

---

# Minimax Sample Complexity for Turn-based Stochastic Game (Supplementary material)

---

Qiwen Cui<sup>1</sup>

Lin F. Yang<sup>2</sup>

<sup>1</sup>School of Mathematical Sciences, Peking University

<sup>2</sup>Electrical and Computer Engineering Department, University of California, Los Angeles

## Abstract

The empirical success of multi-agent reinforcement learning is encouraging, while few theoretical guarantees have been revealed. In this work, we prove that the plug-in solver approach, probably the most natural reinforcement learning algorithm, achieves minimax sample complexity for turn-based stochastic game (TBSG). Specifically, we perform planning in an empirical TBSG by utilizing a ‘simulator’ that allows sampling from arbitrary state-action pair. We show that the empirical Nash equilibrium strategy is an approximate Nash equilibrium strategy in the true TBSG and give both problem-dependent and problem-independent bound. We develop reward perturbation techniques to tackle the non-stationarity in the game and Taylor-expansion-type analysis to improve the dependence on approximation error. With these novel techniques, we prove the minimax sample complexity of turn-based stochastic game.

## 1 INTRODUCTION

Reinforcement learning (RL) [Sutton and Barto, 2018], a framework where agents learn to make sequential decisions in an unknown environment, has received tremendous attention. An interesting branch is multi-agent reinforcement learning (MARL) that multiple agents exist and they interact with the environment as well as the others, which bridges RL and game theory. In general, each agent attempts to maximize its own reward by utilizing the data collected from the environment and also inferring other agents’ strategies. Impressive successes have been achieved in games such as backgammon [Tesauro, 1995], Go [Silver et al., 2017] and strategy games [Ye et al., 2020]. MARL has shown the potential for superhuman performance, but theoretical

guarantees are rather rare due to complex interaction between agents that makes the problem considerably harder than single agent reinforcement learning. This is also known as non-stationarity in MARL, which means when multiple agents alter their strategies based on samples collected from previous strategy, the system becomes non-stationary for each agent and the improvement can not be guaranteed. One fundamental question in MBRL is that how to design efficient algorithms to overcome non-stationarity.

Two-players turn-based stochastic game (TBSG) is a two-agents generalization of Markov decision process (MDP), where two agents choose actions in turn and one agent wants to maximize the total reward while the other wants to minimize it. As a zero-sum game, TBSG is known to have Nash equilibrium strategy [Shapley, 1953], which means there exists a strategy pair that both agents will not benefit from changing its strategy alone, and our target is to find the (approximate) Nash equilibrium strategy. Dynamic programming type algorithms is a basic but powerful approach to solve TBSG. Strategy iteration, a counterpart of policy iteration in MDP, is known to be a polynomial complexity algorithm to solve TBSG with known transition kernel [Hansen et al., 2013, Jia et al., 2020]. However, these algorithms suffer from high computational cost and require full knowledge of the transition dynamic. Reinforcement learning is a promising alternative, which has demonstrated its potential in solving sequential decision making problems. However, non-stationarity pose a large obstacle against the convergence of model-free algorithms. To tackle this challenge, sophisticated algorithms have been proposed [Jia et al., 2019, Sidford et al., 2020]. In this work, we focus on another promising but insufficiently developed method, model-based algorithms, where agents learn the model of the environment and plan in the empirical model.

Model-based RL is long perceived to be the cure to the sample inefficiency in RL, which is also justified by empirical advances [Kaiser et al., 2019, Wang et al., 2019]. However, the theoretical understanding of model-based RL is still far from complete. Recently, a line of research focuses on

analyzing model-based algorithm under generative model setting [Azar et al., 2013, Agarwal et al., 2019, Cui and Yang, 2020, Zhang et al., 2020, Li et al., 2020]. In this work, we aim to prove that the simplest model-based algorithm, plug-in solver approach, enjoys minimax sample complexity for TBSG by utilizing novel absorbing TBSG and reward perturbation techniques.

Specifically, we assume that we have access to a generative model oracle [Kakade et al., 2003], which allows sampling from arbitrary state-action pair. It is known that exploration and exploitation tradeoff is simplified under this setting, i.e., sampling from each state-action pair equally can already yield the minimax sample complexity result. With the generative model, the most intuitive approach is to learn an empirical TBSG and then use a planning algorithm to find the empirical Nash equilibrium strategy as the solution, which separates learning and planning. This kind of algorithm is known as the ‘plug-in solver approach’, which means *arbitrary* planning algorithm can be used. We will show that this simple plug-in solver approach enjoys minimax sample complexity. For MDP, this line of research is studied in the seminal work by Azar et al. [2013]. In the recent work [Li et al., 2020], they almost completely solves this problem by proving the minimax complexity with full range of error  $\epsilon$ . However, the result for TBSG is still unknown and largely due to interaction between agents. In this work, we give the sample complexity upper bounds of TBSG for both problem-dependent and problem-independent cases.

Suboptimality gap is a widely studied notion in the bandit theory, which stands for a constant gap between the optimal arm and second optimal arm. This notion has received increasing focus in MDP and we generalize it to TBSG. To start with, we prove that  $\tilde{O}(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-3}\Delta^{-2})$  samples are enough to recover Nash equilibrium strategy accurately for  $\Delta \in (0, (1-\gamma)^{-\frac{1}{2}}]$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $\gamma$  is the discount factor and  $\Delta$  is the suboptimality gap. The key analysis tool is the absorbing TBSG technique that helps to show the empirical optimal Q value is close to true optimal Q value. With the absorbing TBSG, the statistical dependence on  $\hat{P}(s, a)$  is moved to a single parameter  $u$ , which can be approximated by an covering argument. Therefore, standard concentration arguments combined with union bound can be applied. Suboptimality gap plays an critical role in showing the empirical Nash equilibrium strategy is exactly the same as true Nash equilibrium strategy.

The main contribution is in the second part, where we give our problem-independent bound which meets the existing lower bound  $O(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-3}\epsilon^{-2})$  [Azar et al., 2013, Sidford et al., 2020]. Note that the problem-dependent result becomes meaningless when the suboptimality gap is sufficiently small and also the gap may not even exist. In the problem-independence case, we develop the reward perturbation technique to create such a gap in the estimated

TBSG, which is inspired by the technique developed in [Li et al., 2020]. The key observation is that if  $r(s, a)$  increases,  $Q^*(s, a)$  increases much faster than  $Q^*(s, a')$  for  $a' \neq a$ , thus the perturbed TBSG enjoys a suboptimality gap with high probability. Different from the usage in the first part, here the suboptimality gap is used to ensure the empirical Nash equilibrium strategy lies in a finite set so that union bound can be applied. Combining the reward perturbation technique with the absorbing TBSG technique, we are able to prove more subtle concentration arguments and finally show that the empirical Nash equilibrium strategy is an  $\epsilon$ -approximate Nash equilibrium strategy in the true TBSG with minimax sample complexity  $\tilde{O}(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-3}\epsilon^{-2})$  for  $\epsilon \in (0, (1-\gamma)^{-1}]$ . In addition, we can recover the problem-dependent bound by simply setting  $\epsilon = \Delta$ .

Recently Zhang et al. [2020] proposes a similar result for simultaneous stochastic game. However, their result requires a planning oracle for regularized stochastic game, which is computationally intractable. Our algorithm only requires planning in a standard turn-based stochastic game, which can be performed efficiently by utilizing strategy iteration [Hansen et al., 2013] or learning algorithms [Sidford et al., 2020]. In addition, their sample complexity result  $N = \tilde{O}(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-3}\epsilon^{-2})$  only holds for  $\epsilon \in (0, (1-\gamma)^{-\frac{1}{2}}]$ , which means  $N = \tilde{O}(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-2})$ . our result fills the blank in the sample region  $\tilde{O}(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-1}) \leq N \leq \tilde{O}(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-2})$ . Note that the lower bound [Azar et al., 2013] indicates that  $N = O(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-1})$  is insufficient to learn a policy. As both parts of our analysis heavily rely on the suboptimality gap, we hope our work can provide more understanding about this notion and TBSG.

## 2 PRELIMINARY

**Turn-based Stochastic Game** Turn-based two-player zero-sum stochastic game (TBSG) is a generalized version of Markov decision process (MDP) which includes two players competing with each other. Player 1 aims to maximize the total reward while player 2 aims to minimize it. TBSG is described by the tuple  $\mathcal{G} = (\mathcal{S} = \mathcal{S}_{\max} \cup \mathcal{S}_{\min}, \mathcal{A}, P, r, \gamma)$ , where  $\mathcal{S}_{\max}$  is the state space of player 1,  $\mathcal{S}_{\min}$  is the state space of player 2,  $\mathcal{A}$  is the action space of both players,  $P \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$  is the transition probability matrix,  $r \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$  is the reward vector and  $\gamma$  is the discount factor. In each step, only one player plays an action and a transition happens. For instance, if the state  $s \in \mathcal{S}_{\max}$ , player 1 needs to select an action  $a \in \mathcal{A}$ . After selecting the action, the state will transit to  $s' \in \mathcal{S}_{\min}$  according to the distribution  $P(\cdot|s, a)$  with reward  $r(s, a)$  and player 2 needs to choose the action. For representation simplicity and without loss of generality, we assume that  $r$  is known and only  $P$  is

unknown<sup>1</sup>.

We denote a strategy pair as  $\pi := (\mu, \nu)$ , where  $\mu$  is the strategy of player 1 and  $\nu$  is the strategy of player 2. Given strategy  $\pi$ , the value function and Q-function can be defined similarly as in MDP:

$$\begin{aligned} V^\pi(s) &:= \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(s^t, \pi(s^t)) \mid s^0 = s \right], \\ Q^\pi(s, a) &:= \mathbb{E} \left[ r(s^0, a^0) + \sum_{t=1}^{\infty} \gamma^t r(s^t, \pi(s^t)) \mid s^0 = s, a^0 = a \right] \\ &= r(s, a) + \gamma P(s, a) V^\pi. \end{aligned}$$

From the perspective of player 1, if the strategy  $\nu$  of player 2 is given, TBSG degenerates to an MDP, so the optimal policy against  $\nu$  exists, which we called as counterstrategy and use  $c_{\max}(\nu)$  to denote it. Similarly we can define  $c_{\min}(\mu)$  as the counterstrategy of  $\mu$  for player 2. For simplicity, we ignore the subscript in  $c_{\max}$  and  $c_{\min}$  when it is clear in the context. In addition, we define  $V^{*,\nu} := V^{c(\nu),\nu}$  and  $V^{\mu,*} := V^{\mu,c(\mu)}$  and the same for  $Q$ . By definition and property of optimal policy in MDP, we have

$$\begin{aligned} Q^{*,\nu}(s, a) &= \max_{\mu} Q^{\mu,\nu}(s, a), \forall s \in \mathcal{S}, \\ Q^{\mu,*}(s, a) &= \min_{\nu} Q^{\mu,\nu}(s, a), \forall s \in \mathcal{S}, \\ Q^{*,\nu}(s, c_{\max}(\nu)(a)) &= \max_{a'} Q^{*,\nu}(s, a'), \forall s \in \mathcal{S}_{\max}, \\ Q^{\mu,*}(s, c_{\min}(\mu)(s)) &= \min_{a'} Q^{\mu,*}(s, a'), \forall s \in \mathcal{S}_{\min}. \end{aligned}$$

Note that these are the sufficient and necessary condition of counterstrategy, which will be utilized repeatedly in our analysis.

To solve a TBSG, our goal is to find the Nash equilibrium strategy  $\pi^* = (\mu^*, \nu^*)$ , where  $\mu^* = c(\nu^*)$ ,  $\nu^* = c(\mu^*)$ . For Nash equilibrium strategy, neither player can benefit from changing its policy alone. As  $\mu^*$  and  $\nu^*$  are counterstrategy to each other, they inherit properties of counterstrategy given above. For simplicity, we will not repeat here. It is well known that in TBSG, there always exists a pure strategy as the Nash equilibrium strategy. In addition, all Nash equilibrium strategy share the same state-action value, which makes pure Nash equilibrium strategy unique given some tie selection rule, so we only consider pure strategies in our analysis.

Specifically, our target is to find an  $\epsilon$ -approximate Nash equilibrium strategy  $\pi = (\mu, \nu)$  such that

$$|Q^{\mu,*}(s, a) - Q^*(s, a)| \leq \epsilon, \forall (s, a),$$

<sup>1</sup>Our proof can be easily adapted to show that the sample complexity of learning the reward  $r$  is an order of  $\frac{1}{1-\gamma}$  smaller than learning the transition  $P$ .

$$|Q^{*,\nu}(s, a) - Q^*(s, a)| \leq \epsilon, \forall (s, a),$$

for some  $\epsilon > 0$  with as few samples as possible. Note that this is different from and stronger than the MDP analogue, which should be  $|Q^\pi(s, a) - Q^*(s, a)| \leq \epsilon$ . This slight difference makes subtle difficulty as we will show later.

### Generative Model Oracle and Plug-in Solver Approach

We assume that we have access to a generative model, where we can input an arbitrary state action pair  $(s, a)$  and receive a sampled state form  $P(\cdot|s, a)$ . Generative model oracle was introduced in [Kearns and Singh, 1999, Kakade et al., 2003]. This setting is different from the offline oracle where we can only sample trajectories via a behaviour policy and online oracle where we adaptively change the policy to explore and exploit. The advantage of generative model setting is that the exploration and exploitation is simplified, as previous work shows that treating all state-action pair equally is already optimal [Azar et al., 2013, Sidford et al., 2018]. In particular, we call the generative model  $N/|\mathcal{S}||\mathcal{A}|$  times on each state-action pair and construct the empirical TBSG  $\hat{G} = (\mathcal{S} = \mathcal{S}_{\max} \cup \mathcal{S}_{\min}, \mathcal{A}, \hat{P}, r, \gamma)$ :

$$\hat{P}(s'|s, a) = \frac{\text{count}(s, a, s')}{N/|\mathcal{S}||\mathcal{A}|}, \forall s, s' \in \mathcal{S}_{\max} \cup \mathcal{S}_{\min}, a \in \mathcal{A},$$

where  $\text{count}(s, a, s')$  is the number of times that  $s'$  is sampled from state-action pair  $(s, a)$ . It is straightforward that  $\hat{P}$  is an unbiased and maximum likelihood estimation of the true transition kernel  $P$ . We use  $\hat{\pi}^* = (\hat{\mu}^*, \hat{\nu}^*)$  to denote the Nash equilibrium strategy in the empirical MDP as well as  $\hat{V}$  and  $\hat{Q}$ . The algorithm is given in Algorithm 1.

As the transition kernel in  $\hat{G}$  is known, arbitrary planning algorithm can be used to find the empirical Nash equilibrium strategy  $\hat{\pi}^*$ . One choice is to use strategy iteration, which finds  $\hat{\pi}^*$  with in  $\tilde{O}(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-1})$  iterations [Hansen et al., 2013]. In addition, algorithms that find approximate Nash equilibrium strategy can be applied, such as QVI-MIVSS in [Sidford et al., 2020]. Note that our analysis is for the exact Nash equilibrium strategy  $\hat{\pi}^*$ , but it can be generalized to approximate Nash equilibrium strategy by using techniques in [Agarwal et al., 2019].

**Suboptimality Gap** Suboptimality gap is originated in bandit theory, which is the gap between the mean reward of the best arm and second best arm. In TBSG, we define the suboptimality gap based on optimal  $Q$ -value and second optimal  $Q$ -value.

**Definition 1.** (Suboptimality gap for Nash equilibrium strategy) A TBSG enjoys a suboptimality gap of  $\Delta$  for Nash equilibrium strategy if and only if

$$\forall s \in \mathcal{S}_{\max}, a \neq \mu^*(s) : Q^*(s, \mu^*(s)) - Q^*(s, a) \geq \Delta,$$

$$\forall s \in \mathcal{S}_{\min}, a \neq \nu^*(s) : Q^*(s, \nu^*(s)) - Q^*(s, a) \leq -\Delta,$$

and if there are two optimal actions, the gap is zero.

**Definition 2.** (Suboptimality gap for counterstrategy) A TBSG enjoys a suboptimality gap of  $\Delta$  for counter strategy to the strategy  $\nu$  of player 2 if and only if

$$\forall s \in \mathcal{S}_{\max}, a \neq c(\nu)(s), Q^*(s, c(\nu)(s)) - Q^*(s, a) \geq \Delta.$$

A TBSG enjoys a suboptimality gap of  $\Delta$  for counter strategy to the strategy  $\mu$  of player 1 if and only if

$$\forall s \in \mathcal{S}_{\min}, a \neq c(\mu)(s), Q^*(s, c(\mu)(s)) - Q^*(s, a) \leq -\Delta.$$

The suboptimality gap means that following the Nash equilibrium strategy, the expected total reward of the best action and the second best action differs at least  $\Delta$ . Intuitively, this gap quantifies the difficulty of learning the optimal action and a small gap hinders finding out the optimal action.

**Notations**  $f(x) = O(g(x))$  means that there exists a constant  $C$  such that  $f \leq Cg$  and  $f(x) = \Omega(g(x))$  means that  $g(x) = O(f(x))$ .  $\tilde{O}$  and  $\tilde{\Omega}$  is same as  $O$  and  $\Omega$  except that logarithmic factors are ignored. We use  $f \gtrsim g$  to denote that there exist some constant  $C$  such that  $f \geq Cg$  and  $f \lesssim g$  means  $f \leq Cg$  for some constant  $C$ . For a strategy  $\pi$ ,  $P^\pi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|}$  is the transition matrix induced by policy  $\pi$  and  $P^\pi(s, a)(s', a') = P(s'|s, a)\mathbf{1}(a' = \pi(s'))$  where  $\mathbf{1}(\cdot)$  is the indicator function.  $P(s, a)$  is the row vector of  $P$  that correspond to  $s, a$ . We use  $\|\cdot\|$  to denote the infinity norm  $\|\cdot\|_\infty$  and  $\sqrt{\cdot}, \leq, \geq$  are entry-wise operators.

---

**Algorithm 1: Solving TBSG via Plug-in Solver**

---

**Input:** A generative model that can output samples from distribution  $P(\cdot|s, a)$  for query  $(s, a)$ , a plug-in solver.

**Initial:** Sample size:  $N$ ;

**for**  $(s, a)$  in  $\mathcal{S} \times \mathcal{A}$  **do**

Collect  $N/|\mathcal{S}||\mathcal{A}|$  samples from  $P(\cdot|s, a)$ ;

Compute  $\hat{P}(s'|s, a) = \frac{\text{count}(s, a, s')}{N/|\mathcal{S}||\mathcal{A}|}$ ;

**end**

Construct the (perturbed) empirical TBSG;

$$\hat{\mathcal{G}} = (\mathcal{S}, \mathcal{A}, \hat{P}, r, \gamma);$$

$$\hat{\mathcal{G}}_p = (\mathcal{S}, \mathcal{A}, \hat{P}, r_p, \gamma);$$

Plan with an arbitrary plug-in solver and receive the (perturbed) empirical Nash equilibrium strategy  $\hat{\pi}^*$  ( $\hat{\pi}_p^*$ );

**Output:**  $\hat{\pi}^*$  ( $\hat{\pi}_p^*$ )

---

### 3 TECHNICAL LEMMAS FROM MDP

In this section, we present several technical lemmas that is originated in MDP analysis [Azar et al., 2013, Agarwal et al., 2019, Li et al., 2020]. These lemmas can be easily adapted to TBSG and we present them to give some intuition on our TBSG analysis.

**Lemma 1.** For any strategy  $\pi$ , we have

$$\begin{aligned} Q^\pi - \hat{Q}^\pi &= \gamma(I - P^\pi)^{-1}(\hat{P} - P)\hat{V}^\pi \\ &= \gamma(I - \hat{P}^\pi)^{-1}(P - \hat{P})V^\pi. \end{aligned}$$

This lemma portrays the concentration of  $\hat{Q}^\pi$ . Note that there are two kinds of factorization and requires different analysis. In the problem-dependent bound, we use the first one and in the problem-independent bound, we use the second one.

**Definition 3.** (One-step Variance) We define the one-step variance of a state-action pair  $(s, a)$  with respect to a certain value function  $V$  to be

$$\text{Var}_{s,a}(V) := P(s, a)V^2 - (P(s, a)V)^2,$$

which is the variance of the next state value. We define  $\text{Var}_P(V) \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$  to be a vector consisting of all  $\text{Var}_{s,a}(V)$ .

We define the one-step variance to facilitate the usage of Bernstein's inequality on term  $(\hat{P} - P)\hat{V}^\pi$ . A detailed introduction of variance in MDP can be found in Azar et al. [2013]. The following two lemmas show how to bound the one-step variance term.

**Lemma 2.** For any policy  $\pi$  and  $V^\pi$  is the value function in a MDP with transition  $P$ , we have

$$|(I - \gamma P^\pi)^{-1} \sqrt{\text{Var}_P(\hat{V}^\pi)}| \leq \sqrt{\frac{2}{(1 - \gamma)^3} + \frac{|Q^\pi - \hat{Q}^\pi|}{1 - \gamma}}.$$

**Lemma 3.** For any policy  $\pi$  and  $V^\pi$  is the value function in a MDP with transition  $P$ , if  $N \gtrsim \frac{|\mathcal{S}||\mathcal{A}|}{1 - \gamma} \log\left(\frac{1}{(1 - \gamma)\delta}\right)$ , with probability larger than  $1 - \delta$ , we have

$$|(I - \gamma \hat{P}^\pi)^{-1} \sqrt{\text{Var}_P(V^\pi)}| \leq \frac{16}{\sqrt{(1 - \gamma)^3}}.$$

Lemma 2 correspond to the first factorization in Lemma 1 and Lemma 3 is correspond to the second one. Lemma 2 is derived by utilizing the variance-Bellman-equation and Lemma 3 is derived by a Taylor expansion type analysis. If we can apply Bernstein's inequality to  $(\hat{P} - P)\hat{V}^\pi$  or  $(\hat{P} - P)V^\pi$  to generate the  $\sqrt{\text{Var}_P(\hat{V}^\pi)}$  or  $\sqrt{\text{Var}_P(V^\pi)}$  term, then with Lemma 2 or Lemma 3, we can bound  $|Q^\pi - \hat{Q}^\pi|$ .

By concentration inequalities, we can bound  $(P - \hat{P})V^\pi$  for fixed  $\pi$ . However, if  $\pi = (\hat{\mu}^*, c(\hat{\mu}^*))$  or  $\pi = (\mu^*, \hat{c}(\mu))$ , the complex statistical dependence between  $\pi$  and  $\hat{P}$  hinders the conventional concentration and subtle techniques are needed. In addition,  $(P - \hat{P})\hat{V}^\pi$  suffers from the dependence even for fixed  $\pi$ . The key contribution in the next two section is to use novel TBSG techniques to make the concentration arguments applicable to these two terms.



## 4 WARM UP: PROBLEM-DEPENDENT UPPER BOUND

In this section, we show that if a TBSG  $\mathcal{G}$  enjoys a suboptimality gap of  $\Delta$ , then with  $\tilde{O}(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-3}\Delta^{-2})$  samples, the empirical Nash equilibrium strategy  $\hat{\pi}^*$  in  $\hat{\mathcal{G}}$  is exactly the Nash equilibrium strategy  $\pi^*$  in  $\mathcal{G}$  with high probability. To begin with, we introduce a novel absorbing TBSG technique, which is motivated by the absorbing MDP technique developed in [Agarwal et al., 2019]. An absorbing TBSG  $\tilde{\mathcal{G}}_{s,a,u}$  is identical to the empirical TBSG  $\hat{\mathcal{G}}$  except that the transition distribution of a specific state-action pair  $(s, a)$  is set to be absorbing. Similar techniques have been developed for MDP [Li et al., 2020] and simultaneous game [Zhang et al., 2020].

**Definition 4.** (Absorbing TBSG) For a TBSG  $\mathcal{G} = (\mathcal{S}_{max}, \mathcal{S}_{min}, \mathcal{A}, P, r, \gamma)$  and a given state-action pair  $(s, a) \in (\mathcal{S}_{max} \cup \mathcal{S}_{min}) \times \mathcal{A}$ , the absorbing TBSG  $\tilde{\mathcal{G}}_{s,a,u} = (\mathcal{S}_{max}, \mathcal{S}_{min}, \mathcal{A}, \tilde{P}, \tilde{r}, \gamma)$ , where

$$\begin{aligned} \tilde{P}(s|s, a) &= 1, \tilde{P}(\cdot|s', a') = P(\cdot|s', a'), \forall (s', a') \neq (s, a), \\ \tilde{r}(s, a) &= u, \tilde{r}(s', a') = r(s', a'), \forall (s', a') \neq (s, a). \end{aligned}$$

**Remark 1.** Absorbing TBSG is independent of  $\hat{P}(s, a)$ , which is a kind of leave-one-out analysis. Note that  $\hat{P}(s, a)$  can be set to arbitrary fixed distribution. We use the absorbing distribution for simplicity and correspondence to its name.

For simplicity, we ignore  $(s, a)$  in absorbing TBSG when there is no misunderstanding. We use  $\tilde{\pi}_u, \tilde{Q}_u, \tilde{V}_u$  to denote the strategy, state-action value and state value in the absorbing TBSG. These terms actually depend on  $(s, a)$  and we omit  $(s, a)$  when there is no confusion. Note that  $\tilde{\mathcal{G}}_u$  has no dependence on  $\hat{P}(s, a)$ , which makes the concentration on  $(\hat{P}(s, a) - P(s, a))\tilde{V}_u$  possible. The key property of absorbing TBSG is that it can recover the Q-value of the empirical MDP by tuning the parameter  $u$ . Moreover, the Q-value of absorbing TBSG is  $\frac{1}{1-\gamma}$ -lipschitz to  $u$ , which means we can use an  $\epsilon$ -cover on the range of  $u$  to approximately recover the Q value in the empirical TBSG.

**Lemma 4.** (Properties of absorbing TBSG) Set  $u^* = r(s, a) + \gamma(P(s, a)V^*) - \gamma V^*(s)$ ,  $u^\mu = r(s, a) + \gamma(P(s, a)V^\mu) - \gamma V^\mu(s)$ . Then we have

$$\begin{aligned} Q_{u^*} &= Q^*, Q_{u^\mu} = Q^{\mu,*}, \\ |Q_{u^*} - Q_{u'}| &\leq \frac{|u - u'|}{1-\gamma}, |Q_{u^\mu} - Q_{u'^\mu}| \leq \frac{|u - u'|}{1-\gamma}. \end{aligned}$$

We can concentrate  $(\hat{P}(s, a) - P(s, a))\hat{V}^{\mu,*}$  by a union bound and an additional approximation error term as we

construct an absorbing TBSG  $\tilde{\mathcal{G}}$  on empirical TBSG  $\hat{\mathcal{G}}$ . Combining Lemma 1 and Lemma 2, we can bound  $|Q^{\mu, \hat{c}(\mu)} - \hat{Q}^{\mu, c(\mu)}|$ .

**Lemma 5.**

$$|Q^* - \hat{Q}^*| \leq \max\{|Q^{\mu, \hat{c}(\mu)} - \hat{Q}^{\mu,*}|, |Q^{\hat{c}(\nu), \nu} - \hat{Q}^{*, \nu}|\}.$$

Lemma 5 shows that we can bound the error of estimate the optimal state-action value  $|Q^* - \hat{Q}^*|$  by policy evaluation error  $|Q^{\mu, \hat{c}(\mu)} - \hat{Q}^{\mu,*}|$  and  $|Q^{\hat{c}(\nu), \nu} - \hat{Q}^{*, \nu}|$ .

Now we show how to use absorbing TBSG and suboptimality gap to prove that  $\hat{\pi}^* = \pi^*$  with large probability. The key is to prove that if  $|Q^* - \hat{Q}^*| \leq \frac{\Delta}{2}$ , then by the definition of suboptimality gap,  $\hat{Q}^*$  enjoys the gap for policy  $\pi^*$ . As  $\hat{\pi}^*$  is the Nash equilibrium policy in  $\hat{\mathcal{G}}$ , it is the only policy that can enjoy the gap, which means  $\hat{\pi}^* = \pi^*$ .

**Lemma 6.** If  $Q^*$  enjoys a suboptimality gap of  $\Delta$  and  $|Q^* - \hat{Q}^*| \leq \frac{\Delta}{2}$ , then we have  $\hat{\pi}^* = \pi^*$ .

Combining all the parts together, we get our problem-dependent result. Theorem 1 indicates that with a suboptimality gap, we can accurately recover the Nash equilibrium strategy with polynomial sample complexity. The detailed proof is provided in Appendix B.

**Theorem 1.** If  $\mathcal{G}$  enjoys a suboptimality gap of  $\Delta$  and the number of samples satisfies

$$N \geq \frac{C|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\Delta^2} \log\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)\delta\Delta}\right)$$

for some constant  $C$  and  $\Delta \in (0, (1-\gamma)^{-\frac{1}{2}}]$ , then with probability at least  $1-\delta$ , we have  $\hat{\pi}^* = \pi^*$ , which means the empirical Nash equilibrium strategy we obtained is exactly the Nash equilibrium strategy in the true TBSG.

Note that this result can recover  $\pi^*$  exactly, but at a price of restrictive sample complexity. We need to know the suboptimality gap  $\Delta$  beforehand and we have little knowledge about the empirical Nash equilibrium strategy if the sample complexity is smaller than  $\frac{C|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\Delta^2} \log\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)\delta\Delta}\right)$ . In the next section, we will give a problem-independent result, which has no dependence on the suboptimality gap of  $\mathcal{G}$ .

## 5 PROBLEM-INDEPENDENT UPPER BOUND

Two flaws exist in our problem-dependent result. One is if the suboptimality gap is considerably small, the bound  $\tilde{O}(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-3}\Delta^{-2})$  becomes meaningless. In addition, when  $N < |\mathcal{S}||\mathcal{A}|(1-\gamma)^{-3}\Delta^{-2}$ , the quality of empirical

Nash equilibrium strategy  $\hat{\pi}^*$  is completely unknown. In this section, we aim to give a minimax sample complexity result without the assumption of suboptimality gap. Interestingly, though the suboptimality gap is not assumed, we artificially create the suboptimality gap instead and the role of suboptimality gap is different from the first part analysis, which we will specify later. We now introduce the reward perturbation technique that can artificially create a suboptimality gap.

## 5.1 REWARD PERTURBATION TECHNIQUE

Here we use a reward perturbation technique to create a suboptimality gap in TBSG, which is inspired by a similar argument in MDP analysis [Li et al., 2020]. We give a proof that is different from the one in [Li et al., 2020] and our analysis for TBSG can generalize to MDP automatically as MDP is a degenerated version of TBSG. We show that by randomly perturb the reward function, with large probability, the perturbed TBSG enjoys a suboptimality gap. First we define the perturbed TBSG.

**Definition 5.** (Perturbed TBSG) For a TBSG  $\mathcal{G} = (\mathcal{S}, \mathcal{A}, P, r, \gamma)$ , the perturbed TBSG is  $\mathcal{G}_p = (\mathcal{S}, \mathcal{A}, P, r_p, \gamma)$ , where

$$r_p = r + \zeta$$

and  $\zeta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$  is a vector composed of independent random variables following uniform distribution on  $[0, \xi]$ .

We use subscript  $\pi_p$  to denote the strategy in perturbed TBSG as well as  $V_p$  and  $Q_p$ . The key property of perturbed TBSG is that it enjoy a suboptimality gap of  $\frac{\xi\delta(1-\gamma)}{4|\mathcal{S}|^2|\mathcal{A}|}$  with probability at least  $1 - \delta$ .

**Lemma 7.** (TBSG version of suboptimality gap lemma in [Li et al., 2020]) For a fixed policy  $\nu$ , with probability at least  $1 - \delta$ , the perturbed TBSG  $\mathcal{G}_p$  enjoys a suboptimality gap of  $\frac{\xi\delta(1-\gamma)}{4|\mathcal{S}|^2|\mathcal{A}|}$  for the Nash equilibrium strategy and a same gap for counterstrategy of  $\nu$ .

Our proof is substantially different from [Li et al., 2020], which consists of two important observation. Here we consider the case where only  $r(s, a_1)$  is perturbed to  $r(s, a_1) + \tau$ . First, the Nash equilibrium strategy  $\pi_\tau^*$  is a piecewise constant function of  $\tau$ , if some tie breaking rule is given. Second, we show that  $Q_\tau^*(s, a_1) = k_1(\pi_\tau)\tau + b_1(\pi_\tau)$  and  $Q_\tau^*(s, a_2) = k_2(\pi_\tau)\tau + b_2(\pi_\tau)$  are piecewise linear function and  $k_2(\pi_\tau) \leq \gamma k_1(\pi_\tau), \forall \tau$ . Intuitively, these two observation means that  $Q_\tau^*(s, a_2)$  grows at most  $\gamma$  times the speed of  $Q_\tau^*(s, a_1)$ , which further implies  $|Q_\tau^*(s, a_1) - Q_\tau^*(s, a_2)| \leq w$  can only holds for a small interval of  $\tau$ .

**Lemma 8.** Consider a TBSG  $\mathcal{G}_\tau = (\mathcal{S}, \mathcal{A}, P, r_\tau, \gamma)$  where  $r_\tau = r + \tau \mathbf{1}_{s,a}$ . Then the following facts hold.

- Given a rule to select optimal action when there are multiple optimal actions,  $\pi_\tau^*$  is a constant (vector) function of  $\tau$ .
- $Q_\tau^*$  is a piecewise linear (vector) function of  $\tau$ .
- $Q_\tau^*(s, a') = kQ_\tau^*(s, a) + b$ , where  $0 \leq k \leq \gamma$  and  $b$  are a function of  $\pi_\tau^*$ .

With the above argument, we can prove that for arbitrary  $s, a, a'$ , if we increase  $r(s, a)$ , then the growth of  $Q_\tau^*(s, a')$  is at most  $\gamma$  times of  $Q_\tau^*(s, a)$ , which means  $|Q_\tau^*(s, a') - Q_\tau^*(s, a)|$  only holds for a small range of  $r_\tau(s, a)$ . With a union bound argument, we prove the existence of suboptimality gap. The proof for counterstrategy is similar and the details are given in the appendix.

## 5.2 MINIMAX SAMPLE COMPLEXITY

In this section, we show how to use the reward perturbation technique to derive the minimax sample complexity result. First, we show that the optimal strategy in perturbed empirical TBSG is contained in a finite set that has no dependence on  $\hat{P}(s, a)$ .

**Lemma 9.** Set  $U$  to be a set of equally spaced points in  $[-\frac{1}{1-\gamma}, \frac{1}{1-\gamma}]$  and  $|U| = \frac{16|\mathcal{S}|^2|\mathcal{A}|}{(1-\gamma)^2\xi\delta}$ . We define

$$\mathcal{M}^* = \{\tilde{\mu}_{p,u}^* : u \in U\}, \mathcal{M}^{\mu^*} = \{\tilde{c}_{p,u}(\mu^*) : u \in U\}.$$

With probability at least  $1 - \delta$ , we have  $\hat{\mu}_p^* \in \mathcal{M}^*$  and  $\hat{c}_p(\mu^*) \in \mathcal{M}^{\mu^*}$ .

Lemma 9 means with large probability  $\hat{\mu}_p^*$  and  $\hat{c}_p(\mu^*)$  lie in a finite set, which is independent of  $\hat{P}(s, a)$ . This independence allows the usage of Bernstein's inequality and with the union bound, we can prove the concentration of  $(\hat{P}(s, a) - P(s, a))V_p^{\hat{\mu}_p^*,*}$  and  $(\hat{P}(s, a) - P(s, a))V_p^{\mu^*, \hat{c}_p(\mu^*)}$ . Then, with Lemma 1 and Lemma 3, we can bound  $|Q_p^{\hat{\mu}_p^*,*} - \hat{Q}_p^{\hat{\mu}_p^*, c(\hat{\mu}_p^*)}|$  and  $|Q_p^{\mu^*, \hat{c}_p(\mu^*)} - \hat{Q}_p^{\mu^*,*}|$ .

**Lemma 10.** For perturbed empirical Nash equilibrium strategy  $\hat{\pi}_p^* = (\hat{\mu}_p^*, \hat{v}_p^*)$ , we have

$$\begin{aligned} & |Q^{\hat{\mu}_p^*,*} - Q^*| \\ & \leq |Q_p^{\hat{\mu}_p^*,*} - \hat{Q}_p^{\hat{\mu}_p^*, c(\hat{\mu}_p^*)}| + |Q_p^{\mu^*, \hat{c}_p(\mu^*)} - \hat{Q}_p^{\mu^*,*}| + \frac{4\xi}{1-\gamma}. \end{aligned}$$

Lemma 10 shows how to bound  $|Q^{\hat{\mu}_p^*,*} - Q^*|$  by perturbed TBSG. In the same manner, we can bound  $|Q^{*, \hat{v}_p^*} - Q^*|$ . Selecting an appropriate  $\xi$ , we can show that perturbed empirical Nash equilibrium strategy  $\hat{\pi}_p^*$  is an  $\epsilon$ -Nash equilibrium strategy. The detailed proof is provided in Appendix C.

**Theorem 2.** *If the number of samples satisfies*

$$N \geq \frac{C|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\epsilon^2} \log\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)\delta\epsilon}\right)$$

for some constant  $C$ , then with probability at least  $1 - \delta$  and  $\epsilon \in (0, (1-\gamma)^{-1}]$ , we have that  $\hat{\pi}_p^*$  is an  $\epsilon$ -approximate Nash equilibrium strategy in  $\mathcal{G}$ .

**Remark 2.** *Compared with Theorem 3.6 in [Zhang et al., 2020], Theorem 2 extends the range of epsilon to  $(0, (1-\gamma)^{-1}]$ , which is the full possible range. In addition, Zhang et al. [2020] requires a smooth planning oracle, which is computationally inefficient. Our algorithm only needs a standard planning oracle to solve the empirical TBSG.*

In addition, we can derive the following improved result for problem-dependent bound by choosing  $\epsilon = \frac{\Delta}{2}$  and further analysis on suboptimality gap, which is provided in the Appendix C.

**Theorem 3.** *If  $\mathcal{G}$  enjoys a suboptimality gap of  $\Delta$  and the number of samples satisfies*

$$N \geq \frac{C|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\Delta^2} \log\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)\delta\Delta}\right)$$

for some constant  $C$  and  $\Delta \in (0, (1-\gamma)^{-1}]$ , then with probability at least  $1 - \delta$ , we have  $\hat{\pi}_p^* = \pi^*$ , which means the empirical Nash equilibrium strategy we obtained is exactly the Nash equilibrium strategy in the true TBSG.

Theorem 3 is strictly stronger than Theorem 1. The difference between Theorem 1 and Theorem 3 is that the range of suboptimality gap  $\Delta$  is extended to  $(0, (1-\gamma)^{-1}]$ , which is the full possible range of suboptimality gap.

## 6 RELATED LITERATURE

**TBSG** TBSG has been widely studied since [Shapley, 1953]. For a detailed introduction of stochastic game, readers can refer to [Neyman et al., 2003]. In the old days, people focus on dynamic programming type algorithms to solve TBSG. Strategy iteration, as the counterpart of value iteration in MDP and parallelized simplex method, is proved to be a strong polynomial time algorithm [Hansen et al., 2013, Jia et al., 2020]. Reinforcement learning approach has been studied recently for TBSG to relieve the high computational cost of dynamic programming. Several works are proposed for the generative model setting, [Sidford et al., 2020, Jia et al., 2019, Zhang et al., 2020]. [Sidford et al., 2020] first gives a sample efficient algorithm for tabular TBSG, while their result achieves minimax sample complexity  $\tilde{O}(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-3}\epsilon^{-2})$  only for  $\epsilon \in (0, 1]$ . [Jia et al., 2019] adapts the MDP algorithm in [Sidford et al., 2018] for feature-based TBSG, but leaves a gap of  $\frac{1}{1-\gamma}$  between optimal sample complexity. [Cui and Yang, 2020] uses a similar

algorithm as ours in feature-based TBSG, while their result only holds for  $\epsilon$ -Nash equilibrium value, which results in a  $\frac{1}{(1-\gamma)^2}$  gap in finding  $\epsilon$ -Nash equilibrium strategy. [Zhang et al., 2020] considers simultaneous stochastic game, and their approach consists of solving a regularized simultaneous stochastic game, which is computationally costly. For the online sampling setting, a recent work [Bai et al., 2020] uses an upper confidence bound algorithms that can find an approximate Nash equilibrium strategy in  $\tilde{O}(|\mathcal{S}||\mathcal{A}||\mathcal{B}|)$  steps.

**Generative Model** Generative model is a sampling oracle setting in MDP, which has been shown to simplify the exploration and exploitation tradeoff. This concept is formalized in [Kakade et al., 2003] and a  $\tilde{O}(|\mathcal{S}||\mathcal{A}|\text{poly}((1-\gamma)^{-1})\epsilon^{-2})$  sample complexity has been proved there. [Azar et al., 2013] proves the minimax sample complexity  $\tilde{O}(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-3}\epsilon^{-2})$ . However, the upper bound there is only for  $\epsilon \in (1-\gamma)^{-1/2}|\mathcal{S}|^{-1/2}$ . Many works have devoted to improve the dependence on  $\epsilon$ . Recently, [Sidford et al., 2018] gives a minimax model-free algorithm for  $\epsilon \in (0, 1]$  and [Agarwal et al., 2019] gives a minimax model-based algorithm for  $\epsilon \in (0, (1-\gamma)^{-1/2}]$ . Finally, [Li et al., 2020] uses a perturbed MDP technique to prove minimax sample complexity with full range of  $\epsilon$ .

**Suboptimality Gap** Suboptimality gap originated in bandit theory. Multi-armed bandits and linear bandits enjoy a logarithmic gap-dependent regret and a square root gap-independent regret [Auer et al., 2002, Abbasi-Yadkori et al., 2011]. MDP with suboptimality gap have been studied in [Auer et al., 2009] and recently  $\tilde{O}(|\mathcal{S}||\mathcal{A}|\text{poly}(H)\log(T))$  regret has been proved for both model-based and model free algorithms [Simchowitz and Jamieson, 2019, Yang et al., 2020]. [Du et al., 2020] utilized the suboptimality gap in general function approximation setting and proved that optimal policy can be found in  $\tilde{O}(\text{dim}_E)$  trajectories in deterministic MDP. Most of the gap-dependent analysis in MDP focus on online RL and to the best of our knowledge, we are the first to study this notion in TBSG with a generative model.

## 7 CONCLUSION

In this work, we completely solve the sample complexity problem of TBSG with generative model oracle. We prove that the simplest model-based algorithm, plug-in solver approach, is minimax sample optimal for full range of  $\epsilon$  by using absorbing TBSG and reward perturbation techniques. Our proof is based on suboptimality gap, a notion originated from bandit theory and receives great attention in RL. We believe that our work can shed some light on suboptimality gap and TBSG.

## References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011.
- Alekh Agarwal, Sham Kakade, and Lin F Yang. On the optimality of sparse model-based planning for markov decision processes. *arXiv preprint arXiv:1906.03804*, 2019.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- Peter Auer, Thomas Jaksch, and Ronald Ortner. Near-optimal regret bounds for reinforcement learning. In *Advances in neural information processing systems*, pages 89–96, 2009.
- Mohammad Gheshlaghi Azar, Rémi Munos, and Hilbert J Kappen. Minimax pac bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, 91(3):325–349, 2013.
- Yu Bai, Chi Jin, and Tiancheng Yu. Near-optimal reinforcement learning with self-play. *arXiv preprint arXiv:2006.12007*, 2020.
- Qiwen Cui and Lin F. Yang. Is plug-in solver sample-efficient for feature-based reinforcement learning?, 2020.
- Simon S Du, Jason D Lee, Gaurav Mahajan, and Ruosong Wang. Agnostic q-learning with function approximation in deterministic systems: Tight bounds on approximation error and sample complexity. *arXiv preprint arXiv:2002.07125*, 2020.
- Thomas Dueholm Hansen, Peter Bro Miltersen, and Uri Zwick. Strategy iteration is strongly polynomial for 2-player turn-based stochastic games with a constant discount factor. *Journal of the ACM (JACM)*, 60(1):1–16, 2013.
- Zeyu Jia, Lin F Yang, and Mengdi Wang. Feature-based q-learning for two-player stochastic games. *arXiv preprint arXiv:1906.00423*, 2019.
- Zeyu Jia, Zaiwen Wen, and Yinyu Ye. Towards solving 2-tbgs efficiently. *Optimization Methods and Software*, 35(4):706–721, 2020.
- Lukasz Kaiser, Mohammad Babaeizadeh, Piotr Milos, Blazej Osinski, Roy H Campbell, Konrad Czechowski, Dumitru Erhan, Chelsea Finn, Piotr Kozakowski, Sergey Levine, et al. Model-based reinforcement learning for atari. *arXiv preprint arXiv:1903.00374*, 2019.
- Sham Machandranath Kakade et al. *On the sample complexity of reinforcement learning*. PhD thesis, University of London London, England, 2003.
- Michael J Kearns and Satinder P Singh. Finite-sample convergence rates for q-learning and indirect algorithms. In *Advances in neural information processing systems*, pages 996–1002, 1999.
- Gen Li, Yuting Wei, Yuejie Chi, Yuantao Gu, and Yuxin Chen. Breaking the sample size barrier in model-based reinforcement learning with a generative model. *arXiv preprint arXiv:2005.12900*, 2020.
- Abraham Neyman, Sylvain Sorin, and S Sorin. *Stochastic games and applications*, volume 570. Springer Science & Business Media, 2003.
- Lloyd S Shapley. Stochastic games. *Proceedings of the national academy of sciences*, 39(10):1095–1100, 1953.
- Aaron Sidford, Mengdi Wang, Xian Wu, and Yinyu Ye. Variance reduced value iteration and faster algorithms for solving markov decision processes. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 770–787. SIAM, 2018.
- Aaron Sidford, Mengdi Wang, Lin Yang, and Yinyu 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.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- Max Simchowitz and Kevin G Jamieson. Non-asymptotic gap-dependent regret bounds for tabular mdps. In *Advances in Neural Information Processing Systems*, pages 1153–1162, 2019.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- Gerald Tesauro. Temporal difference learning and td-gammon. *Communications of the ACM*, 38(3):58–68, 1995.
- Tingwu Wang, Xuchan Bao, Ignasi Clavera, Jerrick Hoang, Yeming Wen, Eric Langlois, Shunshi Zhang, Guodong Zhang, Pieter Abbeel, and Jimmy Ba. Benchmarking model-based reinforcement learning. *arXiv preprint arXiv:1907.02057*, 2019.
- Kunhe Yang, Lin F Yang, and Simon S Du. q-learning with logarithmic regret. *arXiv preprint arXiv:2006.09118*, 2020.



Deheng Ye, Zhao Liu, Mingfei Sun, Bei Shi, Peilin Zhao, Hao Wu, Hongsheng Yu, Shaojie Yang, Xipeng Wu, Qingwei Guo, et al. Mastering complex control in moba games with deep reinforcement learning. In *AAAI*, pages 6672–6679, 2020.

Kaiqing Zhang, Sham M Kakade, Tamer Başar, and Lin F Yang. Model-based multi-agent rl in zero-sum markov games with near-optimal sample complexity. *arXiv preprint arXiv:2007.07461*, 2020.

## A TECHNICAL LEMMAS FROM MDP

**Additional Notations** We use  $\Pi^\pi : \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} \rightarrow \mathbb{R}^{|\mathcal{S}|}$  to denote the projection operator with respect to policy  $\pi$ , which means if  $V = \Pi^\pi Q$ , then  $V(s) = Q(s, \pi(s)), \forall s \in \mathcal{S}$ . We use  $\hat{\mathcal{G}}$  to denote the empirical TBSG,  $\tilde{\mathcal{G}}$  to denote the absorbing TBSG on empirical TBSG and  $\tilde{\mathcal{G}}_p$  to denote the perturbed TBSG. We use  $\mathbf{1}$  to denote a vector with all entries to be 1 and  $\mathbf{1}_{s,a}$  to denote a zero vector with only  $(s, a)$  entry to be 1.

*Proof of Lemma 1.* By Bellman equation, we have  $Q^\pi = (I - \gamma P^\pi)^{-1}r$  and  $\hat{Q}^\pi = (I - \gamma \hat{P}^\pi)^{-1}r$ . For any policy  $\pi$ , we have

$$\begin{aligned} Q^\pi - \hat{Q}^\pi &= (I - \gamma P^\pi)^{-1}r - (I - \gamma \hat{P}^\pi)^{-1}r \\ &= (I - \gamma P^\pi)^{-1}((I - \gamma \hat{P}^\pi) - (I - \gamma P^\pi))\hat{Q}^\pi \\ &= \gamma(I - \gamma P^\pi)^{-1}(P^\pi - \hat{P}^\pi)\hat{Q}^\pi \\ &= \gamma(I - \gamma P^\pi)^{-1}(P - \hat{P})\hat{V}^\pi. \end{aligned}$$

Similarly, we have  $Q^\pi - \hat{Q}^\pi = \gamma(I - \gamma \hat{P}^\pi)^{-1}(\hat{P} - P)V^\pi$ .

□

**Lemma 11.** (Lemma 3 in [Sidford et al., 2020]) For any policy  $\pi$ , we have

$$\left| (I - \gamma P^\pi)^{-1} \sqrt{\text{Var}_P(V^\pi)} \right| \leq \sqrt{\frac{2}{(1-\gamma)^3}}.$$

*Proof of Lemma 2.* We have

$$\begin{aligned} \left| (I - \gamma P^\pi)^{-1} \sqrt{\text{Var}_P(\hat{V}^\pi)} \right| &\leq \left| (I - \gamma P^\pi)^{-1} \sqrt{\text{Var}_P(V^\pi)} \right| + \left| (I - \gamma P^\pi)^{-1} \sqrt{\text{Var}_P(\hat{V}^\pi - V^\pi)} \right| \\ &\leq \sqrt{\frac{2}{(1-\gamma)^3}} + \frac{|\hat{V}^\pi - V^\pi|}{1-\gamma} \\ &\leq \sqrt{\frac{2}{(1-\gamma)^3}} + \frac{|\hat{Q}^\pi - Q^\pi|}{1-\gamma}, \end{aligned}$$

where the first inequality is the triangle inequality, the second one is from Lemma 11 and the fact that  $|(I - \gamma P^\pi)^{-1}V| \leq \frac{|V|}{1-\gamma}$ , and the last one is because  $V^\pi$  is a subset of  $Q^\pi$ .

□

**Lemma 12.** (Lemma 8 in [Li et al., 2020]) Let  $Q$  be a vector obeying  $Q = (I - \gamma P^\pi)^{-1}r$  for some vector  $r > 0$  and  $V = \Pi^\pi Q$ , then we have

$$\left| (I - \gamma P^\pi)^{-1} \sqrt{\text{Var}_P(V)} \right| \leq \frac{4}{\gamma\sqrt{1-\gamma}}|Q|.$$

*Proof of Lemma 3.* We use a Taylor expansion type analysis. For simplicity, we define

$$r^{(n)} = \sqrt{\text{Var}_P(V^{(n-1)})}, Q^{(n)} = (I - \gamma P^\pi)^{-1}r^{(n)}, V^{(n)} = \Pi^\pi Q^{(n)}, n = 1, 2, \dots$$

and  $V^{(0)} = V^\pi$ . Then with large probability we have

$$\begin{aligned}
|(I - \gamma \widehat{P}^\pi)^{-1} \sqrt{\text{Var}_P(V^\pi)}| &= |(I - \gamma \widehat{P}^\pi)^{-1} r^{(1)}| \\
&= |(I - \gamma P^\pi)^{-1} r^{(1)} + (I - \gamma \widehat{P}^\pi)^{-1} (\gamma \widehat{P}^\pi - \gamma P^\pi) (I - \gamma P^\pi)^{-1} r^{(1)}| \\
&= |(I - \gamma P^\pi)^{-1} r^{(1)} + (I - \gamma \widehat{P}^\pi)^{-1} (\gamma \widehat{P}^\pi - \gamma P^\pi) Q^{(1)}| \\
&\lesssim |(I - \gamma P^\pi)^{-1} r^{(1)}| + \frac{\gamma}{\sqrt{N}} |(I - \gamma \widehat{P}^\pi)^{-1} \sqrt{\text{Var}_P(V^{(1)})}| \\
&= |(I - \gamma P^\pi)^{-1} r^{(1)}| + \frac{\gamma}{\sqrt{N}} |(I - \gamma \widehat{P}^\pi)^{-1} r^{(2)}| \\
&\lesssim \dots \\
&\lesssim \sum_{i=1}^n \left(\frac{\gamma}{\sqrt{N}}\right)^{i-1} |(I - \gamma P^\pi)^{-1} r^{(i)}| + \left(\frac{\gamma}{\sqrt{N}}\right)^n |(I - \gamma \widehat{P}^\pi)^{-1} r^{(n+1)}|
\end{aligned}$$

In addition, by Lemma B, we have

$$\begin{aligned}
|(I - \gamma P^\pi)^{-1} r^{(i)}| &= \left| (I - \gamma P^\pi)^{-1} \sqrt{\text{Var}_P(V^{(i-1)})} \right| \\
&\leq \frac{4}{\gamma \sqrt{1-\gamma}} |Q^{i-1}| \\
&\leq \frac{4}{\gamma \sqrt{1-\gamma}} \left| (I - \gamma \widehat{P}^\pi)^{-1} r^{(i-1)} \right|.
\end{aligned}$$

By induction, we have  $|(I - \gamma P^\pi)^{-1} r^{(i)}| \leq \left(\frac{4}{\gamma \sqrt{1-\gamma}}\right)^i \left| (I - \gamma \widehat{P}^\pi)^{-1} r^{(0)} \right| \leq \left(\frac{4}{\gamma \sqrt{1-\gamma}}\right)^i \frac{1}{1-\gamma}$ . If we set  $N \gtrsim \frac{64}{1-\gamma}$  and  $n \gtrsim \log\left(\frac{1}{1-\gamma}\right)$ , we have

$$\left| (I - \gamma \widehat{P}^\pi)^{-1} \sqrt{\text{Var}_P(V^\pi)} \right| \leq \frac{4}{\gamma \sqrt{(1-\gamma)^3}} \sum_{i=1}^n \left(\frac{4}{\sqrt{N(1-\gamma)}}\right)^{i-1} + \frac{1}{(1-\gamma)^2} \left(\frac{4}{\sqrt{N(1-\gamma)}}\right)^n \leq \frac{16}{\sqrt{(1-\gamma)^3}}$$

For a more detailed proof on lower order term, one can refer to [Li et al., 2020].  $\square$

## B PROBLEM-DEPENDENT UPPER BOUND

*Proof of Lemma 4.* First, we show that choosing  $u = u^*$  in absorbing TBSG can recover  $\widehat{Q}^*$ .

$$\begin{aligned}
\widetilde{Q}_{u^*}^{\widehat{\pi}^*} &= (I - \gamma \widetilde{P}^{\widehat{\pi}^*})^{-1} (r + (u^* - r(s, a)) \mathbf{1}_{s,a}) \\
&= (I - \gamma \widetilde{P}^{\widehat{\pi}^*})^{-1} ((I - \gamma \widehat{P}^{\widehat{\pi}^*}) \widehat{Q}^* + \gamma (\widehat{P} \widehat{V}^* - \widetilde{P} \widehat{V}^*)) \\
&= \widehat{Q}^*.
\end{aligned}$$

By the property of Nash equilibrium strategy, we have  $\widehat{Q}^* = \widetilde{Q}_{u^*}^{\widehat{\pi}^*} = \widetilde{Q}_{u^*}^{\pi^*}$ . Specifically,  $\widehat{\pi}^*$  select the optimal value in  $\widehat{Q}^*$ , which means it select the optimal value in  $\widetilde{Q}_{u^*}^{\widehat{\pi}^*}$ , and this is equivalent to  $\widehat{\pi}^* = \widetilde{\pi}_{u^*}^*$ .

Second, we show that the optimal Q-value in absorbing TBSG is  $\frac{1}{1-\gamma}$ -lipschitz to  $u$ . This is proved by bounding the distance between an upper bound and an lower bound.

$$\begin{aligned}
\left| \widetilde{Q}_{u^*, \widetilde{V}_{u'}}^{\mu_u^*, \widetilde{V}_{u'}} - \widetilde{Q}_{u', \widetilde{V}_{u'}}^{\mu_u^*, \widetilde{V}_{u'}} \right| &= \left| (I - \widetilde{P}^{\mu_u^*, \widetilde{V}_{u'}})^{-1} (r_u - r_{u'}) \right| \\
&= \left| (u - u') (I - \widetilde{P}^{\mu_u^*, \widetilde{V}_{u'}})^{-1} \mathbf{1}_{s,a} \right| \\
&\leq \frac{|u - u'|}{1-\gamma}.
\end{aligned}$$

Similarly, we have  $\left| \tilde{Q}_{u'}^{\mu^*, \tilde{\nu}^*} - \tilde{Q}_{u'}^{\mu^*, \tilde{\nu}^*} \right| \leq \frac{|u-u'|}{1-\gamma}$ . By the definition of  $Q^*$ , we have  $\tilde{Q}_{u'}^{\mu^*, \tilde{\nu}^*} \leq \tilde{Q}_u^* \leq \tilde{Q}_{u'}^{\mu^*, \tilde{\nu}^*}$  and  $\tilde{Q}_{u'}^{\mu^*, \tilde{\nu}^*} \leq \tilde{Q}_{u'}^* \leq \tilde{Q}_{u'}^{\mu^*, \tilde{\nu}^*}$ . Thus we have  $|\tilde{Q}_u^* - \tilde{Q}_{u'}^*| \leq \frac{|u-u'|}{1-\gamma}$ . The proof for counterstrategy is the same.  $\square$

**Lemma 13.** For any fixed value vector  $V$ , with probability larger than  $1 - \delta$ , we have

$$|(P(s, a) - \hat{P}(s, a))V| \leq \sqrt{2 \log(4/\delta)} \sqrt{\frac{\text{Var}_{s,a}(V)}{N}} + \frac{2 \log(4/\delta)}{3(1-\gamma)N},$$

where  $N$  is the sample size from distribution  $P(\cdot|s, a)$ .

*Proof.* The proof is a direct application of Bernstein's inequality.  $\square$

**Lemma 14.** For  $C = 2 \log(\frac{32}{(1-\gamma)^2 \epsilon \delta})$ , with probability larger than  $1 - \delta$ , we have

$$\left| (P(s, a) - \hat{P}(s, a)) \hat{V}^{\mu^*, *} \right| \leq \sqrt{\frac{C \text{Var}_{s,a}(\hat{V}^{\mu^*, *})}{N}} + \frac{C}{3(1-\gamma)N} + \left( \sqrt{\frac{C}{N}} + 1 \right) \frac{\epsilon(1-\gamma)}{4},$$

where  $N$  is the sample size from distribution  $P(\cdot|s, a)$ .

*Proof.* We define a fixed set  $U$  to be evenly spaced points in  $[-\frac{1}{1-\gamma}, \frac{1}{1-\gamma}]$  such that  $|U| = \frac{8}{(1-\gamma)^2 \epsilon}$ . Combining Lemma with union bound, with probability larger than  $1 - \delta$ , we have for all  $u \in U$ ,

$$\left| (P(s, a) - \hat{P}(s, a)) \tilde{V}_u^{\mu^*, *} \right| \leq \sqrt{2 \log(4|U|/\delta)} \sqrt{\frac{\text{Var}_{s,a}(\tilde{V}_u^{\mu^*, *})}{N}} + \frac{2 \log(4|U|/\delta)}{3(1-\gamma)N}.$$

For each  $u \in U$ , we have

$$\begin{aligned} \left| (P(s, a) - \hat{P}(s, a)) \hat{V}^{\mu^*, *} \right| &= \left| (P(s, a) - \hat{P}(s, a)) \tilde{V}_{u^{\mu^*}}^{\mu^*, *} \right| \\ &\leq \left| (P(s, a) - \hat{P}(s, a)) \tilde{V}_u^{\mu^*, *} \right| + \frac{|u - u^{\mu^*}|}{1-\gamma} \\ &\leq \sqrt{2 \log(4|U|/\delta)} \sqrt{\frac{\text{Var}_{s,a}(\tilde{V}_u^{\mu^*, *})}{N}} + \frac{2 \log(4|U|/\delta)}{3(1-\gamma)N} + \frac{|u - u^{\mu^*}|}{1-\gamma} \\ &\leq \sqrt{2 \log(4|U|/\delta)} \sqrt{\frac{\text{Var}_{s,a}(\tilde{V}_{u^{\mu^*}}^{\mu^*, *})}{N}} + \frac{2 \log(4|U|/\delta)}{3(1-\gamma)N} + \left( 1 + \sqrt{\frac{2 \log(4|U|/\delta)}{N}} \right) \frac{|u - u^{\mu^*}|}{1-\gamma} \end{aligned}$$

As there exist a  $u \in U$  such that  $|u - u^{\mu^*}| \leq \frac{\epsilon(1-\gamma)^2}{4}$ , we can prove the argument.  $\square$

*Proof of Lemma 5.* By the property of counterstrategy, we have

$$\begin{aligned} Q^* - \hat{Q}^* &= Q^* - Q^{\hat{c}(\nu^*), \nu^*} + Q^{\hat{c}(\nu^*), \nu^*} - \hat{Q}^{*, \nu^*} + \hat{Q}^{*, \nu^*} - \hat{Q}^* \\ &\geq Q^{\hat{c}(\nu^*), \nu^*} - \hat{Q}^{*, \nu^*}. \end{aligned}$$

Similarly, we have  $Q^* - \hat{Q}^* \leq Q^{\mu^*, \hat{c}(\mu^*)} - \hat{Q}^{\mu^*, *}$ . Together we can prove Lemma 5.  $\square$

*Proof of Lemma 6.* By the definition of suboptimality gap, for all  $s \in \mathcal{S}_{\max}$ ,  $a \in \mathcal{A}$  we have

$$\begin{aligned} \hat{Q}^*(s, \mu^*(s)) - \hat{Q}^*(s, a) &= \hat{Q}^*(s, \mu^*(s)) - Q^*(s, \mu^*(s)) + Q^*(s, \mu^*(s)) - Q^*(s, a) + Q^*(s, a) - \hat{Q}^*(s, a) \\ &\geq -\frac{\Delta}{2} + \Delta - \frac{\Delta}{2} \\ &\geq 0. \end{aligned}$$



Similar, for all  $s \in \mathcal{S}_{\min}$ ,  $a \in \mathcal{A}$ , we have

$$\widehat{Q}^*(s, \nu^*(s)) - \widehat{Q}^*(s, a) < 0.$$

Note that the empirical Nash equilibrium strategy  $\widehat{\pi}^*$  is the only policy that can satisfy the above two conditions, which means  $\pi^* = \widehat{\pi}^*$ .  $\square$

*Proof of Theorem 1.* We set  $N = \frac{32|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\epsilon^2} \log(\frac{32}{(1-\gamma)^2\epsilon\delta})$ . With probability larger than  $1 - \delta$ , We have

$$\begin{aligned} & \left| Q^{\mu^*, \widehat{c}(\mu^*)} - \widehat{Q}^{\mu^*, *} \right| \\ &= \left| \gamma(I - \gamma P^{\mu^*, \widehat{c}(\mu^*)})^{-1} (P - \widehat{P}) \widehat{V}^{\mu^*, \widehat{c}(\mu^*)} \right| \\ &\leq \left| \gamma(I - \gamma P^{\mu^*, \widehat{c}(\mu^*)})^{-1} \left\{ \sqrt{\frac{C \text{Var}_P(\widehat{V}^{\mu^*, *})}{N/|\mathcal{S}||\mathcal{A}|}} + \left[ \frac{C}{3(1-\gamma)N/|\mathcal{S}||\mathcal{A}|} + \left( \sqrt{\frac{C}{N/|\mathcal{S}||\mathcal{A}|}} + 1 \right) \frac{\epsilon(1-\gamma)}{4} \right] \mathbf{1} \right\} \right| \\ &\leq \sqrt{\frac{2C}{(1-\gamma)^3 N/|\mathcal{S}||\mathcal{A}|}} + \sqrt{\frac{C}{N/|\mathcal{S}||\mathcal{A}|}} \frac{|Q^{\mu^*, \widehat{c}(\mu^*)} - \widehat{Q}^{\mu^*, *}|}{1-\gamma} + \frac{C}{3(1-\gamma)^2 N/|\mathcal{S}||\mathcal{A}|} + \left( \sqrt{\frac{C}{N/|\mathcal{S}||\mathcal{A}|}} + 1 \right) \frac{\epsilon}{4} \end{aligned}$$

where the first equality is due to Lemma 1, the first inequality is due to Lemma 14, the second inequality is due to Lemma 2.

Thus for  $\epsilon \leq \frac{1}{\sqrt{1-\gamma}}$  and  $N \geq \frac{32|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\epsilon^2} \log(\frac{32}{(1-\gamma)^2\epsilon\delta})$ , we have

$$\left| Q^{\mu^*, \widehat{c}(\mu^*)} - \widehat{Q}^{\mu^*, *} \right| \leq \epsilon.$$

Similarly, we have  $|Q^{\widehat{c}(\nu^*), \nu^*} - \widehat{Q}^{*, \nu^*}| \leq \epsilon$ . Finally, we set  $\epsilon = \frac{\Delta}{2}$ . Then by Lemma 5 and Lemma 6, we can conclude that with probability  $1 - \delta$  and  $N = \frac{128|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\Delta^2} \log(\frac{128}{(1-\gamma)^2\Delta\delta})$ , the empirical Nash equilibrium strategy  $\widehat{\pi}^*$  is exactly the true Nash equilibrium strategy  $\pi^*$ .  $\square$

## C PROBLEM-INDEPENDENT UPPER BOUND

*Proof of Lemma 8.* As  $Q_\tau^* = \max_\pi Q_\tau^\pi = \max_\pi (I - \gamma P^\pi)^{-1} r_\tau$  is a continuous function of  $\tau$ , with a tie breaking rule to select optimal action if there are multiple optimal actions, the optimal policy  $\pi_\tau^*$  is a piecewise constant function.

For a given policy  $\pi$ , the corresponding Q-function satisfies  $Q_\tau^\pi = (I - \gamma P^\pi)^{-1} (r + \tau \mathbf{1}_{s,a})$ , which means  $Q_\tau^\pi$  is a linear function of  $\tau$ . As  $\pi_\tau^*$  is a piecewise constant function, we have that  $Q_\tau^*$  is a piecewise linear function.

For the last argument, we factorize  $Q_\tau^*(s, a')$  according to the hitting time of state  $s$ .

$$\begin{aligned} Q_\tau^*(s, a') &= r(s, a') + \gamma P(s, a') V_\tau^* \\ &= \gamma p_1 V_\tau(s) + \gamma \sum_{s' \neq s} P(s'|s, a') V_\tau^*(s') + r(s, a') \\ &= \gamma p_1 V_\tau(s) + \gamma^2 p_2 V_\tau(s) + \gamma^2 \sum_{s'', s' \neq s} P(s'|s, a') P(s''|s', \pi^*(s')) V_\tau^*(s'') + r(s, a') + \gamma \sum_{s' \neq s} P(s'|s, a) r(s', \pi^*(s')) \\ &= \dots \\ &= \sum_{n=1}^{\infty} \gamma^n p_n V_\tau^*(s) + b(\pi_\tau^*), \end{aligned}$$

where  $p_n$  is the probability of first visiting state  $s$  in step  $n$  under optimal policy and  $b(\pi_\tau^*)$  is a function of  $\pi_\tau^*$ . For  $a = \pi_\tau^*(s)$ , we define  $k = \sum_{n=1}^{\infty} \gamma^n p_n \leq \gamma$ , then  $Q_\tau^*(s, a') = k Q_\tau^*(s, a) + b(\pi_\tau^*)$ . When  $\pi_\tau^*(s) \neq a$ ,  $Q_\tau^*(s, a') = r(s, a') + \gamma P(s, a') (I - \gamma P_{\pi_\tau^*})^{-1} r_{\pi_\tau^*}$  is a function of  $\pi_\tau^*$ , which means we have  $Q_\tau^*(s, a') = 0 Q_\tau^*(s, a) + b(\pi_\tau^*)$ .  $\square$

*Proof of Lemma 7.* We prove that for any state  $s$  and actions  $a, a'$ , with large probability, a gap between  $Q_p^*(s, a)$  and  $Q_p^*(s, a')$  exist. Then with a union bound, Lemma 7 holds with large probability.

Now we fix all rewards except  $r_p(s, a) = r(s, a) + \tau$ . We define  $\mathcal{I}_w = \{\tau \mid |Q_\tau^*(s, a) - Q_\tau^*(s, a')| \leq w\}$ . First, we have  $Q_{\tau_1}^*(s, a) - Q_{\tau_2}^*(s, a) = \tau_1 - \tau_2 + P(s, a)(V_{\tau_1}^* - V_{\tau_2}^*) \geq \tau_1 - \tau_2, \forall \tau_1 \geq \tau_2$ . In addition, Lemma 8 implies that the growth rate of  $Q_\tau^*(s, a')$  is at most  $\gamma$  times the rate of  $Q_\tau(s, a)$ . Thus we can conclude that the length of  $\mathcal{I}_w$  is at most  $\frac{2w}{1-\gamma}$ . We set  $w = \frac{\delta\xi(1-\gamma)}{2|\mathcal{S}||\mathcal{A}|^2}$ , then gaps between all  $Q_p^*(s, a)$  and  $Q_p^*(s, a')$  exist, which implies Lemma 7. The proof for counterstrategy is the same.  $\square$

*Proof of Lemma 9.* Select  $u \in U$  such that  $u - u^* \leq \frac{\delta\xi(1-\gamma)^2}{4|\mathcal{S}||\mathcal{A}|^2}$ . Thus we have  $|\widehat{Q}_p^* - \widetilde{Q}_{p,u}^*| \leq \frac{\delta\xi(1-\gamma)}{4|\mathcal{S}||\mathcal{A}|^2}$ . As with probability  $1 - \delta$ , a gap of  $\frac{\delta\xi(1-\gamma)}{2|\mathcal{S}||\mathcal{A}|^2}$  exist in  $\widehat{Q}_p^*$ , by Lemma 6, we have  $\widehat{\pi}_p^* = \widetilde{\pi}_{p,u}^*$ . The proof for  $\widehat{c}_p(\mu^*)$  is the same.  $\square$

*Proof of Lemma 10.* First, we consider the unperturbed case.

$$\begin{aligned} 0 &\leq Q^* - Q^{\widehat{\mu}^*,*} \\ &= Q^* - Q^{\mu^*, \widehat{c}(\mu^*)} + Q^{\mu^*, \widehat{c}(\mu^*)} - \widehat{Q}^{\mu^*,*} + \widehat{Q}^{\mu^*,*} - \widehat{Q}^* + \widehat{Q}^* - \widehat{Q}^{\widehat{\mu}^*, c(\widehat{\mu}^*)} + \widehat{Q}^{\widehat{\mu}^*, c(\widehat{\mu}^*)} - Q^{\widehat{\mu}^*,*} \\ &\leq Q^{\mu^*, \widehat{c}(\mu^*)} - \widehat{Q}^{\mu^*,*} + \widehat{Q}^{\widehat{\mu}^*, c(\widehat{\mu}^*)} - Q^{\widehat{\mu}^*,*}. \end{aligned}$$

Then we show that the perturbation only induce an error of  $\frac{\xi}{1-\gamma}$ .

$$\begin{aligned} \left| Q^{\mu^*, \nu_p^*} - Q_p^{\mu^*, \nu_p^*} \right| &= \left| (I - \gamma P^{\mu^*, \nu_p^*})^{-1} (r - r_p) \right| \\ &\leq \frac{|r - r_p|}{1-\gamma} \\ &\leq \frac{\xi}{1-\gamma} \end{aligned}$$

Similarly, we have  $\left| Q_p^{\mu_p^*, \nu^*} - Q^{\mu_p^*, \nu^*} \right| \leq \frac{\xi}{1-\gamma}$ . As  $Q_p^{\mu_p^*, \nu^*} \leq Q^* \leq Q^{\mu_p^*, \nu_p^*}$  and  $Q_p^{\mu_p^*, \nu^*} \geq Q_p^* \geq Q_p^{\mu_p^*, \nu_p^*}$ , we have  $|Q_p^* - Q^*| \leq \frac{\xi}{1-\gamma}$ . Thus  $|Q_p^* - Q^*| \leq \frac{\xi}{1-\gamma}$ . Similarly, we have  $|Q_p^{\widehat{\mu}^*,*} - Q^{\widehat{\mu}^*,*}| \leq \frac{\xi}{1-\gamma}$ . Combining all parts together, we have

$$\begin{aligned} \left| Q^{\widehat{\mu}^*,*} - Q^* \right| &\leq \left| Q^{\widehat{\mu}^*,*} - \widehat{Q}_p^{\widehat{\mu}^*,*} \right| + |Q^* - Q_p^*| + \left| Q_p^{\widehat{\mu}^*,*} - Q^* \right| \\ &\leq \frac{2\xi}{1-\gamma} + \left| Q_p^{\mu_p^*, \widehat{c}_p(\mu_p^*)} - \widehat{Q}_p^{\mu_p^*,*} \right| + \left| \widehat{Q}_p^{\widehat{\mu}^*, c_p(\widehat{\mu}^*)} - \widehat{Q}_p^{\widehat{\mu}^*,*} \right| \end{aligned}$$

Combining these inequalities and we can get the proof.  $\square$

*Proof of Theorem 2.* By Lemma 1, we have

$$Q_p^{\mu_p^*, \widehat{c}_p(\mu_p^*)} - \widehat{Q}_p^{\mu_p^*,*} = \gamma(I - \gamma \widehat{P}^{\mu_p^*, \widehat{c}_p(\mu_p^*)})^{-1} (P - \widehat{P}) V_p^{\mu_p^*, \widehat{c}_p(\mu_p^*)}$$

By uniform bound, with probability  $1 - \delta$  for all  $u \in U$  defined in Lemma 9, we have

$$\left| (P(s, a) - \widehat{P}(s, a)) V_p^{\mu_p^*, \widehat{c}_p(\mu_p^*)} \right| \leq \sqrt{2 \log(4|U|/\delta)} \sqrt{\frac{\text{Var}_{s,a}(V_p^{\mu_p^*, \widehat{c}_p(\mu_p^*)})}{N/|\mathcal{S}||\mathcal{A}|}} + \frac{2 \log(4|U|/\delta)}{3(1-\gamma)N/|\mathcal{S}||\mathcal{A}|}.$$

Condition on the event defined in Lemma 9, with probability  $1 - 2\delta$ , we have

$$\left| (P(s, a) - \widehat{P}(s, a)) V_p^{\mu_p^*, \widehat{c}_p(\mu_p^*)} \right| \leq \sqrt{2 \log(4|U|/\delta)} \sqrt{\frac{\text{Var}_{s,a}(V_p^{\mu_p^*, \widehat{c}_p(\mu_p^*)})}{N/|\mathcal{S}||\mathcal{A}}} + \frac{2 \log(4|U|/\delta)}{3(1-\gamma)N/|\mathcal{S}||\mathcal{A}}.$$

With Lemma 3, we have

$$\left| Q_p^{\mu_p^*, \widehat{c}_p(\mu_p^*)} - \widehat{Q}_p^{\mu_p^*, *} \right| \leq \frac{512 \log(4|U|/\delta)}{(1-\gamma)^3 N/|\mathcal{S}||\mathcal{A}}} + \frac{2 \log(4|U|/\delta)}{3(1-\gamma)^2 N/|\mathcal{S}||\mathcal{A}}.$$

If we set  $N = \frac{8096|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3 \epsilon^2} \log\left(\frac{256|\mathcal{S}|^2|\mathcal{A}|}{(1-\gamma)^2 \xi \delta}\right)$ , then with probability larger than  $1 - \delta/2$ , we have

$$\left| Q_p^{\mu_p^*, \widehat{c}_p(\mu_p^*)} - \widehat{Q}_p^{\mu_p^*, *} \right| \leq \frac{\epsilon}{4}.$$

Similarly, with probability larger than  $1 - \delta/2$ , we have

$$\left| \widehat{Q}_p^{\widehat{\mu}_p^*, c_p(\widehat{\mu}_p^*)} - Q_p^{\widehat{\mu}_p^*, *} \right| \leq \frac{\epsilon}{4}.$$

Using Lemma 10, with probability larger than  $1 - \delta$ , we have

$$\left| Q^{\widehat{\mu}^*, *} - Q^* \right| \leq \frac{2\xi}{1-\gamma} + \frac{\epsilon}{2}$$

Setting  $\xi = \frac{\epsilon(1-\gamma)}{2}$ , we have  $|Q^{\widehat{\mu}^*, *} - Q^*| \leq \epsilon$ . Similarly we can prove  $|Q^{*, \widehat{\nu}^*} - Q^*| \leq \epsilon$ . Finally we have that if  $N \geq \frac{8096|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3 \epsilon^2} \log\left(\frac{256|\mathcal{S}|^2|\mathcal{A}|}{(1-\gamma)^3 \epsilon \delta}\right)$ ,  $\widehat{\pi}_p^* = (\widehat{\mu}_p^*, \widehat{\nu}_p^*)$  is an  $\epsilon$ -Nash equilibrium strategy.  $\square$

*Proof of Theorem 3.* Set  $\epsilon = \frac{\Delta}{2}$  and by Theorem 2, we have that if  $N = \frac{C}{(1-\gamma)^3 \Delta^2} |\mathcal{S}||\mathcal{A}| \log\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)\delta\epsilon}\right)$ , then  $\widehat{\pi}^*$  is an  $\frac{\Delta}{2}$ -optimal strategy, which means

$$|Q^* - Q^{\widehat{\pi}^*}| \leq \frac{\Delta}{2}.$$

By the definition of suboptimality gap, we have

$$\begin{aligned} \forall s \in \mathcal{S}_{\max}, a \neq \mu^*(s) : Q^*(s, \mu^*(s)) - Q^*(s, a) &\geq \Delta, \\ \forall s \in \mathcal{S}_{\min}, a \neq \nu^*(s) : Q^*(s, \nu^*(s)) - Q^*(s, a) &\leq -\Delta. \end{aligned}$$

As  $|Q^* - Q^{\widehat{\pi}^*}| \leq \frac{\Delta}{2}$ , we have

$$\begin{aligned} \forall s \in \mathcal{S}_{\max}, a \neq \mu^*(s) : Q^{\widehat{\pi}^*}(s, \mu^*(s)) - Q^{\widehat{\pi}^*}(s, a) &\geq 0, \\ \forall s \in \mathcal{S}_{\min}, a \neq \nu^*(s) : Q^{\widehat{\pi}^*}(s, \nu^*(s)) - Q^{\widehat{\pi}^*}(s, a) &\leq 0, \end{aligned}$$

which means  $\widehat{\pi}^* = \pi^*$ .  $\square$