{count}(s'|s, a)/n$ , where  $\text{count}(s'|s, a)$  is the number of times  $s'$  is sampled with the input is  $(s, a)$ . The MDP  $\widehat{\mathcal{M}}_n := \langle \mathcal{S}, \mathcal{A}, \widehat{P}_n, r, \gamma \rangle$  may be viewed as an approximation of MDP  $\mathcal{M}$ . As in the game setting, we assume that there is a planning oracle, which takes the approximate model  $\widehat{\mathcal{M}}_n$  as input and generates an  $\alpha_{\text{opt}}$ -optimal strategy  $\hat{\pi}_n$ , where  $\alpha_{\text{opt}} \in \mathbb{R}_{>0}$  is a property of the planning oracle.

As in Sec. 2.4, we are interested in the following question. Given an  $\alpha > 0$ , how many samples  $n$  are needed from the generative model to ensure that  $\hat{\pi}_n$  is an  $\alpha$ -optimal MDP. This is called the *sample complexity* of learning and a simple upper bound can be obtained using the approximation bounds of Theorem 4. For that matter, we state standard results on concentration inequalities.

We state two bounds on the sample complexity of generative models: one based on Hoeffding inequality and the other based on Bernstein inequality.

Normally both these bounds are stated in terms of  $\|r\|_\infty$  but for the first bound, we provide a slightly tighter bound in terms of  $\text{span}(r)$ . For the sake of completeness, we provide a complete proof of the first bound.

Suppose  $X \in \mathcal{X}$  is a random variable with distribution  $\mu$ . Suppose  $\{X_1, \dots, X_n\}$  is a sequence of random variables sampled according to  $\mu$ . Let  $\hat{\mu}_n$  denote the empirical measure constructed from  $\{X_1, \dots, X_n\}$ , i.e.,

$$\hat{\mu}_n(x) := \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{\{X_k=x\}},$$

where  $\mathbb{1}_{\{E\}}$  denotes the indicator function of the event  $E$ . Then, we have the following.

**Lemma 6** *For a given  $H > 0$ , let  $\mathfrak{F}_H$  denote the set of functions  $f: \mathcal{X} \rightarrow \mathbb{R}$  such that  $\text{span}(f) \leq H$ . Then, for any  $f \in \mathfrak{F}_H$  and  $\Delta > 0$ ,*

$$\mathbb{P}\left(\left|\sum_{x \in \mathcal{X}} [\mu(x)f(x) - \mu_n(x)f(x)]\right| \geq \Delta\right) \leq 2 \exp\left(-\frac{2n\Delta^2}{H^2}\right).$$

*Proof* Let  $Z := f(X)$  and  $Z_i := f(X_i)$ . Then,  $\{Z_1, \dots, Z_n\}$  is an i.i.d. sequence and  $\text{Supp}(Z_i) \leq H$ . Then, by the Hoeffding inequality (Cesa-Bianchi and Lugosi 2006, Corollary A.1),

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n Z_i - \mathbb{E}[Z]\right| \geq \Delta\right) \leq 2 \exp\left(-\frac{2n\Delta^2}{H^2}\right).$$

The result then follows from observing that  $\mathbb{E}[Z] = \sum_{x \in \mathcal{X}} \mu(x)f(x)$  and  $(\sum_{i=1}^n Z_i)/n = \sum_{x \in \mathcal{X}} \mu_n(x)f(x)$ .  $\square$

By combining Lemma 6 with Theorem 4, we get the following.

**Theorem 5 (Hoeffding-Type Bound)** *For any  $\alpha_{\text{opt}} > 0$ ,  $\alpha > \alpha_{\text{opt}}$  and  $p > 0$ , let*

$$n \geq \left\lceil \left(\frac{\gamma}{1-\gamma} \text{span}(r)\right)^2 \frac{2 \log(4|\mathcal{S}| |\mathcal{A}| p^{-1})}{(\alpha - \alpha_{\text{opt}})^2} \right\rceil.$$

*Then, with probability  $1 - p$ , the  $\alpha_{\text{opt}}$ -optimal strategy  $\hat{\pi}_n$  of MDP  $\widehat{\mathcal{M}}_n$  is  $(\alpha + \alpha_{\text{opt}})$ -optimal for  $\mathcal{M}$ .*

*Proof* Let  $H := \text{span}(r)$ . Let  $\hat{\pi}_{*,n}$  and  $\hat{\pi}_n$  be an optimal and  $\alpha_{\text{opt}}$ -optimal policies of  $\widehat{\mathcal{M}}_n$ . Let  $\hat{V}_{\hat{\pi}_{*,n}}$  and  $\hat{V}_{\hat{\pi}_n}$  be the corresponding value functions for  $\widehat{\mathcal{M}}_n$ . From Lemma 3, we know that  $\hat{V}_{\hat{\pi}_{*,n}}, \hat{V}_{\hat{\pi}_n} \in \mathfrak{F}_H$ . Therefore, from Lemma 6, we have that for a given state-action pair  $(s, a)$  and  $\Delta > 0$ ,

$$\mathbb{P}\left(\left|\sum_{s' \in \mathcal{S}} \mathbb{P}(s'|s, a) \hat{V}_{\hat{\pi}_{*,n}}(s') - \sum_{s' \in \mathcal{S}} \hat{\mathbb{P}}(s'|s, a) \hat{V}_{\hat{\pi}_n}(s')\right| \geq \Delta\right) \leq 2 \exp\left(-\frac{2n\Delta^2}{H^2}\right)$$

and

$$\mathbb{P}\left(\left|\sum_{s' \in \mathcal{S}} \mathbb{P}(s'|s, a) \hat{V}_{\hat{\pi}_n}(s') - \sum_{s' \in \mathcal{S}} \hat{\mathbb{P}}(s'|s, a) \hat{V}_{\hat{\pi}_n}(s')\right| \geq \Delta\right) \leq 2 \exp\left(-\frac{2n\Delta^2}{H^2}\right).$$

Therefore, by the union bound,

$$\begin{aligned} & \mathbb{P}\left(\left\{\max_{(s,a) \in \mathcal{S} \times \mathcal{A}} \left|\sum_{s' \in \mathcal{S}} \mathbb{P}(s'|s, a) \hat{V}_{\hat{\pi}_{*,n}}(s') - \sum_{s' \in \mathcal{S}} \hat{\mathbb{P}}(s'|s, a) \hat{V}_{\hat{\pi}_{*,n}}(s')\right| \geq \Delta\right\}\right. \\ & \quad \left.\text{and } \left\{\max_{(s,a) \in \mathcal{S} \times \mathcal{A}} \left|\sum_{s' \in \mathcal{S}} \mathbb{P}(s'|s, a) \hat{V}_{\hat{\pi}_n}(s') - \sum_{s' \in \mathcal{S}} \hat{\mathbb{P}}(s'|s, a) \hat{V}_{\hat{\pi}_n}(s')\right| \geq \Delta\right\}\right) \\ & \leq 4|\mathcal{S}||\mathcal{A}| \exp\left(-\frac{2n\Delta^2}{H^2}\right). \end{aligned} \quad (26)$$

Now, choose  $\Delta$  such that the right hand side of (26) equals  $p$ , i.e.,

$$\Delta := H \sqrt{\frac{\log(4|\mathcal{S}||\mathcal{A}|p^{-1})}{2n}}.$$

Then, Theorem 4 implies that with probability  $1 - p$ ,  $\hat{\pi}_n$  is a  $\alpha_n$ -optimal strategy for  $\mathcal{M}$ , where

$$\alpha_n = 2 \frac{\gamma \Delta}{1 - \gamma} + \alpha_{\text{opt}} = \frac{\gamma}{1 - \gamma} H \sqrt{\frac{2 \log(4|\mathcal{S}||\mathcal{A}|p^{-1})}{n}} + \alpha_{\text{opt}}.$$

The result now follows by substituting the value of  $n$ .  $\square$

*Remark 10* The sample complexity bound in Theorem 5 is not tight. It is shown in Azar et al. (2013) that finding an  $\alpha$ -optimal policy with probability  $1 - p$  requires at least

$$\Omega\left(|\mathcal{S}||\mathcal{A}| \frac{\log(|\mathcal{S}||\mathcal{A}|p^{-1})}{(1 - \gamma)\alpha^2}\right)$$

samples.<sup>5</sup> The upper bound in Theorem 5 is loose by a factor of  $1/(1 - \gamma)$ . For the case of MDPs, tighter upper bounds which match the lower bound of Azar et al. (2013) (up to logarithmic factors) have been obtained in Sidford et al. (2018); Agarwal et al. (2020) by using the Bernstein inequality rather than the Hoeffding inequality. We present these bounds below.

**Theorem 6** For any  $\alpha_{\text{opt}} > 0$ ,  $\alpha \in (5\alpha_{\text{opt}}, \|r\|_{\infty} \sqrt{(1 - \gamma)} + 5\alpha_{\text{opt}})$  and  $p > 0$ , let

$$n \geq \frac{c\gamma \|r\|_{\infty}^2 \log(8|\mathcal{S}||\mathcal{A}|(1 - \gamma)^{-2}p^{-1})}{(1 - \gamma)(\alpha - 5\alpha_{\text{opt}})^2},$$

where  $c$  is an absolute constant. Then, with probability  $1 - p$ , any  $\alpha_{\text{opt}}$ -optimal strategy  $\hat{\pi}_n$  of MDP  $\widehat{\mathcal{M}}_n$  is  $\alpha$ -optimal for  $\mathcal{M}$ .

<sup>5</sup> Recall that we are working with normalized total expected reward (see Remark 1), while the results Azar et al. (2013) are derived for the unnormalized total reward. In the discussion above, we have normalized the results of Azar et al. (2013).

*Proof* The proof follows from (Agarwal et al. 2020, Theorem 1) with appropriate scaling to account for normalization of rewards. However, the expression for  $n$  above differs slightly from the expression in (Agarwal et al. 2020, Theorem 1). In particular, let  $L = \log(8|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-1}p^{-1})$  and  $L' = \log(8|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-2}p^{-1})$  (the exponent of  $(1-\gamma)$  is different in the two expressions). The expression in (Agarwal et al. 2020, Theorem 1) (rescaled for normalized rewards) states that

$$n \geq \frac{c\gamma\|r\|_\infty^2 L}{(1-\gamma)\alpha^2},$$

while the expression in the statement of Theorem 6 states that

$$n \geq \frac{c\gamma\|r\|_\infty^2 L'}{(1-\gamma)\alpha^2}.$$

The reason for this difference is that there is a typo in (Agarwal et al. 2020, Lemma 11), which is carried over in all the results. In particular, in the proof of (Agarwal et al. 2020, Lemma 11) it is claimed that for  $|U_s| = 1/(1-\gamma)^2$  and  $p' = p/(2|\mathcal{S}||\mathcal{A}|)$  (in Agarwal et al. (2020), the symbol  $\delta$  is used instead of  $p$ ), we have  $\log(4|U_s|/p') = L$  but elementary algebra shows that  $\log(4|U_s|/p') = L'$ .

Also note that in Agarwal et al. (2020), the bound of Theorem 3 was simplified as

$$n \geq \frac{c\gamma\|r\|_\infty^2 \log(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-2}p^{-1})}{(1-\gamma)(\alpha - 5\alpha_{\text{opt}})^2},$$

i.e., the multiplicative factor of 8 inside the log was removed. Since we want to compare the Bernstein-type bound with the Hoeffding-type bound in Sec. 2.5.2, we carry the multiplicative factor of 8 in our expression.  $\square$

## 4 Proof of the main results

### 4.1 Road map of the proof

As mentioned in the beginning of Sec. 3, the main idea of the proofs is to look at the best response of a player to the pre-specified strategy profile of other players. To formally establish that the problem of finding the best response is an MDP, we first start by characterizing the best response in terms of Bellman operators and stating existing results that characterize MPE and approximate MPE in terms of Bellman operators (Propositions 5 and 6). Then, we formally define a “best response MDP” and show that MPE and approximate MPE can be stated in terms of such “best response MDPs” (Corollaries 7 and 8).

To establish the robustness results, we show that the “best response MDP” corresponding to an  $(\varepsilon, \delta)$ -approximation of a game is an  $(\varepsilon, \delta)$ -approximation of the “best response MDP” of the original game (Lemma 8). This allows us to generalize the approximation results of MDPs to games. We then build on this relationship to generalize the sample complexity results of MDPs to games.

#### 4.2 Bellman operators and characterization of Markov perfect equilibrium

Given a Markov strategy profile  $\pi := (\pi^i)_{i \in \mathcal{N}}$ , state  $s \in \mathcal{S}$ , and action profile  $a := (a^i)_{i \in \mathcal{N}} \in \mathcal{A}$ , we use the notation

$$\begin{aligned} \pi(a|s) &:= \prod_{i \in \mathcal{N}} \pi^i(a^i|s) \quad \text{and} \\ \pi^{-i}(a^{-i}|s) &:= \prod_{j \in \mathcal{N} \setminus \{i\}} \pi^j(a^j|s). \end{aligned} \quad (27)$$

Given a player  $i \in \mathcal{N}$  and a Markov strategy profile  $\pi := (\pi^i, \pi^{-i})$ , we define two Bellman operators as follows:

1. An operator  $\mathcal{B}_{(\pi^i, \pi^{-i})}^i : \mathbb{R}^{|\mathcal{S}|} \rightarrow \mathbb{R}^{|\mathcal{S}|}$  given as follows: for any  $v \in \mathbb{R}^{|\mathcal{S}|}$  and  $s \in \mathcal{S}$ ,

$$[\mathcal{B}_{(\pi^i, \pi^{-i})}^i v](s) := \sum_{a \in \mathcal{A}} \pi(a|s) \left[ (1 - \gamma) r^i(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) v(s') \right].$$

2. An operator  $\mathcal{B}_{*, \pi^{-i}}^i : \mathbb{R}^{|\mathcal{S}|} \rightarrow \mathbb{R}^{|\mathcal{S}|}$  given as follows: for any  $v \in \mathbb{R}^{|\mathcal{S}|}$  and  $s \in \mathcal{S}$ ,

$$[\mathcal{B}_{*, \pi^{-i}}^i v](s) := \max_{a^i \in \mathcal{A}^i} \left[ \sum_{a^{-i} \in \mathcal{A}^{-i}} \pi^{-i}(a^{-i}|s) \left[ (1 - \gamma) r^i(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) v(s') \right] \right].$$

Now, MPE and approximate MPE can be characterized using the Bellman operators. These are standard results. See, for example, Filar and Vrieze (1996).

**Proposition 5** *A Markov strategy profile  $\pi := (\pi^i)_{i \in \mathcal{N}}$  is an MPE if and only if there exist **value functions**  $V^i \in \mathbb{R}^{|\mathcal{S}|}$ ,  $i \in \mathcal{N}$ , such that*

$$V^i = \mathcal{B}_{(\pi^i, \pi^{-i})}^i V^i \quad \text{and} \quad V^i = \mathcal{B}_{*, \pi^{-i}}^i V^i, \quad \forall i \in \mathcal{N}. \quad (28)$$

An immediate consequence of Proposition 5 and the definition of approximation MPE is the following.

**Proposition 6** *Given a Markov strategy profile  $\pi := (\pi^i)_{i \in \mathcal{N}}$ , for any  $i \in \mathcal{N}$ , let  $V_\pi^i$  be the unique fixed point of  $V_\pi^i = \mathcal{B}_{(\pi^i, \pi^{-i})}^i V_\pi^i$  and let  $V_{(*, \pi^{-i})}^i$  be the unique fixed point of  $V_{(*, \pi^{-i})}^i = \mathcal{B}_{(*, \pi^{-i})}^i V_{(*, \pi^{-i})}^i$ . Then, the strategy profile  $\pi$  is an  $\alpha$ -MPE,  $\alpha := (\alpha^i)_{i \in \mathcal{N}}$ , if and only if*

$$V_\pi^i \geq V_{(*, \pi^{-i})}^i - \alpha^i \mathbf{1}, \quad \forall i \in \mathcal{N}. \quad (29)$$

*Proof* The proof follows from arguments similar to the proof of Prop. 4.



### 4.3 Relationship between games and MDPs

Given a game  $\mathcal{G} := \langle \mathcal{N}, \mathcal{S}, (\mathcal{A}^i)_{i \in \mathcal{N}}, \mathbf{P}, (r^i)_{i \in \mathcal{N}}, \gamma \rangle$  and a Markov strategy  $\pi := (\pi^i)_{i \in \mathcal{N}}$ , we can define MDPs  $\{\mathcal{M}_{\pi^{-i}}^i\}_{i \in \mathcal{N}}$  as follows. For player  $i \in \mathcal{N}$ , MDP  $\mathcal{M}_{\pi^{-i}}^i := \langle \mathcal{S}, \mathcal{A}^i, \mathbf{P}_{\pi^{-i}}^i, r_{\pi^{-i}}^i, \gamma \rangle$ , where the transition matrix  $\mathbf{P}_{\pi^{-i}}^i : \mathcal{S} \times \mathcal{A}^i \rightarrow \mathcal{P}(\mathcal{S})$  is given by

$$\mathbf{P}_{\pi^{-i}}^i(s'|s, a^i) := \sum_{a^{-i} \in \mathcal{A}^{-i}} \pi^{-i}(a^{-i}|s) \mathbf{P}(s'|s, (a^i, a^{-i})), \quad (30)$$

and the reward function  $r_{\pi^{-i}}^i : \mathcal{S} \times \mathcal{A}^i \rightarrow \mathbb{R}$  is given by

$$r_{\pi^{-i}}^i(s, a^i) := \sum_{a^{-i} \in \mathcal{A}^{-i}} \pi^{-i}(a^{-i}|s) r^i(s, (a^i, a^{-i})). \quad (31)$$

In other words, in  $\mathcal{M}_{\pi^{-i}}^i$ , the strategy of player  $i$  may be chosen freely while the strategies of all other players are fixed at those specified in  $\pi^{-i}$ . Note the Bellman operators  $\mathcal{B}_{(\pi^i, \pi^{-i})}^i$  and  $\mathcal{B}_{(*, \pi^{-i})}^i$  corresponding to game  $\mathcal{G}$  and strategy  $\pi$  are the same as Bellman operators of MDP  $\mathcal{M}_{\pi^{-i}}^i$ . Therefore, by combining Propositions 3 and 5, we have the following:

**Corollary 7** *A Markov strategy profile  $\pi := (\pi^i)_{i \in \mathcal{N}}$  is an MPE if and only if for every  $i \in \mathcal{N}$ , the strategy  $\pi^i$  is an optimal strategy for MDP  $\mathcal{M}_{\pi^{-i}}^i$ .*

*Proof* This is an immediate consequence of the definition of  $\mathcal{M}_{\pi^{-i}}^i$ . For the sake of completeness, we provide a formal proof. Arbitrarily pick a player  $i \in \mathcal{N}$  and consider any Markov policy  $\tilde{\pi}^i$  for player  $i$ . The Bellman operator  $\mathcal{B}_{(\tilde{\pi}^i, \pi^{-i})}^i$  of game  $\mathcal{G}$  is the same as the Bellman operator for evaluating policy  $\tilde{\pi}^i$  in MDP  $\mathcal{M}_{\pi^{-i}}^i$ . Thus, the value function  $V_{\tilde{\pi}^i, \pi^{-i}}$  (which is the fixed point of  $\mathcal{B}_{(\tilde{\pi}^i, \pi^{-i})}^i$ ) is equal to the value of policy  $\tilde{\pi}^i$  in MDP  $\mathcal{M}_{\pi^{-i}}^i$ . Now, we prove the two directions separately.

- ( $\Rightarrow$ ) Suppose  $\pi$  is an MPE of game  $\mathcal{G}$ . By the definition of MPE, for any player  $i \in \mathcal{N}$  and any policy  $\tilde{\pi}^i$ , we have  $V_{(\pi^i, \pi^{-i})}(s) \geq V_{(\tilde{\pi}^i, \pi^{-i})}(s)$ , for all  $s \in \mathcal{S}$ . This means that in MDP  $\mathcal{M}_{\pi^{-i}}^i$ , the performance of policy  $\pi^i$  is at least as good as the performance of any other policy  $\tilde{\pi}^i$ . Hence, policy  $\pi^i$  is optimal for MDP  $\mathcal{M}_{\pi^{-i}}^i$ .
- ( $\Leftarrow$ ) Suppose for all player  $i \in \mathcal{N}$ , the policy  $\pi^i$  is optimal for MDP  $\mathcal{M}_{\pi^{-i}}^i$ . This means that for any other policy  $\tilde{\pi}^i$  for player  $i$ , the performance of policy  $\pi^i$  in MDP  $\mathcal{M}_{\pi^{-i}}^i$  is at least as good as the performance of policy  $\tilde{\pi}^i$ . Thus, we have  $V_{(\pi^i, \pi^{-i})}(s) \geq V_{(\tilde{\pi}^i, \pi^{-i})}(s)$ , for all  $s \in \mathcal{S}$ . Since this is true for every player  $i \in \mathcal{N}$ , the policy  $\pi$  is an MPE.  $\square$

Similarly, by combining Propositions 4 and 6, we have the following:

**Corollary 8** *Given approximate levels  $\alpha := (\alpha^i)_{i \in \mathcal{N}}$ ,  $\alpha^i \in \mathbb{R}_{\geq 0}$ , a Markov strategy profile  $\pi := (\pi^i)_{i \in \mathcal{N}}$ , is an  $\alpha$ -MPE if and only if for every  $i \in \mathcal{N}$ , the strategy  $\pi^i$  is an  $\alpha^i$ -optimal strategy for MDP  $\mathcal{M}_{\pi^{-i}}^i$ .*

*Proof* The proof argument is almost the same as the proof of Corollary 7. As argued in the proof of Corollary 7, the value function  $V_{\tilde{\pi}^i, \pi^{-i}}$  (which is the fixed point of  $\mathcal{B}_{(\tilde{\pi}^i, \pi^{-i})}^i$ ) is equal to the value of policy  $\tilde{\pi}^i$  in MDP  $\mathcal{M}_{\pi^{-i}}^i$ . Now, we prove the two directions separately.

- ( $\Rightarrow$ ) Suppose  $\pi$  is an  $\alpha$ -MPE of game  $\mathcal{G}$ . By the definition of MPE, for any player  $i \in \mathcal{N}$  and any policy  $\tilde{\pi}^i$ , we have  $V_{(\pi^i, \pi^{-i})}(s) \geq V_{(\tilde{\pi}^i, \pi^{-i})}(s) - \alpha^i$ , for all  $s \in \mathcal{S}$ . This means that in MDP  $\mathcal{M}_{\pi^{-i}}^i$ , the performance of policy  $\pi^i$  is at least as good as the performance of any other policy  $\tilde{\pi}^i$  minus  $\alpha^i$ . Hence, policy  $\pi^i$  is  $\alpha^i$ -optimal for MDP  $\mathcal{M}_{\pi^{-i}}^i$ .
- ( $\Leftarrow$ ) Suppose for all player  $i \in \mathcal{N}$ , the policy  $\pi^i$  is  $\alpha^i$ -optimal for MDP  $\mathcal{M}_{\pi^{-i}}^i$ . This means that for any other policy  $\tilde{\pi}^i$  for player  $i$ , the performance of policy  $\pi^i$  in MDP  $\mathcal{M}_{\pi^{-i}}^i$  is at least as good as the performance of policy  $\tilde{\pi}^i$  minus  $\alpha^i$ . Thus, we have  $V_{(\pi^i, \pi^{-i})}(s) \geq V_{(\tilde{\pi}^i, \pi^{-i})}(s) - \alpha^i$ , for all  $s \in \mathcal{S}$ . Since this is true for every player  $i \in \mathcal{N}$ , the policy  $\pi$  is an  $\alpha$ -MPE.  $\square$

#### 4.4 Relationship between MDPs corresponding to a strategy profile

We first provide a preliminary result.

**Lemma 7** For any function  $f: \mathcal{S} \rightarrow \mathbb{R}$ , transitions  $P, \hat{P}: \mathcal{S} \times (\mathcal{A}^i)_{i \in \mathcal{N}} \rightarrow \Delta(\mathcal{S})$ , player  $i \in \mathcal{N}$ , strategy  $\pi^{-i}$  for players other than  $i$ ,  $(s, a^i) \in \mathcal{S} \times \mathcal{A}^i$  and transitions  $P_{\pi^{-i}}, \hat{P}_{\pi^{-i}}: \mathcal{S}^i \times \mathcal{A}^i \rightarrow \Delta(\mathcal{S}^i)$  defined as in (30), we have

$$\begin{aligned} & \left| \sum_{s' \in \mathcal{S}} f(s') P_{\pi^{-i}}^i(s'|s, a^i) - \sum_{s' \in \mathcal{S}} f(s') \hat{P}_{\pi^{-i}}^i(s'|s, a^i) \right| \\ & \leq \max_{a^{-i} \in \mathcal{A}^{-i}} \left| \sum_{s' \in \mathcal{S}} f(s') P(s'|s, (a^i, a^{-i})) - \sum_{s' \in \mathcal{S}} f(s') \hat{P}(s'|s, (a^i, a^{-i})) \right|. \end{aligned}$$

Therefore,

$$\begin{aligned} & \max_{s \in \mathcal{S}, a^i \in \mathcal{A}^i} \left| \sum_{s' \in \mathcal{S}} f(s') P_{\pi^{-i}}^i(s'|s, a^i) - \sum_{s' \in \mathcal{S}} f(s') \hat{P}_{\pi^{-i}}^i(s'|s, a^i) \right| \\ & \leq \max_{s \in \mathcal{S}, (a^i, a^{-i}) \in \mathcal{A}} \left| \sum_{s' \in \mathcal{S}} f(s') P(s'|s, (a^i, a^{-i})) - \sum_{s' \in \mathcal{S}} f(s') \hat{P}(s'|s, (a^i, a^{-i})) \right|. \end{aligned}$$

*Proof* For the first part, from definition of  $\hat{P}_{\pi^{-i}}^i$ , we have

$$\begin{aligned} & \left| \sum_{s' \in \mathcal{S}} f(s') P_{\pi^{-i}}^i(s'|s, a^i) - \sum_{s' \in \mathcal{S}} f(s') \hat{P}_{\pi^{-i}}^i(s'|s, a^i) \right| \\ & = \left| \sum_{s' \in \mathcal{S}} \sum_{a^{-i} \in \mathcal{A}^{-i}} f(s') \pi^{-i}(a^{-i}|s) P(s'|s, (a^i, a^{-i})) \right. \\ & \quad \left. - \sum_{s' \in \mathcal{S}} \sum_{a^{-i} \in \mathcal{A}^{-i}} f(s') \pi^{-i}(a^{-i}|s) \hat{P}(s'|s, (a^i, a^{-i})) \right| \end{aligned}$$

$$\begin{aligned}
&\leq \left| \sum_{a^{-i} \in \mathcal{A}^{-i}} \pi^{-i}(a^{-i}|s) \right. \\
&\quad \times \left. \left[ \sum_{s' \in \mathcal{S}} f(s') (P(s'|s, (a^i, a^{-i})) - \hat{P}(s'|s, (a^i, a^{-i}))) \right] \right| \\
&\leq \sum_{a^{-i} \in \mathcal{A}^{-i}} \pi^{-i}(a^{-i}|s) \\
&\quad \times \left| \sum_{s' \in \mathcal{S}} f(s') (P(s'|s, (a^i, a^{-i})) - \hat{P}(s'|s, (a^i, a^{-i}))) \right| \\
&\leq \sum_{a^{-i} \in \mathcal{A}^{-i}} \pi^{-i}(a^{-i}|s) \\
&\quad \times \max_{\tilde{a}^{-i} \in \mathcal{A}^{-i}} \left| \sum_{s' \in \mathcal{S}} f(s') (P(s'|s, (a^i, \tilde{a}^{-i})) - \hat{P}(s'|s, (a^i, \tilde{a}^{-i}))) \right| \\
&= \max_{\tilde{a}^{-i} \in \mathcal{A}^{-i}} \left| \sum_{s' \in \mathcal{S}} f(s') (P(s'|s, (a^i, \tilde{a}^{-i})) - \hat{P}(s'|s, (a^i, \tilde{a}^{-i}))) \right|.
\end{aligned}$$

The second part following by taking a maximum over  $(s, a^i)$ .  $\square$

Suppose we are given a game  $\mathcal{G}$  and its  $(\varepsilon, \delta)$  approximation  $\hat{\mathcal{G}}$ . Moreover, suppose  $\hat{\pi} := (\hat{\pi}^i)_{i \in \mathcal{N}}$  is an MPE of  $\hat{\mathcal{G}}$ .

Let  $\{\hat{\mathcal{M}}_{\hat{\pi}^{-i}}^i\}$  be the MDPs corresponding to game  $\hat{\mathcal{G}}$  and strategy  $\hat{\pi}$ . Similarly, let  $\{\mathcal{M}_{\hat{\pi}^{-i}}^i\}$  be the MDPs corresponding to game  $\mathcal{G}$  and strategy  $\hat{\pi}$ . An immediate implication of Lemma 7 is the following.

**Lemma 8** *For any player  $i \in \mathcal{N}$ , MDP  $\hat{\mathcal{M}}_{\hat{\pi}^{-i}}^i$  is an  $(\varepsilon, \delta)$  approximation of MDP  $\mathcal{M}_{\hat{\pi}^{-i}}^i$ .*

*Proof* Consider

$$\begin{aligned}
&|r_{\hat{\pi}^{-i}}^i(s, a^i) - \hat{r}_{\hat{\pi}^{-i}}^i(s, a^i)| \\
&\stackrel{(a)}{\leq} \sum_{a^{-i} \in \mathcal{A}^{-i}} \hat{\pi}^{-i}(a^{-i}|s) |r^i(s, (a^i, a^{-i})) - \hat{r}^i(s, (a^i, a^{-i}))| \\
&\stackrel{(b)}{\leq} \sum_{a^{-i} \in \mathcal{A}^{-i}} \hat{\pi}^{-i}(a^{-i}|s) \varepsilon \\
&\stackrel{(c)}{=} \varepsilon,
\end{aligned} \tag{32}$$

where (a) follows from (31), (b) follows from (6) and (c) follows as  $\varepsilon$  is independent of  $a^{-i}$ . Furthermore,

$$\begin{aligned}
&\max_{s \in \mathcal{S}, a^i \in \mathcal{A}^i} d_{\mathcal{F}}(P_{\hat{\pi}^{-i}}^i(\cdot|s, a^i), \hat{P}_{\hat{\pi}^{-i}}^i(\cdot|s, a^i)) \\
&\stackrel{(d)}{=} \sup_{f \in \mathcal{F}} \max_{s \in \mathcal{S}, a^i \in \mathcal{A}^i} \left| \sum_{s' \in \mathcal{S}} f(s') P_{\hat{\pi}^{-i}}^i(s'|s, a^i) - \sum_{s' \in \mathcal{S}} f(s') \hat{P}_{\hat{\pi}^{-i}}^i(s'|s, a^i) \right|
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(e)}{\leq} \sup_{f \in \mathcal{F}} \max_{\substack{s \in \mathcal{S} \\ (a^i, a^{-i}) \in \mathcal{A}}} \left| \sum_{s' \in \mathcal{S}} f(s') (P(s'|s, (a^i, a^{-i})) - \hat{P}(s'|s, (a^i, a^{-i}))) \right| \\
&\stackrel{(f)}{=} \max_{\substack{s \in \mathcal{S} \\ (a^i, a^{-i}) \in \mathcal{A}}} d_{\mathcal{F}}(P(\cdot|s, (a^i, a^{-i})), \hat{P}(\cdot|s, (a^i, a^{-i}))) \\
&\stackrel{(g)}{=} \delta
\end{aligned} \tag{33}$$

where (d) and (f) follows from Definition 3, (e) follows from Lemma 7, and (g) follows from Definition 4.

Equations (32) and (33) imply that MDP  $\widehat{\mathcal{M}}_{\hat{\pi}^{-i}}^i$  is an  $(\varepsilon, \delta)$ -approximation of  $\mathcal{M}_{\hat{\pi}^{-i}}^i$  (see Definition 8).  $\square$

#### 4.5 Proof of Theorem 1

Arbitrarily fix a player  $i \in \mathcal{N}$ . Then, we have the following.

1. From Corollary 8, since  $\hat{\pi}$  is an  $\alpha_{\text{opt}}$ -MPE of  $\widehat{\mathcal{G}}$ , we have that the strategy  $\hat{\pi}^i$  is  $\alpha_{\text{opt}}$ -optimal for MDP  $\widehat{\mathcal{M}}_{\hat{\pi}^{-i}}^i$ .
2. From Lemma 8, we know that MDP  $\widehat{\mathcal{M}}_{\hat{\pi}^{-i}}^i$  is an  $(\varepsilon, \delta)$  approximation of MDP  $\mathcal{M}_{\hat{\pi}^{-i}}^i$ . Then, by Theorem 4, we get that strategy  $\hat{\pi}^i$  is an  $\alpha^i$ -optimal strategy for MDP  $\mathcal{M}_{\hat{\pi}^{-i}}^i$ , where  $\alpha^i$  is given by Theorem 1. Lemma 7 shows that

$$\max_{s \in \mathcal{S}, a^i \in \mathcal{A}^i} \left| \sum_{s' \in \mathcal{S}} \hat{V}_{(\hat{\pi}^i, \hat{\pi}^{-i})}^i(s') P_{\hat{\pi}^{-i}}^i(s'|s, a^i) - \sum_{s' \in \mathcal{S}} \hat{V}_{(\hat{\pi}^i, \hat{\pi}^{-i})}^i(s') \hat{P}_{\hat{\pi}^{-i}}^i(s'|s, a^i) \right|$$

is upper bounded by  $\Delta_{(\hat{\pi}^i, \hat{\pi}^{-i})}^i$  given in Theorem 1. By a similar argument, we can show that

$$\max_{s \in \mathcal{S}, a^i \in \mathcal{A}^i} \left| \sum_{s' \in \mathcal{S}} \hat{V}_{(\hat{\pi}_*^i, \hat{\pi}^{-i})}^i(s') P_{\hat{\pi}^{-i}}^i(s'|s, a^i) - \sum_{s' \in \mathcal{S}} \hat{V}_{(\hat{\pi}_*^i, \hat{\pi}^{-i})}^i(s') \hat{P}_{\hat{\pi}^{-i}}^i(s'|s, a^i) \right|$$

is upper bounded by  $\Delta_{(\hat{\pi}_*^i, \hat{\pi}^{-i})}^i$  given in Theorem 1.

3. Since the above results hold for all  $i \in \mathcal{N}$ , Corollary 8 implies that strategy profile  $\hat{\pi}$  is an  $\alpha$ -MPE of  $\mathcal{G}$ , where  $\alpha := (\alpha^i)_{i \in \mathcal{N}}$  and  $\alpha^i$  is given by Theorem 1.
4. The specific formulas for  $\alpha$  in Corollaries 1 and 2 follow from Corollaries 5 and 6.

#### 4.6 Proofs of Theorem 2 and Theorem 3

The proof argument for Theorems 2 and 3 are similar, so we prove them together. Theorems 5 and 6 show that for any  $p > 0$ ,  $\alpha_{\text{opt}} > 0$  and  $\alpha$  which satisfies appropriate conditions, there exists a function  $N(p)$ , such that for

$n \geq N(p)$  any  $\alpha_{\text{opt}}$ -optimal policy  $\hat{\pi}_n$  of MDP  $\widehat{\mathcal{M}}_n$  is an  $\alpha$ -optimal for MDP  $\mathcal{M}$ .

Now consider the approximate game  $\widehat{\mathcal{G}}_n$  and let  $\pi$  be an  $\alpha_{\text{opt}}$ -MPE of game  $\widehat{\mathcal{G}}$ . For any player  $i \in \mathcal{N}$ , let  $\mathcal{E}_p^i$  denote the event that policy  $\pi^i$  is *not*  $\alpha$ -optimal for  $\mathcal{M}_{\pi^{-i}}^i$ . The above restatement of Theorems 5 and 6 imply that for  $n \geq N(p/|\mathcal{N}|)$ ,  $\mathbb{P}(\mathcal{E}_{p/|\mathcal{N}|}^i) \leq p/|\mathcal{N}|$ . Therefore, by the union bound,

$$\mathbb{P}\left(\bigcup_{i \in \mathcal{N}} \mathcal{E}_{p/|\mathcal{N}|}^i\right) \leq \sum_{i \in \mathcal{N}} \mathbb{P}(\mathcal{E}_{p/|\mathcal{N}|}^i) \leq p.$$

Thus, for  $n \geq N(p/|\mathcal{N}|)$ , the policy  $\pi^i$  is  $\alpha$ -optimal for MDP  $\mathcal{M}_{\pi^{-i}}^i$ , for all  $i \in \mathcal{N}$ . Thus, by Corollary 8,  $\pi$  is an  $\alpha$ -MPE of game  $\mathcal{G}$ .

The results of Theorems 2 and 3 then follow by substituting the expressions for  $N(p/|\mathcal{N}|)$  from Theorems 5 and 6, respectively.

## 5 Numerical examples

In this section, we present two numerical examples to demonstrate the main results of Theorem 1 and 2.

### 5.1 Robustness of Markov perfect equilibrium

Consider a setting where  $\mathcal{N} = \{1, 2\}$ ,  $\mathcal{S} = \{1, 2, 3\}$ ,  $\mathcal{A}_1 = \mathcal{A}_2 = \{1, 2\}$ , and  $\gamma = 0.9$ . We consider two games: original game  $\mathcal{G}$  and approximate game  $\widehat{\mathcal{G}}$  which differ in their reward functions and transition matrices. We describe the transition matrices as  $\{\mathbf{P}(a)\}_{a \in \mathcal{A}}$ , where  $\mathbf{P}(a) = [\mathbf{P}(s' | s, a)]_{s, s' \in \mathcal{S}}$  and describe the reward functions as  $\{r(s)\}_{s \in \mathcal{S}}$  where  $r(s)$  is the bi-matrix  $[(r_1(s, (a_1, a_2)), r_2(s, (a_1, a_2)))]_{(a_1, a_2) \in \mathcal{A}}$ .

For the original game  $\mathcal{G}$ , we have

$$r(1) = \begin{bmatrix} (1.0, 0.4) & (0.7, 1.0) \\ (0.3, 1.0) & (0.8, 0.7) \end{bmatrix}, \quad r(2) = \begin{bmatrix} (0.6, 0.7) & (0.7, 0.6) \\ (0.3, 0.8) & (0.2, 0.2) \end{bmatrix},$$

$$r(3) = \begin{bmatrix} (0.2, 0.6) & (0.1, 0.7) \\ (0.6, 0.7) & (0.5, 0.3) \end{bmatrix},$$

and

$$\mathbf{P}((1, 1)) = \begin{bmatrix} 0.40 & 0.40 & 0.20 \\ 0.10 & 0.50 & 0.40 \\ 0.40 & 0.10 & 0.50 \end{bmatrix}, \quad \mathbf{P}((1, 2)) = \begin{bmatrix} 0.30 & 0.40 & 0.30 \\ 0.20 & 0.20 & 0.60 \\ 0.30 & 0.35 & 0.35 \end{bmatrix},$$

$$\mathbf{P}((2, 1)) = \begin{bmatrix} 0.25 & 0.25 & 0.50 \\ 0.30 & 0.30 & 0.40 \\ 0.20 & 0.20 & 0.60 \end{bmatrix}, \quad \mathbf{P}((2, 2)) = \begin{bmatrix} 0.10 & 0.20 & 0.70 \\ 0.20 & 0.10 & 0.70 \\ 0.40 & 0.20 & 0.40 \end{bmatrix}.$$

For the approximate game  $\widehat{\mathcal{G}}$ , we have

$$\begin{aligned}\hat{r}(1) &= \begin{bmatrix} (0.99, 0.40) & (0.69, 1.00) \\ (0.30, 0.99) & (0.81, 0.71) \end{bmatrix}, & \hat{r}(2) &= \begin{bmatrix} (0.59, 0.70) & (0.69, 0.61) \\ (0.30, 0.80) & (0.19, 0.21) \end{bmatrix}, \\ \hat{r}(3) &= \begin{bmatrix} (0.19, 0.59) & (0.09, 0.70) \\ (0.59, 0.69) & (0.50, 0.30) \end{bmatrix},\end{aligned}$$

and

$$\begin{aligned}\hat{P}((1, 1)) &= \begin{bmatrix} 0.45 & 0.35 & 0.20 \\ 0.15 & 0.45 & 0.40 \\ 0.45 & 0.10 & 0.45 \end{bmatrix}, & \hat{P}((1, 2)) &= \begin{bmatrix} 0.25 & 0.45 & 0.30 \\ 0.25 & 0.15 & 0.60 \\ 0.35 & 0.30 & 0.35 \end{bmatrix}, \\ \hat{P}((2, 1)) &= \begin{bmatrix} 0.25 & 0.30 & 0.45 \\ 0.35 & 0.30 & 0.35 \\ 0.25 & 0.20 & 0.55 \end{bmatrix}, & \hat{P}((2, 2)) &= \begin{bmatrix} 0.15 & 0.15 & 0.70 \\ 0.25 & 0.10 & 0.65 \\ 0.40 & 0.25 & 0.35 \end{bmatrix}.\end{aligned}$$

A MPE of  $\widehat{\mathcal{G}}$  and the corresponding value functions (computed by solving a non-linear program as described in Filar et al. (1991)) are as follows:

$$\hat{\pi}^1 = \begin{bmatrix} 0.33 & 0.67 \\ 1.00 & 0.00 \\ 0.00 & 1.00 \end{bmatrix}, \quad \hat{\pi}^2 = \begin{bmatrix} 0.13 & 0.87 \\ 1.00 & 0.00 \\ 1.00 & 0.00 \end{bmatrix}, \quad (34)$$

$$\hat{V}_{\hat{\pi}}^1 = \begin{bmatrix} 0.6327 \\ 0.6170 \\ 0.6187 \end{bmatrix}, \quad \hat{V}_{\hat{\pi}}^2 = \begin{bmatrix} 0.7258 \\ 0.7148 \\ 0.7148 \end{bmatrix}. \quad (35)$$

In (34), the strategy is described as  $\hat{\pi}^i = [\hat{\pi}^i(a^i|s)]_{s \in \mathcal{S}, a^i \in \mathcal{A}^i}$ . For strategy  $\hat{\pi}$  in (34), we compute the value functions  $V_{\hat{\pi}}^i$  for game  $\mathcal{G}$  as described in Proposition 5 and the value functions  $V_{(*, \hat{\pi}^{-i})}^i$  as described in Proposition 6 (see Sec. 4). These are given by

$$V_{\hat{\pi}}^1 = \begin{bmatrix} 0.6341 \\ 0.6192 \\ 0.6209 \end{bmatrix}, \quad V_{\hat{\pi}}^2 = \begin{bmatrix} 0.7252 \\ 0.7142 \\ 0.7154 \end{bmatrix}, \quad (36)$$

$$V_{(*, \hat{\pi}^2)}^1 = \begin{bmatrix} 0.6394 \\ 0.6222 \\ 0.6241 \end{bmatrix}, \quad V_{(\hat{\pi}^1, *)}^2 = \begin{bmatrix} 0.7280 \\ 0.7158 \\ 0.7171 \end{bmatrix}. \quad (37)$$

Note that

$$\alpha_*^1 = \|V_{(*, \hat{\pi}^2)}^1 - V_{\hat{\pi}}^1\|_{\infty} = 0.005300, \quad (38a)$$

$$\alpha_*^2 = \|V_{(\hat{\pi}^1, *)}^2 - V_{\hat{\pi}}^2\|_{\infty} = 0.002785. \quad (38b)$$

Thus,  $\hat{\pi}$  is a (0.005300, 0.002785)-MPE of  $\mathcal{G}$ .

Now, we compare  $\alpha_*$  with the bounds that we obtain using Theorem 1. Note that since  $\alpha_{\text{opt}} = 0$ ,  $\hat{\pi}^i$  defined in Theorem 1 is equal to  $\hat{\pi}^i$ . Therefore, the first upper bound on  $\alpha$  is given by  $2\varepsilon + 2\gamma\Delta_{\hat{\pi}}^i/(1-\gamma)$ . Note that

$$\max_{a \in \mathcal{A}} \max_{s \in \mathcal{S}} |r(s, a) - \hat{r}(s, a)| = 0.01.$$

Thus,  $\varepsilon = 0.01$ . Moreover,

$$\begin{aligned} \Delta_{\hat{\pi}}^1 &= \max_{s \in \mathcal{S}, a \in \mathcal{A}} \left| \sum_{s' \in \mathcal{S}} \left[ P(s'|s, a) \hat{V}_{\hat{\pi}}^1(s') - \hat{P}(s'|s, a) \hat{V}_{\hat{\pi}}^1(s') \right] \right| = 0.000784 \\ \Delta_{\hat{\pi}}^2 &= \max_{s \in \mathcal{S}, a \in \mathcal{A}} \left| \sum_{s' \in \mathcal{S}} \left[ P(s'|s, a) \hat{V}_{\hat{\pi}}^2(s') - \hat{P}(s'|s, a) \hat{V}_{\hat{\pi}}^2(s') \right] \right| = 0.000550 \end{aligned}$$

Then, by Theorem 1, we have that

$$\alpha \leq 2 \left( \varepsilon + \frac{\gamma \Delta_{\hat{\pi}}}{1 - \gamma} \right) = 2 \times 0.01 + \frac{2 \times 0.9}{0.1} \begin{bmatrix} 0.000784 \\ 0.000550 \end{bmatrix} = \begin{bmatrix} 0.034112 \\ 0.029900 \end{bmatrix}.$$

Now, we consider the upper bound on  $\Delta_{\hat{\pi}}^i$  in terms of  $\rho_F(\hat{V}_{\hat{\pi}}^i)$ .

1. We first consider the case when  $\mathfrak{F} = \mathfrak{F}^{\text{TV}}$ . Note that

$$\max_{a \in \mathcal{A}} \max_{s \in \mathcal{S}} d_{\mathfrak{F}^{\text{TV}}}(P(\cdot|s, a), \hat{P}(\cdot|s, a)) = 0.05,$$

Thus when  $\mathfrak{F} = \mathfrak{F}^{\text{TV}}$ ,  $\hat{\mathcal{G}}$  is a  $(0.01, 0.05)$ -approximation of game  $\mathcal{G}$ . Also note that  $\text{span}(\hat{V}_{\hat{\pi}}^1) = 0.015684$  and  $\text{span}(\hat{V}_{\hat{\pi}}^2) = 0.010990$ . Then, from Theorem 1, we have that

$$\alpha \leq 2 \left( \varepsilon + \frac{\gamma \delta \text{span}(\hat{V}_{\hat{\pi}}^i)}{1 - \gamma} \right) = 2 \times 0.01 + \frac{2 \times 0.9 \times 0.05}{0.1} \begin{bmatrix} 0.015684 \\ 0.010990 \end{bmatrix} = \begin{bmatrix} 0.034116 \\ 0.029903 \end{bmatrix}.$$

2. Now we equip the state space  $\mathcal{S}$  with a metric  $d$  where  $d(s, s') = |s - s'|$  and consider the case  $\mathfrak{F} = \mathfrak{F}^{\text{W}}$ . Note that

$$\max_{a \in \mathcal{A}} \max_{s \in \mathcal{S}} d_{\mathfrak{F}^{\text{W}}}(P(\cdot|s, a), \hat{P}(\cdot|s, a)) = 0.10.$$

Thus when  $\mathfrak{F} = \mathfrak{F}^{\text{W}}$ ,  $\hat{\mathcal{G}}$  is a  $(0.01, 0.10)$ -approximation of game  $\mathcal{G}$ . Also note that  $\text{Lip}(\hat{V}_{\hat{\pi}}^1) = 0.015684$  and  $\text{Lip}(\hat{V}_{\hat{\pi}}^2) = 0.010990$ . Then, from Theorem 1, we have that

$$\alpha \leq 2 \left( \varepsilon + \frac{\gamma \delta \text{Lip}(\hat{V}_{\hat{\pi}}^i)}{1 - \gamma} \right) = 2 \times 0.01 + \frac{2 \times 0.9 \times 0.10}{0.1} \begin{bmatrix} 0.015684 \\ 0.010990 \end{bmatrix} = \begin{bmatrix} 0.048231 \\ 0.039782 \end{bmatrix}.$$

The above example shows that with the given  $(\varepsilon, \delta)$ , the bound of Theorem 1 is loose by only a small multiplicative factor of approximately 6 to 15.

## 5.2 Sample complexity of generative models

We now consider the setting for model based MARL. Consider the game  $\mathcal{G}$  described in Sec. 5.1 but suppose that the transition matrix  $\mathbf{P}$  is not known but we have access to a generative model which can generate samples from  $\mathbf{P}$ . We assume that we are interested in identifying an  $\alpha$ -MPE, where  $\alpha = 0.1$  with probability  $1 - p = 0.99$ . We assume that we can exactly compute the MPE of the approximated game, i.e.,  $\alpha_{\text{opt}} = 0$ .

From Theorem 2, we have an upper bound on the number  $n$  of samples for each state-action pair as

$$\begin{aligned} n &\geq \left\lceil \left( \frac{\gamma}{1-\gamma} \text{span}(r) \right)^2 \frac{2 \log(4|\mathcal{S}|(\prod_{i \in \mathcal{N}} |\mathcal{A}^i|) |\mathcal{N}| p^{-1})}{\alpha^2} \right\rceil \\ &= \left\lceil \left( \frac{0.9}{1-0.9} \times 0.9 \right)^2 \frac{2 \log(4 \times 3 \times 4 \times 2/0.01)}{0.1^2} \right\rceil = 120,323. \end{aligned}$$

We now verify this result via simulation. We run  $M = 1,000$  experiments. For each experiment, we generate  $n = 120,323$  samples for each state-action pair, and estimate an empirical model  $\hat{\mathbf{P}}_n(s'|s, a) = \text{count}(s'|s, a)/n$ . We compute the MPE  $\hat{\pi}_n$  for the approximate game  $\hat{\mathcal{G}} = \langle \mathcal{S}, \{\mathcal{A}\}_{i \in \mathcal{N}}, \hat{\mathbf{P}}, r, \gamma \rangle$ . Then, using Proposition 6, we compute  $\alpha_n = (\alpha_n^1, \alpha_n^2)$  such that  $\hat{\pi}_n$  is a  $\alpha$ -MPE of game  $\mathcal{G}$ . The scatter plot of  $\alpha_n = (\alpha_n^1, \alpha_n^2)$  along with the empirical distribution of  $\alpha_n^1$  and  $\alpha_n^2$  are shown in Fig. 1. The values for  $\alpha_n^1$  are typically larger than the values for  $\alpha_n^2$  because of the parameters of the specific game  $\mathcal{G}$  chosen in the example. This trend is also apparent in the analytical  $\alpha$  calculated previously based on worst case errors. Note that for most cases, both  $\alpha_n^1$  and  $\alpha_n^2$  are smaller than  $10^{-4}$ , which is much less than our target  $\alpha$  of 0.1. This highlights the looseness of the upper bound in Theorem 2.

Note that  $\gamma = 0.9 \geq 0.8889$ , which is the lower bound on the critical discount factor  $\gamma^\circ$  calculated in Lemma 2. So, we don't know upfront whether the Hoeffding-type bound is tighter than the Bernstein-type bound. We compute the number of samples  $n$  for each state-action pair based on Theorem 3, which is given by

$$\begin{aligned} n &\geq \left\lceil \frac{c\gamma \|r\|_\infty^2 \log(8|\mathcal{S}|(\prod_{i \in \mathcal{N}} |\mathcal{A}^i|) |\mathcal{N}| (1-\gamma)^{-2} p^{-1})}{(1-\gamma)\alpha^2} \right\rceil \\ &\geq \left\lceil \frac{64 \times 0.9 \times 1.0^2 \log(8 \times 3 \times 4 \times 2/((1-0.9)^2 \times 0.01))}{(1-0.9) \times 0.1^2} \right\rceil = 833,348. \end{aligned}$$

which is larger than the sample complexity bound of  $n \geq 120,323$  calculated using Theorem 2 above.

## 6 Conclusion

In this paper, we quantify how robust MPE are to model approximations given the degree of approximation in the approximate game. We provide bounds on



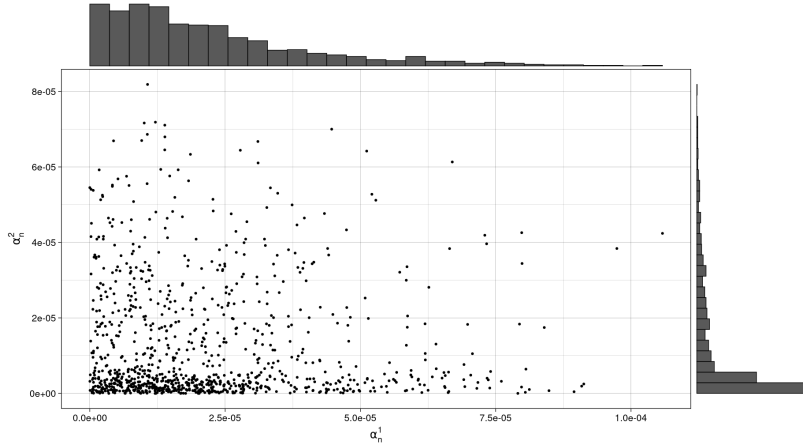


Fig. 1: Scatter plot of  $(\alpha_n^1, \alpha_n^2)$  such that the policy  $\hat{\pi}_n$  of game  $\mathcal{G}_n$  is a  $(\alpha_n^1, \alpha_n^2)$ -MPE of game  $\mathcal{G}$  for  $M = 1,000$  independently generated experiments for  $n = 120, 323$ . The histogram of the marginal probability distribution for  $\alpha_n^1$  and  $\alpha_n^2$  are shown on the top and the right.

the degree of approximation based on the approximation error in the reward and transition functions and properties of the value function of the MPE. We also present coarser, instance independent upper bounds, which do not depend of the value function but only depend on the properties of the reward and transition function of the approximate game. Using these approximation bounds, we provide sample complexity bounds for computing an approximate MPE using a generative model.

An interesting feature of the results is that the approximation bounds depend on the choice of the metric on probability spaces. We work with a class of metrics known as IPMs and specialize our results for two specific choices of IPMs: total variation distance and Wasserstein distance. However, the results are applicable to any IPM. For games with high-dimensional state spaces, metrics such as maximum mean discrepancy (Sriperumbudur et al. 2008) might be more appropriate.

The generative model setting considered in this paper circumvents the exploration problem of learning. Therefore, the sample complexity results presented in this paper should be viewed as a lower bound on the number of samples required by an algorithm to learn an  $\alpha$ -MPE. The proposed algorithm is not practical as it requires storing the transition matrix, which has a size of  $|\mathcal{S}| \prod_{i \in \mathcal{N}} |\mathcal{A}^i|$ , which is exponential in the number of agents. The algorithm also assumes that there is a system planner which computes the approximate MPE. Developing a scalable and distributing algorithm for learning MPE remains a challenging open problem.

We conclude by noting that the results presented in this paper were restricted to Markov games with perfect information. An interesting future di-

rection is to develop similar approximation bounds for Markov games with imperfect information as well as specific classes of dynamic games such as mean-field games and their variants.

### Compliance with Ethical Standards: Conflict of Interest Statement

The authors declare that they have no conflict of interest.

### Data Availability Statement

Data sharing not applicable to this article as no datasets were generated or analysed during the current study.

### A Proof of Theorem 4

Let  $V_*$  denote the optimal value function for MDP  $\mathcal{M}$  and  $V_\pi$  denote the value function for policy  $\pi$  in MDP  $\mathcal{M}$ . Let  $\hat{\pi}_*$  be the optimal policy for  $\widehat{\mathcal{M}}$  and  $\hat{\pi}$  be an  $\alpha_{\text{opt}}$ -optimal policy for  $\widehat{\mathcal{M}}$ . From triangle inequality, we have

$$\|V_* - V_{\hat{\pi}}\|_\infty \leq \|V_* - \hat{V}_{\hat{\pi}_*}\|_\infty + \|\hat{V}_{\hat{\pi}_*} - \hat{V}_{\hat{\pi}}\|_\infty + \|V_{\hat{\pi}} - \hat{V}_{\hat{\pi}}\|_\infty. \quad (39)$$

Now we bound the three terms separately. For the first term, we have

$$\begin{aligned} \|V_* - \hat{V}_{\hat{\pi}_*}\|_\infty &\stackrel{(a)}{\leq} \max_{s \in \mathcal{S}} \left| \max_{a \in \mathcal{A}} \left[ (1 - \gamma)r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathbf{P}(s'|s, a)V_*(s') \right. \right. \\ &\quad \left. \left. - (1 - \gamma)\hat{r}(s, a) - \gamma \sum_{s' \in \mathcal{S}} \hat{\mathbf{P}}(s'|s, a)\hat{V}_{\hat{\pi}_*}(s') \right] \right| \\ &\stackrel{(b)}{\leq} (1 - \gamma) \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} |r(s, a) - \hat{r}(s, a)| \\ &\quad + \gamma \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} \left| \sum_{s' \in \mathcal{S}} \mathbf{P}(s'|s, a)V_*(s') - \mathbf{P}(s'|s, a)\hat{V}_{\hat{\pi}_*}(s') \right| \\ &\quad + \gamma \max_{(s, a) \in \mathcal{S} \times \mathcal{A}} \left| \sum_{s' \in \mathcal{S}} \mathbf{P}(s'|s, a)\hat{V}_{\hat{\pi}_*}(s') - \hat{\mathbf{P}}(s'|s, a)\hat{V}_{\hat{\pi}_*}(s') \right| \\ &\leq (1 - \gamma)\varepsilon + \gamma\|V_* - \hat{V}_{\hat{\pi}_*}\|_\infty + \gamma\Delta_{\hat{\pi}_*}, \end{aligned}$$

where (a) relies on the fact that  $\max f(x) \leq \max |f(x) - g(x)| + \max g(x)$ , (b) follows from triangle inequality, and (c) follows from the definition of  $\varepsilon$  and  $\Delta_{\hat{\pi}_*}$ . Therefore,

$$\|V_* - \hat{V}_{\hat{\pi}_*}\|_\infty \leq \varepsilon + \frac{\gamma\Delta_{\hat{\pi}_*}}{1 - \gamma}. \quad (40)$$

For the second term of (39), we have

$$\|\hat{V}_{\hat{\pi}_*} - \hat{V}_{\hat{\pi}}\|_\infty \leq \alpha_{\text{opt}}, \quad (41)$$

since  $\hat{\pi}$  is an  $\alpha_{\text{opt}}$ -optimal policy of  $\widehat{\mathcal{M}}$ .

For the third term of (39), we have

$$\begin{aligned}
\|V_{\hat{\pi}} - \hat{V}_{\hat{\pi}}\|_{\infty} &= \max_{s \in \mathcal{S}} \left| \sum_{a \in \mathcal{A}} \hat{\pi}(a|s) \left[ (1-\gamma)r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathbf{P}(s'|s, a) V_{\hat{\pi}}(s') \right. \right. \\
&\quad \left. \left. - (1-\gamma)\hat{r}(s, a) - \gamma \sum_{s' \in \mathcal{S}} \hat{\mathbf{P}}(s'|s, a) \hat{V}_{\hat{\pi}}(s') \right] \right| \\
&\stackrel{(d)}{\leq} (1-\gamma) \max_{s \in \mathcal{S}} \left| \sum_{a \in \mathcal{A}} \hat{\pi}(a|s) [r(s, a) - \hat{r}(s, a)] \right| \\
&\quad + \gamma \max_{s \in \mathcal{S}} \left| \sum_{a \in \mathcal{A}} \hat{\pi}(a|s) \left[ \sum_{s' \in \mathcal{S}} [\mathbf{P}(s'|s, a) V_{\hat{\pi}}(s') - \mathbf{P}(s'|s, a) \hat{V}_{\hat{\pi}}(s')] \right] \right| \\
&\quad + \gamma \max_{s \in \mathcal{S}} \left| \sum_{a \in \mathcal{A}} \hat{\pi}(a|s) \left[ \sum_{s' \in \mathcal{S}} [\mathbf{P}(s'|s, a) \hat{V}_{\hat{\pi}}(s') - \hat{\mathbf{P}}(s'|s, a) \hat{V}_{\hat{\pi}}(s')] \right] \right| \\
&\stackrel{(e)}{\leq} (1-\gamma)\varepsilon + \gamma\|V_{\hat{\pi}} - \hat{V}_{\hat{\pi}}\|_{\infty} + \Delta_{\hat{\pi}}
\end{aligned}$$

where (d) follows from triangle inequality and (e) follows from the definition of  $\varepsilon$  and  $\Delta_{\hat{\pi}}$ . Therefore,

$$\|V_{\hat{\pi}} - \hat{V}_{\hat{\pi}}\|_{\infty} \leq \varepsilon + \frac{\gamma\Delta_{\hat{\pi}}}{1-\gamma}. \quad (42)$$

The result then follows by substituting (40)–(41) in (39).  $\square$

## References

- D. Acemoglu and J. A. Robinson. A theory of political transitions. *American Economic Review*, 91(4):938–963, 2001.
- A. Agarwal, S. Kakade, and L. F. Yang. Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pages 67–83. PMLR, 2020.
- V. Aguirregabiria and P. Mira. Sequential estimation of dynamic discrete games. *Econometrica*, 75(1):1–53, 2007.
- N. Akchurina. *Multi-agent reinforcement learning algorithms*. PhD thesis, University of Paderborn, 2010.
- S. C. Albright and W. Winston. A birth–death model of advertising and pricing. *Advances in Applied Probability*, 11(1):134–152, 1979.
- E. Altman. *Constrained Markov decision processes: stochastic modeling*. CRC Press, 1999.
- M. G. Azar, R. Munos, and H. J. Kappen. Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, 91(3):325–349, 2013.
- P. Bajari, C. L. Benkard, and J. Levin. Estimating dynamic models of imperfect competition. *Econometrica*, 75(5):1331–1370, 2007.
- T. Başar and P. Bernhard. *H-infinity optimal control and related minimax design problems: a dynamic game approach*. Springer Science & Business Media, 2008.
- T. Başar and G. Zaccour, editors. *Handbook of Dynamic Game Theory*. Springer International Publishing, 2018.
- D. P. Bertsekas. *Dynamic programming and optimal control*. Belmont, MA: Athena Scientific, 2017.
- M. Breton. *Algorithms for Stochastic Games*, pages 45–57. Springer, 1991.
- L. Busoniu, R. Babuska, and B. De Schutter. A comprehensive survey of multiagent reinforcement learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 38(2):156–172, 2008.
- N. Cesa-Bianchi and G. Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.

- X. Deng, Y. Li, D. H. Mguni, J. Wang, and Y. Yang. On the complexity of computing markov perfect equilibrium in general-sum stochastic games. *arXiv preprint arXiv:2109.01795*, 2021.
- U. Doraszelski and J. F. Escobar. A theory of regular Markov perfect equilibria in dynamic stochastic games: Genericity, stability, and purification. *Theoretical Economics*, 5(3): 369–402, 2010. ISSN 1555-7561.
- R. Ericson and A. Pakes. Markov-perfect industry dynamics: A framework for empirical work. *The Review of economic studies*, 62(1):53–82, 1995.
- C. Fershtiman and A. Pakes. A dynamic oligopoly with collusion and price wars. *The RAND Journal of Economics*, 31(2):207–236, 2000.
- J. Filar and K. Vrieze. *Competitive Markov Decision Processes*. Springer, New York, NY, 1996. ISBN 978-1-4612-8481-9 978-1-4612-4054-9.
- J. A. Filar, T. A. Schultz, F. Thuijsman, and O. Vrieze. Nonlinear programming and stationary equilibria in stochastic games. *Mathematical Programming*, 50(1):227–237, 1991.
- A. M. Fink. Equilibrium in a stochastic  $n$ -person game. *Hiroshima Mathematical Journal*, 28(1), 1964.
- P. J.-J. Herings and R. Peeters. Homotopy methods to compute equilibria in game theory. *Economic Theory*, 42(1):119–156, 2010.
- P. J.-J. Herings, R. J. Peeters, et al. Stationary equilibria in stochastic games: Structure, selection, and computation. *Journal of Economic Theory*, 118(1):32–60, 2004.
- K. Hinderer. Lipschitz continuity of value functions in Markovian decision processes. *Mathematical Methods of Operations Research*, 62(1):3–22, Sep 2005. ISSN 1432-5217.
- A. J. Hoffman and R. M. Karp. On nonterminating stochastic games. *Management Science*, 12(5):359–370, 1966.
- A. Jaśkiewicz and A. S. Nowak. Robust Markov perfect equilibria. *Journal of Mathematical Analysis and Applications*, 419(2):1322–1332, 2014.
- S. M. Kakade. *On the sample complexity of reinforcement learning*. PhD thesis, University College, London, 2003.
- M. Kearns and S. Singh. Finite-sample convergence rates for q-learning and indirect algorithms. In *Advances in neural information processing systems*, pages 996–1002, 1999.
- O. Krupnik, I. Mordatch, and A. Tamar. Multi-agent reinforcement learning with multi-step generative models. *arXiv preprint arXiv:1901.10251*, Nov. 2019.
- S. Leonardos, W. Overman, I. Panageas, and G. Piliouras. Global convergence of multi-agent policy gradient in markov potential games. *arXiv preprint arXiv:2106.01969*, 2021.
- G. Li, Y. Wei, Y. Chi, Y. Gu, and Y. Chen. Breaking the sample size barrier in model-based reinforcement learning with a generative model. *Advances in neural information processing systems*, 33, 2020.
- M. L. Littman. Markov games as a framework for multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 157–163. Elsevier, 1994.
- M. L. Littman. Value-function reinforcement learning in Markov games. *Cognitive systems research*, 2(1):55–66, 2001.
- G. J. Mailath and L. Samuelson. *Repeated Games and Reputations: Long-Run Relationships*. Oxford University Press, 2006.
- E. Maskin and J. Tirole. A theory of dynamic oligopoly, I: Overview and quantity competition with large fixed costs. *Econometrica: Journal of the Econometric Society*, pages 549–569, 1988a.
- E. Maskin and J. Tirole. A theory of dynamic oligopoly, II: Price competition, kinked demand curves, and edgeworth cycles. *Econometrica: Journal of the Econometric Society*, pages 571–599, 1988b.
- E. Maskin and J. Tirole. Markov perfect equilibrium: I. observable actions. *Journal of Economic Theory*, 100(2):191–219, 2001.
- A. Müller. How does the value function of a Markov decision process depend on the transition probabilities? *Mathematics of Operations Research*, 22(4):872–885, 1997.
- A. Müller. Integral probability metrics and their generating classes of functions. *Advances in Applied Probability*, 29(2):429–443, 1997.
- A. Pakes, M. Ostrovsky, and S. Berry. Simple estimators for the parameters of discrete dynamic games (with entry/exit examples). *the RAND Journal of Economics*, 38(2):

- 373–399, 2007.
- J. Pérolat, F. Strub, B. Piot, and O. Pietquin. Learning Nash equilibrium for general-sum Markov games from batch data. In *Artificial Intelligence and Statistics*, pages 232–241. PMLR, 2017.
- M. Pesendorfer and P. Schmidt-Dengler. Asymptotic least squares estimators for dynamic games. *The Review of Economic Studies*, 75(3):901–928, 2008.
- H. Prasad, P. LA, and S. Bhatnagar. Two-timescale algorithms for learning Nash equilibria in general-sum stochastic games. In *Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems*, pages 1371–1379, 2015.
- P. D. Rogers. *Nonzero-sum stochastic games*. PhD thesis, University of California, Berkeley, 1969.
- S. Sengupta, A. Chowdhary, D. Huang, and S. Kambhampati. General sum markov games for strategic detection of advanced persistent threats using moving target defense in cloud networks. In *International Conference on Decision and Game Theory for Security*, pages 492–512. Springer, 2019.
- L. S. Shapley. Stochastic games. *Proceedings of the National Academy of Sciences*, 39(10):1095–1100, 1953.
- Y. Shoham, R. Powers, and T. Grenager. Multi-agent reinforcement learning: a critical survey. Technical report, Stanford University, 2003.
- A. Sidford, M. Wang, X. Wu, L. F. Yang, and Y. Ye. Near-optimal time and sample complexities for solving markov decision processes with a generative model. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 5192–5202, 2018.
- A. Sidford, M. Wang, L. Yang, and Y. Ye. Solving discounted stochastic two-player games with near-optimal time and sample complexity. In *International Conference on Artificial Intelligence and Statistics*, pages 2992–3002. PMLR, 2020.
- E. Solan. *A Course in Stochastic Game Theory*. Cambridge University Press, 2021.
- Z. Song, S. Mei, and Y. Bai. When can we learn general-sum markov games with a large number of players sample-efficiently? *arXiv preprint arXiv:2110.04184*, 2021.
- B. K. Sriperumbudur, A. Gretton, K. Fukumizu, G. R. G. Lanckriet, and B. Schölkopf. Injective Hilbert space embeddings of probability measures. In *Conference on Learning Theory*, 2008.
- J. Subramanian, A. Sinha, R. Seraj, and A. Mahajan. Approximate information state for approximate planning and reinforcement learning in partially observed systems. *J. Mach. Learn. Res.*, 23:12–1, 2022.
- R. S. Sutton. Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In *International Conference on Machine Learning*, pages 216–224. San Francisco (CA), 1990.
- M. Takahashi. Equilibrium points of stochastic non-cooperative  $n$ -person games. *Hiroshima Mathematical Journal*, 28(1), Jan. 1964.
- M. M. Tidball and E. Altman. Approximations in dynamic zero-sum games I. *SIAM Journal on Control and Optimization*, 34(1):311–328, Jan. 1996.
- M. M. Tidball, O. Pourtallier, and E. Altman. Approximations in dynamic zero-sum games II. *SIAM Journal on Control and Optimization*, 35(6):2101–2117, Nov. 1997.
- O. J. Vrieze. *Stochastic games with finite state and action spaces*. CWI, Jan. 1987. ISBN 978-90-6196-313-4.
- T. Wang, X. Bao, I. Clavera, J. Hoang, Y. Wen, E. Langlois, S. Zhang, G. Zhang, P. Abbeel, and J. Ba. Benchmarking Model-Based Reinforcement Learning. *arXiv preprint arXiv:1907.02057*, July 2019.
- W. Whitt. Representation and approximation of noncooperative sequential games. *SIAM Journal on Control and Optimization*, 18(1):33–48, 1980.
- K. Zhang, S. Kakade, T. Basar, and L. Yang. Model-based multi-agent rl in zero-sum Markov games with near-optimal sample complexity. *Advances in Neural Information Processing Systems*, 33, 2020.
- K. Zhang, Z. Yang, and T. Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Handbook of Reinforcement Learning and Control*, pages 321–384, 2021a.

- 
- R. Zhang, Z. Ren, and N. Li. Gradient play in multi-agent markov stochastic games: Stationary points and convergence. *arXiv e-prints*, pages arXiv–2106, 2021b.
- W. Zhang, X. Wang, J. Shen, and M. Zhou. Model-based Multi-agent Policy Optimization with Adaptive Opponent-wise Rollouts. In *International Joint Conference on Artificial Intelligence*. Montreal, Canada, 2021c.
- M. Zinkevich, A. Greenwald, and M. Littman. Cyclic equilibria in Markov games. In *Neural Information Processing Systems*, pages 1641–1648, 2006.