
Solving Discounted Stochastic Two-Player Games with Near-Optimal Time and Sample Complexity

Aaron Sidford
Stanford University
sidford@stanford.edu

Mengdi Wang
Princeton University
mengdiw@princeton.edu

Lin F. Yang
UCLA
linyang@ee.ucla.edu

Yinyu Ye
Stanford University
yye@stanford.edu

Abstract

In this paper we settle the sampling complexity of solving discounted two-player turn-based zero-sum stochastic games up to poly-logarithmic factors. Given a stochastic game with discount factor $\gamma \in (0, 1)$ we provide an algorithm that computes an ϵ -optimal strategy with high-probability given $\tilde{O}((1 - \gamma)^{-3}\epsilon^{-2})$ samples from the transition function for each state-action-pair. Our algorithm runs in time nearly linear in the number of samples and uses space nearly linear in the number of state-action pairs. As stochastic games generalize Markov decision processes (MDPs) our runtime and sample complexities are optimal due to Azar et al. (2013). We achieve our results by showing how to generalize a near-optimal Q-learning based algorithms for MDP, in particular Sidford et al. (2018a), to two-player strategy computation algorithms. This overcomes limitations of standard Q-learning and strategy iteration or alternating minimization based approaches and we hope will pave the way for future reinforcement learning results by facilitating the extension of MDP results to multi-agent settings with little loss.

1 Introduction

In this paper we study the sample complexity of learning a near-optimal strategy in discounted two-player turn-based zero-sum stochastic games Shapley (1953); Hansen et al. (2013), which we refer to more concisely as *stochastic games*. Stochastic games model dynamic strategic settings in which two players take turns and

the state of game evolves stochastically according to some transition law. This model encapsulates a major challenge in multi-agent learning: other agents may be learning and adapting as well. Further, stochastic games are a generalization of the Markov decision process (MDP), a fundamental model for reinforcement learning, to the two-player setting Littman (1994). MDPs can be viewed as degenerate stochastic games in which one of the players has no influence. Consequently, understanding stochastic games is a natural step towards resolving challenges in reinforcement learning of extending single-agent learning to multi-agent settings.

There is a long line of research in both MDPs and stochastic games (for a more thorough introduction, see Filar and Vrieze (2012); Hansen et al. (2013) and references therein). Strikingly, Hansen et al. (2013) showed that there exists a pure-strategy Nash equilibrium which can be computed in strongly polynomial time for stochastic games, if the game matrix is fully accessible and the discount factor is fixed. In reinforcement learning settings, however, the transition function of the game is unknown and a common goal is to obtain an approximately optimal strategy (a function that maps states to actions) that is able to obtain an expected cumulative reward of at least (or at most) the Nash equilibrium value no matter what the other player does. Unfortunately, despite interest in generalizing MDP results to stochastic games, currently the best known running times/sample complexity for solving stochastic games in a variety of settings are worse than for solving MDPs. This may not be surprising since in general stochastic games are harder to solve than MDPs, e.g., whereas MDPs can be solved in (weakly) polynomial time it remains open whether or not the same can be done for stochastic games.

There are two natural approaches towards achieving sample complexity bounds for solving stochastic games. The first is to note that the popular stochastic value iteration, dynamic programming, and Q-learning methods all apply to stochastic games Littman (1994); Hu and Wellman (2003); Littman (2001a); Perolat et al.

(2015). Consequently, recent advances in these methods Kearns and Singh (1999); Sidford et al. (2018b) developed for MDPs can be directly generalized to solving stochastic games (though the sample complexity of these generalized methods has not been analyzed previously). It is tempting to generalize the analysis of sample optimal methods for estimating values Azar et al. (2013) and estimating policies Sidford et al. (2018a) of MDPs to stochastic games. However, this is challenging as these methods rely on monotonicities in MDPs induced by the linear program nature of the problem Azar et al. (2013); Sidford et al. (2018a).

The second approach would be to apply strategy iteration or alternating minimization / maximization to reduce solving stochastic games to approximately solving a sequence of MDPs. Unfortunately, the best analysis of such a method Hansen et al. (2013) requires solving $\Omega(1/(1-\gamma))$ MDPs. Consequently, even if this approach could be carried out with approximate MDP solvers, the resulting sample complexity for solving stochastic games would be larger than that needed for solving MDPs. More discussion of related literatures is given in Section 1.4.

Given the importance of solving stochastic games in reinforcement learning (e.g. Hu et al. (1998); Bowling and Veloso (2000, 2001); Hu and Wellman (2003); Arslan and Yüksel (2017)), this suggests the following fundamental open problem:

Can we design stochastic game learning algorithms that provably match the performance of MDP algorithms and achieve near-optimal sample complexities?

In this paper, we answer this question in the affirmative in the particular case of solving discounted stochastic games with a generative model, i.e. an oracle for sampling from the transition function for state-action pairs. We provide an algorithm with the same near-optimal sample complexity that is known for solving discounted MDPs. Further, we achieve this result by showing how to transform particular MDP algorithms to solving stochastic games that satisfy particular two-sided monotonicity constraints. Therefore, while there is a major gap between MDPs and stochastic games in terms of computation time for obtaining the exact solutions, this gap disappears when considering the sampling complexity between the two. We hope this work opens the door to more generally extend results for MDP to stochastic games and thereby enable the application of the rich research on reinforcement learning to a broader multi-player settings with little overhead.

1.1 The Model

Formally, throughout this paper, we consider *discounted turn-based two-player zero-*

sum stochastic games described as the tuple $\mathcal{G} = (\mathcal{S}_{\min}, \mathcal{S}_{\max}, \mathcal{A}, \mathbf{P}, \mathbf{r}, \gamma)$. In these games there are two players, a *min* or *minimization* player which seeks to minimize the cumulative reward in the game and a *max* or *maximization* player which seeks to maximize the cumulative reward. Here, \mathcal{S}_{\min} and \mathcal{S}_{\max} are disjoint finite sets of *states* controlled by the min-player and the max-player respectively and their union $\mathcal{S} := \mathcal{S}_{\min} \cup \mathcal{S}_{\max}$ is the set of all possible *states of the game*. Further, \mathcal{A} is a finite set of *actions* available at each state, $\mathbf{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto [0, 1]$ is a *transition probability function*, $\mathbf{r} : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ is the payoff or *reward function* and $\gamma \in (0, 1)$ is a discount factor.¹

Stochastic games $\mathcal{G} = (\mathcal{S}_{\min}, \mathcal{S}_{\max}, \mathcal{A}, \mathbf{P}, \mathbf{r}, \gamma)$ are played dynamically in a sequence of turns, $\{t\}_{t=0}^{\infty}$, starting from some initial state $s^0 \in \mathcal{S}$ at turn $t = 0$. In each turn $t \geq 0$, the game is in one of the states $s^t \in \mathcal{S}$ and the player who controls the state s^t chooses or *plays* an action a^t from the action space \mathcal{A} . This action yields reward $r^t := r(s^t, a^t)$ for the turn and causes the next state s^{t+1} to be chosen at random from \mathcal{S} where the transition probability $\Pr[s^{t+1} = s' | s_1, \dots, s_t, a_1, \dots, a_t] = \mathbf{P}(s' | s^t, a^t)$. The goal of the min-player (resp. max-player) is to choose actions to minimize (resp. maximize) the expected infinite-horizon discounted-reward or *value* of the game $\sum_{t=0}^{\infty} \gamma^t r^t$.

In this paper we focus on the case where the players play pure (deterministic) stationary strategies (policies), i.e. strategies which depend only on the current state. That is we wish to compute a *min-player strategy* or *policy* $\pi_{\min} : \mathcal{S}_{\min} \rightarrow \mathcal{A}$ which defines the action the min player chooses at a state in \mathcal{S}_{\min} and *max-player strategy* $\pi_{\max} : \mathcal{S}_{\max} \rightarrow \mathcal{A}$ which defines the action the max player chooses at a state in \mathcal{S}_{\max} . We call a pair of min-player and max-player strategies $\sigma = (\pi_{\min}, \pi_{\max})$ simply a *strategy*. Further, we let $\sigma(s) := \pi_{\min}(s)$ for $s \in \mathcal{S}_{\min}$ and $\sigma(s) := \pi_{\max}(s)$ for $s \in \mathcal{S}_{\max}$ and define the *value function* or *expected discounted cumulative reward* by \mathbf{v}^{σ} where

$$\mathbf{v}^{\sigma}(s) = \mathbf{v}[\sigma](s) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s^t, \sigma(s^t)) \mid s^0 = s \right]$$

for all $s \in \mathcal{S}$

and the expectation is over the random sequence of states, s^0, s^1, s^2, \dots generated according to \mathbf{P} under the strategy σ , i.e. $\Pr[s^{t+1} = s' | s^t, s^{t-1}, \dots, s^0] = \mathbf{P}(s' | s^t, \sigma(s^t))$ for all $t > 0$.

¹Standard reductions allow this result to be applied for rewards of a broader range Sidford et al. (2018a). Further, while we assume there are the same number of actions per-state, our results easily extend to the case where this is non-uniform; in this case our dependencies on $|\mathcal{S}||\mathcal{A}|$ can be replaced with the number of state-action pairs.

Our goal in solving a game is to compute an approximate *Nash equilibrium* restricted to stationary strategies Nash (1951); Maskin and Tirole (2001). We call a strategy $\sigma = (\pi_{\min}, \pi_{\max})$ an *equilibrium strategy* or *optimal* if

$$\max_{\pi'_{\max}: \mathcal{S}_{\max} \rightarrow \mathcal{A}} \mathbf{v}^{(\pi_{\min}, \pi'_{\max})} \leq \mathbf{v}^{\sigma} \leq \min_{\pi'_{\min}: \mathcal{S}_{\min} \rightarrow \mathcal{A}} \mathbf{v}^{(\pi'_{\min}, \pi_{\max})}.$$

and we call it ϵ -optimal if these same inequalities hold up to an additive ϵ entrywise. It is worth noting that the best response strategy to a stationary policy is also stationary Fudenberg and Tirole (1991) and there always exists a pure stationary strategy attaining the Nash equilibrium Shapley (1953). Consequently, it is sufficient to focus on deterministic strategies.

Throughout this paper we focus on solving stochastic games in the learning setting where the game is not fully specified. We assume that a *generative model* is available which given any state-action pair, i.e. $s \in \mathcal{S}$ and $a \in \mathcal{A}$, can sample a random s' independently at random from the transition probability function, i.e. $\Pr[s' = t] = \mathbf{P}(t | s, a)$. Accessibility to a generative model is a standard and natural assumption (Kakade (2003); Azar et al. (2013); Sidford et al. (2018a); Agarwal et al. (2019)) and corresponds to PAC learning. The special case of solving a MDP given a generative model has been studied extensively (Kakade (2003); Azar et al. (2013); Sidford et al. (2018b,a); Agarwal et al. (2019)) and is a natural proving ground towards designing theoretically motivated reinforcement learning algorithms.

1.2 Our Results

In this paper we provide an algorithm that computes an ϵ -optimal strategy using a sample size that matches the best known sample complexity for solving discounted MDPs. Further, our algorithm runs in time proportional to the number of samples and space proportional to $|\mathcal{S}||\mathcal{A}|$. Interestingly, we achieve this result by showing how to run two-player variant of Q-learning such that the value-strategy sequences induced enjoy certain monotonicity properties. Essentially, we show that provided a value improving algorithm is sufficiently stable, then it can be extended to the two-player setting with limited loss. This allows us to leverage recent advances in solving single player games to solve stochastic games with limited overhead. Our main result is given below.

Theorem 1.1 (Main Theorem). *There is an algorithm which given a stochastic game, $\mathcal{G} = (\mathcal{S}_{\min}, \mathcal{S}_{\max}, \mathbf{P}, \mathbf{r}, \gamma)$ with a generative model, outputs, with probability at least $1 - \delta$, an ϵ -optimal strategy σ by querying $Z = \tilde{O}(|\mathcal{S}||\mathcal{A}|(1 - \gamma)^{-3}\epsilon^{-2})$ samples, where $\epsilon \in (0, 1)$ and $\tilde{O}(\cdot)$ hides polylogarithmic factors. The algorithm runs in time $O(Z)$ and uses space $O(|\mathcal{S}||\mathcal{A}|)$.*

Our sample and time complexities are optimal due to a known lower bound in the single player case by Azar et al. (2013). It was shown in Azar et al. (2013) that solving any one-player MDP to ϵ -optimality with high probability needs at least $\Omega(|\mathcal{S}||\mathcal{A}|(1 - \gamma)^{-3}\epsilon^{-2})$ samples. Our sample complexity upper bound generalizes the recent sharp sample complexity results for solving the discounted MDP Sidford et al. (2018a); Agarwal et al. (2019), and tightly matches the information-theoretic sample complexity up to polylogarithmic factors. This result provides the first and near-optimal sample complexity for solving the two-person stochastic game.

1.3 Notations and Preliminaries

Notation: We use $\mathbf{1}$ to denote the all-ones vector whose dimension is adapted to the context. We use the operators $|\cdot|$, $(\cdot)^2$, $\sqrt{\cdot}$, \leq , \geq as entrywise operators on vectors. We identify the transition probability function \mathbf{P} as a matrix in $\mathbb{R}^{(\mathcal{S} \times \mathcal{A}) \times \mathcal{S}}$ and each row $\mathbf{P}(\cdot | s, a) \in \mathbb{R}^{\mathcal{S}}$ as a vector. We denote \mathbf{v} as a vector in $\mathbb{R}^{\mathcal{S}}$ and \mathbf{Q} as a vector in $\mathbb{R}^{\mathcal{S} \times \mathcal{A}}$. Therefore $\mathbf{P}\mathbf{v}$ is a vector in $\mathbb{R}^{\mathcal{S} \times \mathcal{A}}$. We use σ to denote strategy pairs and π for the min-player or max-player strategy. For any strategy σ , we define $\mathbf{Q}_{\sigma} \in \mathbb{R}^{\mathcal{S}}$ as $\mathbf{Q}_{\sigma}(s) := \mathbf{Q}(s, \sigma(s))$ for $\forall s \in \mathcal{S}$. We denote \mathbf{P}^{σ} as a linear operator defined as

$$\begin{aligned} \forall s \in \mathcal{S} : \quad [\mathbf{P}^{\sigma}\mathbf{v}](s) &= \mathbf{P}(\cdot | s, \sigma(s))^{\top} \mathbf{v}, \\ \forall s, a \in \mathcal{S} \times \mathcal{A} : [\mathbf{P}^{\sigma}\mathbf{Q}](s, a) &= \mathbf{P}(\cdot | s, a)^{\top} \mathbf{Q}_{\sigma}. \end{aligned}$$

Min-value and max-value: For a min-player strategy π_{\min} , we define its *value* as

$$\mathbf{v}^{\pi_{\min}} := \max_{\pi_{\max}: \mathcal{S}_{\max} \rightarrow \mathcal{A}} \mathbf{v}^{(\pi_{\min}, \pi_{\max})}, \quad (1)$$

We let $\sigma_{\max}(\pi_{\min})$ denote a maximizing argument of the above and call it an *optimal counter strategy* of π_{\min} . Thus a value of a min-player strategy gives his expected reward in the *worst case*. We say a min-player strategy π_{\min} is ϵ -optimal if

$$\mathbf{v}^{\pi_{\min}} \leq \min_{\pi'_{\min}: \mathcal{S}_{\min} \rightarrow \mathcal{A}} \mathbf{v}^{\pi'_{\min}} + \epsilon \cdot \mathbf{1}, \quad \text{entrywisely.}$$

The value and ϵ -optimality for the max player is defined similarly. We denote by σ^* the optimal strategy and by \mathbf{v}^* the value function of the optimal strategy.

Q-function: For a strategy σ , we denote its *Q-function* (or *action value*) as $\mathbf{Q}^{\sigma} \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ by $\mathbf{Q}^{\sigma} := \mathbf{r} + \gamma \mathbf{P}\mathbf{v}^{\sigma}$. For a vector $\mathbf{v} \in \mathbb{R}^{\mathcal{S}}$ we denote $\mathbf{Q}(\mathbf{v}) := \mathbf{r} + \gamma \mathbf{P}\mathbf{v}$. Given a $\mathbf{Q} \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$, we denote the greedy value of \mathbf{Q} as

$$\begin{aligned} V[\mathbf{Q}](s) &:= \min_{a \in \mathcal{A}} \mathbf{Q}(s, a) \quad \text{if } s \in \mathcal{S}_{\min} \\ \text{and } V[\mathbf{Q}](s) &:= \max_{a \in \mathcal{A}} \mathbf{Q}(s, a) \quad \text{if } s \in \mathcal{S}_{\max}. \end{aligned}$$

Bellman Operator: We denote the Bellman operator, \mathcal{T} , as follows: $\mathcal{T}[v] \in \mathbb{R}^S$, and

$$\mathcal{T}[v](s) := V[r + \gamma P v].$$

We also denote the greedy strategy, $\sigma(v)$ or $\sigma(Q)$, as the maximization/minimization argument of the \mathcal{T} operator. Moreover, for a given strategy σ , we denote $\mathcal{T}_\sigma[v] = Q(v)_\sigma$. For a given min-player strategy π_{\min} , we define the *half* Bellman operator $\mathcal{H}_{\pi_{\min}}$

$$\mathcal{H}_{\pi_{\min}}[v](s) = r(s, \pi_{\min}(s)) + \gamma P(\cdot | s, \pi_{\min}(s))^\top v$$

$$\text{if } s \in \mathcal{S}_{\min};$$

$$\mathcal{H}_{\pi_{\min}}[v](s) \text{ if } s \in \mathcal{S}_{\max}.$$

We define $\mathcal{H}_{\pi_{\max}}$ similarly. Note that v^* is the unique fixed point of the Bellman operator, i.e., $\mathcal{T}[v^*] = v^*$ (known as the Bellman equation Bellman (1957)). Similarly, $v^{\pi_{\min}}$ (resp. $v^{\pi_{\max}}$) is the unique fixed point for $\mathcal{H}_{\pi_{\min}}$ (resp. $\mathcal{H}_{\pi_{\max}}$). The (half) Bellman-operators satisfy the following properties (see. e.g. Hansen et al. (2013); Puterman (2014))

1. *contraction:* $\|\mathcal{T}[v_1] - \mathcal{T}[v_2]\|_\infty \leq \gamma \|v_1 - v_2\|_\infty$;
2. *monotonicity:* $v_1 \leq v_2 \Rightarrow \mathcal{T}[v_1] \leq \mathcal{T}[v_2]$.

High Probability: we say an algorithm has a property “with high probability” if for any δ by increasing the time and sample complexity by $O(\log(1/\delta))$ it has the property with probability $1 - \delta$.

1.4 Previous Work

Here we provide a more detailed survey of previous works related to stochastic games and MDPs. Two-person stochastic games generalize MDPs Shapley (1953). When one of the players has only one action to choose from, the problem reduces to a MDP. A related game is the stochastic game where both players choose their respective actions simultaneously at each state and the process transitions to the next state under the control of both players Shapley (1953). The turn-based stochastic game can be reduced to the game with simultaneous moves Pérolat et al. (2015).

Computing an optimal strategy for a two-player turn-based zero-sum stochastic game is known to be in $\text{NP} \cap \text{co-NP}$ Condon (1992). Later Hansen et al. (2013) showed that the strategy iteration, a generalization of Howard’s policy iteration algorithm Howard (1960), solves the discounted problem in strongly polynomial time when the discount factor is fixed. Their work uses ideas from Ye (2011) which proved that the policy iteration algorithm solves the discounted MDP (DMDP) in strongly polynomial time when the discount factor is fixed. In general (e.g., if the discount factor is part of the input size), it is open if stochastic games can

be solved in polynomial time Littman (1996). This is in contrast to MDPs which can be solved in (weakly) polynomial time as they are a special case of linear programming.

The algorithms and complexity theory for solving two-player stochastic games is closely related to that of solving MDPs. There is vast literature on solving MDPs which dates back to Bellman who developed value iteration in 1957 Bellman (1957). The policy iteration was introduced shortly after by Howard Howard (1960), and its complexity has been extensively studied in Mansour and Singh (1999); Ye (2011); Scherrer (2013). Then d’Epenoux (1963) and De Ghellinck (1960) discovered that MDPs are special cases of a linear program, which leads to the insight that the simplex method, when applied to solving DMDPs, is a simple policy iteration method. Ye Ye (2011) showed that policy iteration (which is a variant of the general simplex method for linear programming) and the simplex method are strongly polynomial for DMDP and terminate in $O(|\mathcal{S}|^2 |\mathcal{A}| (1-\gamma)^{-1} \log(|\mathcal{S}| (1-\gamma)^{-1}))$ iterations. Hansen et al. (2013) and Scherrer (2013) improved the iteration bound to $O(|\mathcal{S}| |\mathcal{A}| (1-\gamma)^{-1} \log(|\mathcal{S}| (1-\gamma)^{-1}))$ for Howard’s policy iteration method. The best known convergence result for policy and strategy iteration are given by Ye (2005) and Hansen et al. (2013). The best known iteration complexities for both problems are of the order $(1-\gamma)^{-1}$, which becomes unbounded as $\gamma \rightarrow 1$. It is worth mentioning that Ye (2005) designed a combinatorial interior-point algorithm (CIPA) that solves the DMDP in strongly polynomial time.

Sample-based algorithms for learning value and policy functions for MDP have been studied in Kearns and Singh (1999); Kakade (2003); Singh and Yee (1994); Azar et al. (2011b, 2013); Sidford et al. (2018b,a); Agarwal et al. (2019) and many others. Among these papers, Azar et al. (2013) obtains the first tight sample bound for finding an ϵ -optimal value function and for finding ϵ -optimal policies in a restricted ϵ regime and Sidford et al. (2018a) obtains the first tight sample bound for finding an ϵ -optimal *policy* for any ϵ . Both sample complexities are of the form $\tilde{O}(|\mathcal{S}| |\mathcal{A}| (1-\gamma)^{-3})$. Lower bounds have been shown in Azar et al. (2011a); Even-Dar et al. (2006) and Azar et al. (2013). Azar et al. (2013) give the first tight lower bound $\Omega(|\mathcal{S}| |\mathcal{A}| (1-\gamma)^{-3})$. For undiscounted average-reward MDP, a primal-dual based method was proposed in Wang (2017) which achieves sample complexity $\tilde{O}(|\mathcal{S}| |\mathcal{A}| t_{\text{mix}}^2 c_{\text{max}}^2 / c_{\text{min}}^2)$, where t_{mix} is the worst-case mixing time and $c_{\text{max}}/c_{\text{min}}$ is the ergodicity ratio. Sampling-based method for two-player stochastic game has been considered in Wei et al. (2017) in an online learning setting. However, their algorithm leads to a sub-optimal sample-complexity when generalized to the generative model setting.

As for general stochastic games, the minimax Q-learning algorithm and the friend-and-foe Q-learning algorithm were introduced in Littman (1994) and Littman (2001a), respectively. The Nash Q-learning algorithm was proposed for zero-sum games in Hu and Wellman (2003) and for general-sum games in Littman (2001b); Hu and Wellman (1999).

2 Technique Overview

Since stochastic games are a generalization of MDPs, many techniques for solving MDPs can be immediately generalized to stochastic games. However, as we have discussed, some of the techniques used to achieve optimal sample complexities for solving MDPs in a generative model do not have a clear generalization to stochastic games. Nevertheless, we show how to design an algorithm that carefully extends particular Q-learning based methods, i.e. methods that always maintain an estimator for the optimal value function (or Q^*), to achieve our goals.

Q-Learning: To motivate our approach we first briefly review previous Q-learning based methods and the core technique that achieves near-optimal sample complexity. To motivate Q-learning, we first recall the value iteration algorithm solving an MDP. Given a full model for the MDP value iteration updates the iterates as follows

$$\mathbf{v}^{(i)} \leftarrow \mathcal{T}[\mathbf{v}^{(i-1)}] := V[\mathbf{Q}(\mathbf{v}^{(i-1)})]$$

where $\mathbf{v}^{(0)}$ can be an arbitrary vector. Since the Bellman operator is contractive and \mathbf{v}^* is a fix point of \mathcal{T} , this method gives an ϵ -optimal value in $O[(1 - \gamma)^{-1} \log(\epsilon^{-1})]$ iterations. In the learning setting, \mathcal{T} cannot be exactly computed. The Q-learning approach estimates \mathcal{T} by its approximate version, i.e., to compute $\mathbf{P}(\cdot | s, a)^\top \mathbf{v}^{(i-1)}$, we obtain samples from $\mathbf{P}(\cdot | s, a)$, and then compute the empirical average. Then we compute the approximate Q-value at the i -th iteration as

$$\begin{aligned} Q^{(i)} &= \widehat{Q}[\mathbf{v}^{(i-1)}] := \mathbf{r} + \widehat{\mathbf{P}}\mathbf{v}^{(i-1)} \\ \text{and } \widehat{\mathcal{T}}(\mathbf{v}^{(i-1)}) &:= V[\widehat{Q}(\mathbf{v}^{(i-1)})], \end{aligned}$$

where

$$\widehat{\mathbf{P}}(\cdot | s, a)^\top \mathbf{v} = \frac{1}{m} \sum_{s_i \sim P(\cdot | s, a), i \in [m]} \mathbf{v}(s_i)$$

for some $m > 0$. Then the estimation error per step is defined as

$$\epsilon^{(i)} = Q[\mathbf{v}^{(i-1)}] - \widehat{Q}[\mathbf{v}^{(i-1)}].$$

Since the exact value iteration takes at least $\Omega[(1 - \gamma)^{-1}]$ iterations to converge, the Q-learning (or approximated value iteration) takes at least $\Omega[(1 - \gamma)^{-1}]$ iterations. The total number of samples used over all the iterations is the sample complexity of the algorithm.

Variance Control and Monotonicity Techniques:

To obtain the optimal sample complexity for one-player MDP, one approach is to carefully bound each entry of $\epsilon^{(i)}$. By Bernstein inequality (Azar et al. (2013); Sidford et al. (2018a); Agarwal et al. (2019)), we have, with high probability,

$$|\epsilon^{(i)}| \lesssim \sqrt{\text{var}(\mathbf{v}^{(i-1)})/m} \leq \sqrt{\text{var}(\mathbf{v}^*)/m} + \text{lower-order terms.}$$

where $\text{var}(\mathbf{v}) = \mathbf{P}\mathbf{v}^2 - (\mathbf{P}\mathbf{v})^2$ is the *variance-of-value* vector and “ \lesssim ” means “approximately less than.” Let $\pi^{(i)}$ be a policy maintained in the i -th iteration (e.g. the greedy policy of the current Q-value). Due to the estimation error $\epsilon^{(i)}$, the per step error bound reads,

$$Q^* - Q^{(i)} \lesssim \gamma \mathbf{P}^{\pi^*} Q^* - \gamma \mathbf{P}^{\pi^{(i-1)}} Q^{(i-1)} + \epsilon^{(i)}.$$

To derive the overall error accumulation, Sidford et al. (2018a) use the crucial *monotonicity* property, i.e., since $\pi^{(i-1)}(s) = \arg \max_a Q^{(i-1)}(s, a)$, we have

$$Q^{(i-1)}(s, \pi^*(s)) \leq Q^{(i)}(s, \pi^{(i-1)}(s)). \quad (2)$$

We thus have

$$Q^* - Q^{(i)} \lesssim \gamma \mathbf{P}^{\pi^*} Q^* - \gamma \mathbf{P}^{\pi^*} Q^{(i-1)} + \epsilon^{(i)}.$$

By induction, we have

$$Q^* - Q^{(i)} \leq (I - \gamma \mathbf{P}^{\pi^*})^{-1} \sqrt{\text{var}(\mathbf{v}^*)/m} + \text{lower-order terms.} \quad (3)$$

The leading-order error accumulation term $(I - \gamma \mathbf{P}^{\pi^*})^{-1} \sqrt{\text{var}(\mathbf{v}^*)/m}$ satisfies the so-called *total variance property*, and can be upper bounded uniformly by $\sqrt{(1 - \gamma)^{-3} m^{-1}}$, resulting the correct dependence on $(1 - \gamma)$. Therefore the monotonicity property allows us to use π^* as a *proxy* policy, which carefully bounds the error accumulation. For the additional subtlety of how to obtain an optimal policy, please refer to Sidford et al. (2018a) for the variance reduction technique and the monotone-policy technique.

Similar observations regarding MDPs was used in Agarwal et al. (2019) as well. This powerful technique, however, does not generalize to the game case due to the *lack of monotonicity*. Indeed, (2) does not hold for stochastic games due to the existence of both minimization and maximization operations in the Bellman operator. This is the critical issue which this paper seeks to overcome.

Finding Monotone Value-Strategy Sequences for Stochastic Games: Analogously to the MDP case, one approach is to bound error accumulation for stochastic games is to bound each entry of the error

vector $\epsilon^{(i)}$ carefully. In fact, our method for solving stochastic games is very much like the MDP method used in Sidford et al. (2018a). However, the analysis is much different in order to resolve the difficulty introduced by the lack of monotonicity.

Since a stochastic game has two players, we modify the variance reduced Q-value iteration (vQVI) method in Sidford et al. (2018a) to obtain a min-player strategy and a max-player strategy respectively. Since the two players are symmetric, let us focus on introducing and analyzing the algorithm for the min-player. By a slight modification of the vQVI method, we can guarantee to obtain a sequence of strategies and values, $\{\mathbf{v}^{(i)}, \mathbf{Q}^{(i)}, \sigma^{(i)}, \epsilon^{(i)}\}_{i=0}^R$, that satisfy, with high probability,

1. $\mathbf{v}^{(0)} \geq \mathbf{v}^{(1)} \geq \dots \geq \mathbf{v}^{(R)} \geq \mathbf{v}^*$;
2. $\mathcal{T}_{\sigma^{(i)}}[\mathbf{v}^{(i)}] \leq \mathbf{v}^{(i)}, \mathcal{T}[\mathbf{v}^{(i)}] \leq \mathbf{v}^{(i)}, \mathcal{H}_{\pi_{\min}^{(i)}}[\mathbf{v}^{(i)}] \leq \mathbf{v}^{(i)}$;
3. $\mathbf{Q}^{(i)} \leq \mathbf{Q}[\mathbf{v}^{(i-1)}] + \epsilon^{(i)}$;
4. $\mathbf{v}^{(i)} \leq V[\mathbf{Q}^{(i)}]$.

where $\sigma^{(i)} = (\pi_{\max}^{(i)}, \pi_{\min}^{(i)})$. The first property guarantees that the value sequences are monotonically decreasing, the second property guarantees $\mathbf{v}^{(i)}$ is always an upper bound of the value $\mathbf{v}^{\pi_{\min}^{(i)}}$, and the third and fourth inequality guarantees that $\mathbf{v}^{(i)}$ is well approximated by $V[\mathbf{Q}^{(i)}]$ and the estimation error satisfy $|\epsilon^{(i)}| \lesssim \sqrt{\text{var}(\mathbf{v}^{(i)})/m}$, where m is the total number of samples used per state-action pair. Note that, as long as we can guarantee that $\mathbf{v}^{(R)} - \mathbf{v}^* \leq \epsilon$, we can guarantee the min-strategy $\pi_{\min}^{(R)}$ is also good: $\mathbf{v}^* \leq \mathbf{v}^{\pi_{\min}^{(R)}} \leq \mathbf{v}^{(R)}$.

Controlling Error Accumulation using Auxiliary Markovian Strategy: Due to the lack of monotonicity (2), we cannot use the optimal strategy σ^* as a proxy strategy to carefully account for the error accumulation. To resolve this issue, we construct a new proxy strategy σ^∞ . This strategy is a Markovian strategy, which is time-dependent but not history dependent, i.e., at time t , the strategy played is a deterministic map $\sigma_t^\infty : \mathcal{S} \rightarrow \mathcal{A}$. The proxy strategy satisfies the following:

Underestimation. its value, $\mathbf{v}[\sigma_i^\infty]$, (expected discounted cumulative reward starting from any time) is upper bounded by \mathbf{v}^* ;

Contraction.

$$\mathbf{v}^{(i)}(s) - \mathbf{v}[\sigma_i^\infty](s) \leq \gamma \mathbf{P}(\cdot | s, \sigma_i^\infty(s))^\top (\mathbf{v}^{(i-1)} - \mathbf{v}[\sigma_{i-1}^\infty]) + \epsilon^{(i)}(s, \sigma_i^\infty(s)),$$

Similarly, we can bound the error $\epsilon^{(i)}(s, \sigma_i^\infty(s))$ by the variance-of-value of the proxy strategy

$$\epsilon^{(i)}(s, \sigma_i^\infty(s)) \leq \sqrt{\text{var}(\mathbf{v}[\sigma_i^\infty])(s, \sigma_i^\infty(s))/m} + \text{lower-order terms.}$$

Based on the first property, we can upper bound

$$\mathbf{v}^{(i)} - \mathbf{v}^* \leq \mathbf{v}^{(i)} - \mathbf{v}[\sigma_i^\infty].$$

Based on the second property, and induction on i , we can now write a new form of error accumulation,

$$\begin{aligned} \mathbf{v}^{(R)} - \mathbf{v}^* &\lesssim \sum_{i=1}^R \gamma^{R-i} \mathbf{P}^{\sigma_R^\infty} \cdot \mathbf{P}^{\sigma_{R-1}^\infty} \cdot \dots \cdot \mathbf{P}^{\sigma_{i+1}^\infty} \\ &\quad \cdot \sqrt{\text{var}(\mathbf{v}[\sigma_{i-1}^\infty])_{\sigma_i^\infty}/m} + \text{lower-order terms,} \end{aligned}$$

where $\text{var}(\mathbf{v}[\sigma_{i-1}^\infty])_{\sigma_i^\infty}(s) := \text{var}(\mathbf{v}[\sigma_i^\infty])(s, \sigma_i^\infty(s))$ for all $s \in \mathcal{S}$. We derive a new *law of total variance* bound for the first term and ultimately prove an error accumulation upper bound:

$$\mathbf{v}^{(R)} - \mathbf{v}^* \lesssim \sqrt{(1-\gamma)^{-3}m} + \text{lower-order terms,}$$

giving the optimal sample bound.

3 Sample Complexity of Stochastic Games

In this section, we provide and analyze our sampling-based algorithm for solving stochastic games. Recall that we have a *generative model* for the game such that we can obtain samples from state-action pairs. Each sample is obtained in time $O(1)$. As such we care about the total number of samples used or the total amount of time consumed by the algorithm. We will provide an efficient algorithm that takes input a generative model and obtains a good strategy for the underlying stochastic game.

We now describe the algorithm. Since the min-player and max-player are symmetric, let us focus on the min-player strategy. For the max player strategy, we can either consider the game $\mathcal{G}' = (\mathcal{S}_{\min}, \mathcal{S}_{\max}, \mathbf{P}, \mathbf{1} - \mathbf{r}, \gamma)$, in which the roles of the max and min players switched, or use the corresponding algorithm for the max-player defined in Section A.4, an algorithm that is a direct generalization from the min-player algorithm.

The Full Algorithm. For simplicity, let us denote $\beta = 1/(1-\gamma)$. Our full algorithm will use the QVI-MDVSS algorithm (Algorithm 1) as a subroutine. As we will show shortly, this subroutine maintains a monotonic value strategy sequence with high probability. Suppose the algorithm is specified by an accuracy parameter $\epsilon \in (0, 1]$. We initialize a value vector $\mathbf{v}^{(0)} = \beta \mathbf{1}$, and an arbitrary strategy $\sigma^{(0)} = (\pi_{\min}^{(0)}, \pi_{\max}^{(0)})$. Let $u^{(0)} = \beta$. Then our initial value and strategy satisfy the requirement of the input specified by Algorithm 1:

$$\begin{aligned} \mathbf{v}^* &\leq \mathbf{v}^{(0)} \leq \mathbf{v}^* + u^{(0)} \mathbf{1}, & \mathbf{v}^{(0)} &\geq \mathcal{T}[\mathbf{v}^{(0)}], \\ \text{and } \mathbf{v}^{(0)} &\geq \mathcal{T}_{\sigma^{(0)}}[\mathbf{v}^{(0)}]; \end{aligned}$$

Algorithm 1 QVI-MDVSS: algorithm for computing monotone decreasing value-strategy sequences.

1: **Input:** A generative model for stochastic game, $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathbf{r}, \mathbf{P}, \gamma)$;
 2: **Input:** Precision parameter $u \in [0, (1 - \gamma)^{-1}]$, and error probability $\delta \in (0, 1)$;
 3: **Input:** Initial values $\mathbf{v}^{+(0)}, \sigma^{+(0)}$ that satisfies monotonicity:

$$\mathbf{v}^* \leq \mathbf{v}^{+(0)} \leq \mathbf{v}^* + u \mathbf{1}, \quad \mathbf{v}^{+(0)} \geq \mathcal{T}[\mathbf{v}^{+(0)}], \quad \text{and} \quad \mathbf{v}^{+(0)} \geq \mathcal{T}_{\sigma^{+(0)}}[\mathbf{v}^{+(0)}]; \quad (5)$$

4: **Output:** $\{\mathbf{v}^{+(i)}, \mathbf{Q}^{+(i)}, \sigma^{+(i)}, \boldsymbol{\xi}^{+(i)}\}_{i=0}^R$ which is an MDVSS with probability at least $1 - \delta$;
 5:
 6: **INITIALIZATION:**
 7: Let c_1, c_2, c_3, c be some tunable absolute constants;
 8: *\\Initialize constants:*
 9: $\beta \leftarrow (1 - \gamma)^{-1}$, and $R \leftarrow \lceil c_1 \beta \ln[\beta u^{-1}] \rceil$; $m_1 \leftarrow c_2 \beta^3 \cdot \min(1, u^{-2}) \cdot \log(8|\mathcal{S}||\mathcal{A}|\delta^{-1})$;
 10: $m_2 \leftarrow c_3 \beta^2 \log[2R|\mathcal{S}||\mathcal{A}|\delta^{-1}]$; $\alpha_1 \leftarrow L/m_1$ where $L = c \log(|\mathcal{S}||\mathcal{A}|\delta^{-1}(1 - \gamma)^{-1}u^{-1})$;
 11: *\\Obtain an initial batch of samples:*
 12: For each $(s, a) \in \mathcal{S} \times \mathcal{A}$: obtain independent samples $s_{s,a}^{(1)}, s_{s,a}^{(2)}, \dots, s_{s,a}^{(m_1)}$ from $\mathbf{P}(\cdot|s, a)$;
 13: Initialize: $\mathbf{w}^+ = \tilde{\mathbf{w}}^+ = \hat{\sigma}^+ = \mathbf{Q}^{+(0)} = \mathbf{Q}^{+(1)} \leftarrow \beta \cdot \mathbf{1}_{\mathcal{S} \times \mathcal{A}}$ and $i \leftarrow 0$;
 14: **for** each $(s, a) \in \mathcal{S} \times \mathcal{A}$ **do**
 15: *\\Compute empirical estimates of $\mathbf{P}_{s,a}^\top \mathbf{v}^{+(0)}$ and $\text{var}(\mathbf{v}^{+(0)})(s, a)$:*
 16: $\tilde{\mathbf{w}}^+(s, a) \leftarrow \frac{1}{m_1} \sum_{j=1}^{m_1} \mathbf{v}^{+(0)}(s_{s,a}^{(j)})$; $\hat{\sigma}^+(s, a) \leftarrow \frac{1}{m_1} \sum_{j=1}^{m_1} (\mathbf{v}^{+(0)}(s_{s,a}^{(j)}))^2 - (\tilde{\mathbf{w}}^+(s, a))^2$;
 17: *\\Shift the empirical estimate to have one-sided error and guarantee monotonicity:*
 18: $\mathbf{w}^+(s, a) \leftarrow \tilde{\mathbf{w}}^+(s, a) + \sqrt{\alpha_1 \hat{\sigma}^+(s, a)} + \alpha_1^{3/4} \beta$
 19: *\\Compute coarse estimate of the Q-function and make sure its value is in $[0, \beta]$:*
 20: $\mathbf{Q}^{+(0)}(s, a) \leftarrow \min[\mathbf{r}(s, a) + \gamma \mathbf{w}^+(s, a), \beta]$
 21: **end for**
 22:
 23: **REPEAT:** *\\successively improve*
 24: **for** $i = 1$ to R **do**
 25: *\\Compute the one-step dynamic programming:*
 26: Let $\mathbf{v}^{+(i)} \leftarrow \tilde{\mathbf{v}}^{+(i)} \leftarrow \mathcal{T}[\mathbf{Q}^{+(i-1)}]$, $\sigma^{+(i)} \leftarrow \tilde{\sigma}^{+(i)} \leftarrow \sigma(\mathbf{Q}^{+(i-1)})$;
 27: *\\Compute strategy and value and maintain monotonicity:*
 28: For each $s \in \mathcal{S}$ if $\mathbf{v}^{+(i)}(s) \geq \mathbf{v}^{+(i-1)}(s)$, then $\mathbf{v}^{+(i)}(s) \leftarrow \mathbf{v}^{+(i-1)}(s)$ and $\sigma^{+(i)}(s) \leftarrow \sigma^{+(i-1)}(s)$;
 29: *\\Obtaining a small batch of samples:*
 30: For each $(s, a) \in \mathcal{S} \times \mathcal{A}$: draw independent samples $\tilde{s}_{s,a}^{(1)}, \tilde{s}_{s,a}^{(2)}, \dots, \tilde{s}_{s,a}^{(m_2)}$ from $\mathbf{P}(\cdot|s, a)$;
 31: *\\Compute the expected value, $\mathbf{g}^{\pm(i)}$, the estimate of $\mathbf{P}[\mathbf{v}^{\pm(i)} - \mathbf{v}^{\pm(0)}]$ with one-sided error:*
 32: Let $\tilde{\mathbf{g}}^{+(i)}(s, a) \leftarrow \frac{1}{m_2} \sum_{j=1}^{m_2} [\mathbf{v}^{+(i)}(\tilde{s}_{s,a}^{(j)}) - \mathbf{v}^{+(0)}(\tilde{s}_{s,a}^{(j)})]$;
 33: Let $\mathbf{g}^{+(i)}(s, a) \leftarrow \tilde{\mathbf{g}}^{+(i)}(s, a) + C(1 - \gamma)u$, where $C > 0$ is an absolute constant;
 34: *\\Estimate the approximation error:*
 35: $\boldsymbol{\xi}^{+(i)} \leftarrow 2\sqrt{\alpha_1 \sigma_{\mathbf{v}^{+(0)}}} + 2[\alpha_1^{3/4} \beta + C(1 - \gamma)u] \cdot \mathbf{1}$
 36: *\\Improve $\mathbf{Q}^{+(i)}$ and make sure its value is in $[0, \beta]$:*
 37: $\mathbf{Q}^{+(i+1)} \leftarrow \min[\mathbf{r} + \gamma \cdot [\mathbf{w}^+ + \mathbf{g}^{+(i)}], \beta]$;
 38: **end for**
 39: **return** $\{\mathbf{v}^{+(i)}, \mathbf{Q}^{+(i)}, \sigma^{+(i)}, \boldsymbol{\xi}^{+(i)}\}_{i=0}^R$

Let $u^{(j)} \leftarrow \beta/2^j$ and $\delta \leftarrow 1/\text{poly}(\log(\beta/\epsilon))$.
 We run Algorithm 1 repeatedly:

$$\begin{aligned} (v^{(j+1)}, \sigma^{(j+1)}) &\leftarrow \text{QVI-MDVSS} \\ &\leftarrow (v^{(j)}, \sigma^{(j)}, u^{(j)}, \delta), \quad (6) \end{aligned}$$

where $\sigma^{(j)} = (\pi_{\min}^{(j)}, \pi_{\max}^{(j)})$ and we take the terminal value and strategy of the output sequence of Algorithm 1 as the input for the next iteration. In total we run (6) $R' = \Theta(\log(\beta/\epsilon))$ iterations. In the end, we output $\pi_{\min}^{(R')}$ from $\sigma^{(R')} = (\pi_{\min}^{(R')}, \pi_{\max}^{(R')})$ as our min-player strategy.

The formal guarantee of the algorithm is presented in the following theorem.

Theorem 3.1 (Restatement of Theorem 1.1). *Given a stochastic game $\mathcal{G} = (\mathcal{S}_{\min}, \mathcal{S}_{\max}, \mathbf{P}, \mathbf{r}, \gamma)$ with a generative model, there exists (constructively) an algorithm that outputs, with probability at least $1 - \delta$, an ϵ -optimal strategy σ by querying $Z := \tilde{O}(|\mathcal{S}||\mathcal{A}|(1 - \gamma)^{-3}\epsilon^{-2})$ samples in time $O(Z)$ using space $O(|\mathcal{S}||\mathcal{A}|)$ where $\epsilon \in (0, 1)$ and $\tilde{O}(\cdot)$ hides $\text{poly} \log[|\mathcal{S}||\mathcal{A}|/(1 - \gamma)/\epsilon/\delta]$ factors.*

The formal proof of Theorem 3.1 is given in the next section. Here we give a sketch of the proof.

Proof Sketch of Theorem 3.1: We first show the

high-level idea. Considering one iteration of (6), we claim that if the input value and strategy $\sigma^{(j)}, \mathbf{v}^{(j)}, u^{(j)}$ satisfies the input condition (5), then with probability at least $1 - \delta$, the terminal value and strategy of the output sequence, $\sigma^{(j+1)}, \mathbf{v}^{(j+1)}$, satisfies,

$$\mathbf{v}^{\pi_{\min}^{j+1}} \leq \mathbf{v}^{j+1} \leq \mathbf{v}^* + u^{(j)} \mathbf{1} / 2 =: \mathbf{v}^* + u^{(j+1)} \mathbf{1}; \quad (7)$$

and $(\sigma^{(j+1)}, \mathbf{v}^{(j+1)}, u^{(j+1)})$ satisfies the the input condition (5). Namely, with high probability, the error of the output is decreased by at least half and the output can be used as an input to the QVI-MDVSS algorithm again. Suppose we run the subroutine of Algorithm 1 for R' times, and conditioning on the event that all the instances of QVI-MDVSS succeed, the final error of $\pi_{\min}^{(R')}$ is then at most $u^{(R')} = 2^{-R'} \beta = \epsilon$, as desired. By setting $\delta = \delta' / R'$ for some $\delta' > 0$, we have that all QVI-MDVSS instances succeed with probability at least $1 - \delta'$. It remains to show that the algorithm QVI-MDVSS works as claimed.

High-level Structure of Algorithm 1. To outline the proof, we denote a *monotone decreasing value-strategy sequence* (MDVSS) as $\{\mathbf{v}^{(i)}, \mathbf{Q}^{(i)}, \sigma^{(i)}, \boldsymbol{\epsilon}^{(i)}\}_{i=0}^R$, satisfying (4), where $\mathbf{v}^{(i)}, \boldsymbol{\epsilon}^{(i)} \in \mathbb{R}^{\mathcal{S}}, \mathbf{Q}^{(i)} \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ and $\sigma^{(i)} = (\pi_{\min}^{(i)}, \pi_{\max}^{(i)}) \in \mathcal{A}^{\mathcal{S}}$. A more formal treatment of the sequence is presented in Section A.2.

We next introduce the high-level idea of Algorithm 1. The basic step of the algorithm is to do approximate value-iteration while preserving all monotonic properties required by an MDVSS, i.e., we would like to approximate

$$\begin{aligned} \mathbf{Q}^{(i)} &= \mathbf{Q}[\mathbf{v}^{(i-1)}] := \mathbf{r} + \mathbf{P}\mathbf{v}^{(i-1)} \\ \text{and } \mathcal{T}[\mathbf{v}^{(i-1)}] &:= V[\mathbf{Q}(\mathbf{v}^{(i-1)})]. \end{aligned}$$

We would like to approximate $\mathbf{P}\mathbf{v}^{(i-1)}$ using samples, but we do not want to use the same amount of samples per iteration (as it become costly if the number of iterations is large). Instead, we compute only the *first* iteration (i.e., estimate $\mathbf{P}\mathbf{v}^{(0)}$) up to high accuracy with a large number of samples (m_1 samples, defined in Line 9). These computations are presented in Line 15-20. To maintain an upper bound of the of the estimation error, we also compute the empirical variances of the updates in Line 16. We shift upwards our estimates by the estimation error upper bounds to make our estimators one-sided, which is crucial to maintain the MDVSS properties. For the subsequent steps (Line 26 - 37), we use m_2 samples per iteration ($m_2 \ll m_1$) to estimate $\mathbf{P}(\mathbf{v}^{(i)} - \mathbf{v}^{(0)})$. The expectation is that $(\mathbf{v}^{(i)} - \mathbf{v}^{(0)})$ has a small ℓ_∞ norm, and hence $\mathbf{P}(\mathbf{v}^{(i)} - \mathbf{v}^{(0)})$ can be estimated up to high accuracy with only a small number of samples. The estimator of $\mathbf{P}(\mathbf{v}^{(i)} - \mathbf{v}^{(0)})$ plus the estimator of $\mathbf{P}\mathbf{v}^{(0)}$ in the initialization steps gives a high-accuracy estimator (Line 37) for the value

iteration. Since $m_2 \ll m_1$, the total number of samples per state-action pair is dominated by m_1 . This idea is formally known as *variance-reduction*, firstly proposed for solving MDP in Sidford et al. (2018b). Similarly, we shift our estimators to be one-sided. We additionally maintain carefully-designed strategies in Line 26-28 to preserve monotonicity. Hence the algorithm can be viewed as a value-strategy iteration algorithm.

Correctness of Algorithm 1. We now sketch the proof of correctness for Algorithm 1. Firstly Proposition (A.3) shows that the if an MDVSS, e.g., $\{\mathbf{v}^{+(i)}, \mathbf{Q}^{+(i)}, \sigma^{+(i)}, \boldsymbol{\epsilon}^{+(i)}\}_{i=0}^R$, satisfies $\|\mathbf{v}^{+(R)} - \mathbf{v}^*\|_\infty \leq \epsilon$ for some $\epsilon > 0$ then their terminal strategies and values satisfy

$$\mathbf{v}^{\pi_{\min}^{+(R)}} \leq \mathbf{v}^{+(R)} \leq \mathbf{v}^* + \epsilon \mathbf{1}.$$

This indicates that as long as we can show $\epsilon \leq u/2$, then the *halving-error-property* (7) holds.

Proposition A.4 shows the halving-error-property can be achieved by setting

$$\epsilon^{+(i)} \lesssim \sqrt{\text{var}(\mathbf{v}^{+(0)})/m} + \text{lower-order terms},$$

where $\text{var}(\mathbf{v}^{+(0)})$ is the variance-of-value vector of $\mathbf{v}^{+(0)}$ and $m \gtrsim \sqrt{\beta^3 u^{-2}}$. This proof is based on constructing an auxiliary Markovian strategy for analyzing the error accumulation throughout the value-strategy iterations. The Markovian strategy is a time-dependent strategy used as a proxy for analyzing the entrywise error recursion (Lemmas A.4-A.11).

Proposition A.12 shows, with high probability, Algorithm 1 produces value-strategy sequences $\{\mathbf{v}^{+(i)}, \mathbf{Q}^{+(i)}, \sigma^{+(i)}, \boldsymbol{\xi}^{+(i)}\}_{i=0}^R$, which is indeed an MDVSS and $\boldsymbol{\xi}^{+(i)}$ satisfies Proposition A.4. The proof involves analyzing the probability of “good events” on which monotonicity is preserved at every iteration by using confidence estimates computed during the iterations and concentration arguments. See Lemmas A.13-A.18 for the full proof of Proposition A.12.

Putting Everything Together. Finally by putting together the strategies, we conclude that the terminal strategy of the iteration (6) is always an approximately optimal min-player strategy to the game, with high probability. For implementation, since our algorithm only computes the inner product based on samples, the total computation time is proportional to the number of samples. Moreover, since we can update as samples are drawn and output the monotone sequences as they are generated, we do not need to store samples or the value-strategy sequences, thus the overall space is $O(|\mathcal{S}||\mathcal{A}|)$. \square

References

- Agarwal, A., Kakade, S., and Yang, L. F. (2019). On the optimality of sparse model-based planning for markov decision processes. *arXiv preprint arXiv:1906.03804*.
- Arslan, G. and Yüksel, S. (2017). Decentralized q-learning for stochastic teams and games. *IEEE Transactions on Automatic Control*, 62(4):1545–1558.
- Azar, M. G., Munos, R., Ghavamzadeh, M., and Kappen, H. (2011a). Reinforcement learning with a near optimal rate of convergence.
- Azar, M. G., Munos, R., Ghavamzadeh, M., and Kappen, H. (2011b). Speedy q-learning. In *Advances in neural information processing systems*.
- Azar, M. G., Munos, R., and Kappen, H. J. (2013). Minimax pac bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, 91(3):325–349.
- Bellman, R. (1957). *Dynamic Programming*. Princeton University Press, Princeton, NJ.
- Bowling, M. and Veloso, M. (2000). An analysis of stochastic game theory for multiagent reinforcement learning. Technical report, Carnegie-Mellon Univ Pittsburgh Pa School of Computer Science.
- Bowling, M. and Veloso, M. (2001). Rational and convergent learning in stochastic games. In *International joint conference on artificial intelligence*, volume 17, pages 1021–1026. Lawrence Erlbaum Associates Ltd.
- Condon, A. (1992). The complexity of stochastic games. *Information and Computation*, 96(2):203–224.
- De Ghellinck, G. (1960). Les problemes de decisions sequentielles. *Cahiers du Centre dEtudes de Recherche Opérationnelle*, 2(2):161–179.
- d’Epenoux, F. (1963). A probabilistic production and inventory problem. *Management Science*, 10(1):98–108.
- Even-Dar, E., Mannor, S., and Mansour, Y. (2006). Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *Journal of machine learning research*, 7(Jun):1079–1105.
- Filar, J. and Vrieze, K. (2012). *Competitive Markov decision processes*. Springer Science & Business Media.
- Fudenberg, D. and Tirole, J. (1991). *Game theory*. MIT Press, Cambridge, MA.
- Hansen, T. D., Miltersen, P. B., and Zwick, U. (2013). 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.
- Howard, R. A. (1960). *Dynamic programming and Markov processes*. The MIT press, Cambridge, MA.
- Hu, J. and Wellman, M. P. (1999). Multiagent reinforcement learning in stochastic games. *Submitted for publication*.
- Hu, J. and Wellman, M. P. (2003). Nash q-learning for general-sum stochastic games. *Journal of machine learning research*, 4(Nov):1039–1069.
- Hu, J., Wellman, M. P., et al. (1998). Multiagent reinforcement learning: theoretical framework and an algorithm. In *ICML*, volume 98, pages 242–250. Citeseer.
- Kakade, S. M. (2003). *On the sample complexity of reinforcement learning*. PhD thesis, University of London London, England.
- Kearns, M. J. and Singh, S. P. (1999). Finite-sample convergence rates for q-learning and indirect algorithms. In *Advances in neural information processing systems*, pages 996–1002.
- Littman, M. L. (1994). Markov games as a framework for multi-agent reinforcement learning. In *Proceedings of the 11th International Conference on International Conference on Machine Learning, ICML’94*, pages 157–163, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Littman, M. L. (1996). Algorithms for sequential decision making.
- Littman, M. L. (2001a). Friend-or-foe q-learning in general-sum games. In *ICML*, volume 1, pages 322–328.
- Littman, M. L. (2001b). Value-function reinforcement learning in markov games. *Cognitive Systems Research*, 2(1):55–66.
- Mansour, Y. and Singh, S. (1999). On the complexity of policy iteration. In *Proceedings of the Fifteenth conference on Uncertainty in artificial intelligence*, pages 401–408. Morgan Kaufmann Publishers Inc.
- Maskin, E. and Tirole, J. (2001). Markov perfect equilibrium. I. Observable actions. *J. Econom. Theory*, 100(2):191–219.
- Nash, J. (1951). Non-cooperative games. *Annals of mathematics*, pages 286–295.
- Perolat, J., Scherrer, B., Piot, B., and Pietquin, O. (2015). Approximate dynamic programming for two-player zero-sum markov games. In *International Conference on Machine Learning (ICML 2015)*.
- Pérolat, J., Scherrer, B., Piot, B., and Pietquin, O. (2015). Approximate dynamic programming for two-player zero-sum markov games. In *Proceedings of the 32th International Conference on International Conference on Machine Learning, ICML’15*.

- Puterman, M. L. (2014). *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons.
- Scherrer, B. (2013). Improved and generalized upper bounds on the complexity of policy iteration. In *Advances in Neural Information Processing Systems*, pages 386–394.
- Shapley, L. S. (1953). Stochastic games. *Proceedings of the national academy of sciences*, 39(10):1095–1100.
- Sidford, A., Wang, M., Wu, X., Yang, L., and Ye, Y. (2018a). Near-optimal time and sample complexities for solving markov decision processes with a generative model. In *Advances in Neural Information Processing Systems*, pages 5186–5196.
- Sidford, A., Wang, M., Wu, X., and Ye, Y. (2018b). 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. Society for Industrial and Applied Mathematics.
- Singh, S. P. and Yee, R. C. (1994). An upper bound on the loss from approximate optimal-value functions. *Machine Learning*, 16(3):227–233.
- Wang, M. (2017). Randomized linear programming solves the discounted Markov decision problem in nearly-linear running time. *arXiv preprint arXiv:1704.01869*.
- Wei, C.-Y., Hong, Y.-T., and Lu, C.-J. (2017). Online reinforcement learning in stochastic games. In *Advances in Neural Information Processing Systems*, pages 4994–5004.
- Ye, Y. (2005). A new complexity result on solving the Markov decision problem. *Mathematics of Operations Research*, 30(3):733–749.
- Ye, Y. (2011). The simplex and policy-iteration methods are strongly polynomial for the Markov decision problem with a fixed discount rate. *Mathematics of Operations Research*, 36(4):593–603.

A Proof of Main Results

The remainder of this section is devoted to proving Theorem 1.1. We prove this by formally providing a notion of *monotone value-strategy sequences*. With this, we show if an algorithm outputs some monotone value-strategy sequence, then the terminal strategy of the sequence is always an approximately optimal strategy to the game. We then show that Algorithm 1 produces monotone value-strategy sequences with high probability.

A.1 Additional Notation

First we provide additional notation critical to our proofs.

Markovian Strategies: We denote a Markovian strategy σ^∞ as an infinitely long sequence of pre-defined strategies

$$\sigma^\infty := (\sigma_1, \sigma_2, \dots),$$

where each σ_i is a normal deterministic strategy. We denote

$$\sigma_t^\infty = (\sigma_t, \sigma_{t+1}, \dots)$$

as another Markovian strategy. We denote σ_{\min}^∞ and σ_{\max}^∞ as the min-player strategy and the max-player strategy respectively. When using the strategy, players uses σ_t at time t . The strategy is Markovian because it does not depend on the historical moves. Note that a stationary strategy σ is a special case of the Markovian strategy: $\sigma = (\sigma, \sigma, \dots)$. The value of a Markovian strategy is defined as before, but the states are generated by playing the action $\sigma_t(s^t)$ at time t . Since the strategy has a time dependence, we denote

$$\mathbf{v}_t^{\sigma^\infty} := \mathbf{v}[\sigma_t^\infty] \quad \text{and} \quad \mathbf{Q}_t^{\sigma^\infty} = \mathbf{r} + \gamma \mathbf{P} \mathbf{v}_{t+1}^{\sigma^\infty}.$$

The (half) Bellman operators are defined similarly to that of stationary policies.

A.2 Monotone Value-Strategy Sequence

In this section we formally define monotone strategy value sequences. Such a sequence, although not explicitly stated in Sidford et al. (2018b,a), are crucial for these algorithms to obtain good policy while obtaining a good value for an MDP. In the following sections, we denote $m \geq 1$, $L \geq 1$ and $\epsilon \in [0, (1 - \gamma)^{-1}]$ as parameters. Monotone value-strategy sequences are formally defined as follows.

Definition A.1 (Monotone Decreasing Value-Strategy Sequence). *A monotone decreasing value-strategy sequence (MDVSS) is a sequence of $\{\mathbf{v}^{(i)}, \mathbf{Q}^{(i)}, \sigma^{(i)}, \boldsymbol{\epsilon}^{(i)}\}_{i=0}^R$ where $\mathbf{v}^{(i)}, \boldsymbol{\epsilon}^{(i)} \in \mathbb{R}^{\mathcal{S}}$, $\mathbf{Q}^{(i)} \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ and $\sigma^{(i)} = (\pi_{\min}^{(i)}, \pi_{\max}^{(i)}) \in \mathcal{A}^{\mathcal{S}}$ satisfy*

1. $\mathbf{v}^{(0)} \geq \mathbf{v}^{(1)} \geq \dots \geq \mathbf{v}^{(R)} \geq \mathbf{v}^*$;
2. $\forall i \in [0, R], \mathcal{T}_{\sigma^{(i)}}[\mathbf{v}^{(i)}] \leq \mathbf{v}^{(i)}, \mathcal{T}[\mathbf{v}^{(i)}] \leq \mathbf{v}^{(i)}, \mathcal{H}_{\pi_{\min}^{(i)}}[\mathbf{v}^{(i)}] \leq \mathbf{v}^{(i)}$;
3. $\forall i \in [R], \mathbf{Q}^{(i)} \leq \mathbf{r} + \gamma \mathbf{P} \mathbf{v}^{(i-1)} + \boldsymbol{\epsilon}^{(i)}$;
4. $\forall i \in [R], \mathbf{v}^{(i)} \leq V[\mathbf{Q}^{(i)}]$.

Note that $\mathbf{Q}^{(0)}, \boldsymbol{\epsilon}^{(0)}$ can be arbitrary.

Here, we explain the intuition of the sequence. The first property guarantees that the value-estimator $\mathbf{v}^{(i)}$ s always upper bound the optimal value. The second property guarantees that $\mathbf{v}^{\pi_{\min}} \leq \mathbf{v}^{(i)}$. Indeed

$$\mathbf{v}^{\pi_{\min}} = \lim_{t \rightarrow \infty} \mathcal{H}_{\pi_{\min}}^t[\mathbf{v}^{(i)}] \leq \mathbf{v}^{(i)},$$

where $\mathcal{H}_{\pi_{\min}}^t$ denotes applying $\mathcal{H}_{\pi_{\min}}$ for t times. Therefore, as long as $\mathbf{v}^{(R)} - \mathbf{v}^* \leq \epsilon \mathbf{1}$, we have

$$\mathbf{v}^* \leq \mathbf{v}^{\pi_{\min}} \leq \mathbf{v}^* + \epsilon \mathbf{1}.$$

The third and the fourth property guarantees $\mathbf{v}^{(R)}$ is good by requiring that $\mathbf{v}^{(i)}$ and $\mathbf{Q}^{(i)}$ satisfy the approximate value iteration with one-sided error. However the overall error $\mathbf{v}^{(R)} - \mathbf{v}^*$ is controlled by the per-step error term $\boldsymbol{\epsilon}^{(i)}$.

Similarly, we define *monotone increasing value-strategy sequence*(MIVSS) analogously with every inequality reversed.

Definition A.2 (Monotone Increasing Value-Strategy Sequence). *A monotone increasing value-strategy sequence (MIVSS) is a sequence of $\{\mathbf{v}^{(i)}, \mathbf{Q}^{(i)}, \sigma^{(i)}, \boldsymbol{\epsilon}^{(i)}\}_{i=0}^R$ where $\mathbf{v}^{(i)} \in \mathbb{R}^S$, $\mathbf{Q}^{(i)} \in \mathbb{R}^{S \times \mathcal{A}}$ and $\sigma^{(i)} \in \mathcal{A}^S$ that satisfies,*

1. $\mathbf{v}^{(0)} \leq \mathbf{v}^{(1)} \leq \dots \mathbf{v}^{(R)} \leq \mathbf{v}^*$;
2. $\forall i \in [0, R], \mathcal{T}_{\sigma^{(i)}}[\mathbf{v}^{(i)}] \geq \mathbf{v}^{(i)}, \mathcal{T}\mathbf{v}^{(i)} \geq [\mathbf{v}^{(i)}], \mathcal{H}_{\pi_{\max}^{(i)}}[\mathbf{v}^{(i)}] \geq \mathbf{v}^{(i)}$;
3. $\forall i \in [R], \mathbf{Q}^{(i)} \geq \mathbf{r} + \gamma \mathbf{P}\mathbf{v}^{(i-1)} - \boldsymbol{\epsilon}^{(i)}$;
4. $\forall i \in [R], \mathbf{v}^{(i)} \geq V[\mathbf{Q}^{(i)}]$.

Note that $\mathbf{Q}^{(0)}, \boldsymbol{\epsilon}^{(0)}$ can be arbitrary.

A.3 Monotone Value-Strategy Sequence Implies Good Strategy

Next, we show that MDVSS or MIVSS implies a good terminal value/strategy. First we show that if the terminal value $\mathbf{v}^{(R)}$ is close to the optimal value, then we are guaranteed to have good strategies as well.

Proposition A.3. *Suppose we have an MDVSS, $\{\mathbf{v}^{(i)}, \mathbf{Q}^{(i)}, \sigma^{(i)}, \boldsymbol{\epsilon}^{(i)}\}_{i=0}^R$, with $\|\mathbf{v}^{(R)} - \mathbf{v}^*\|_\infty \leq \epsilon$ for some $\epsilon \geq 0$. Then we have*

$$\mathbf{v}^{\pi_{\min}^{(R)}} \leq \mathbf{v}^* + \epsilon \mathbf{1}.$$

Similarly, suppose $\{\mathbf{v}^{(i)}, \mathbf{Q}^{(i)}, \sigma^{(i)}, \boldsymbol{\epsilon}^{(i)}\}_{i=0}^R$ is an MIVSS, then

$$\mathbf{v}^{\pi_{\max}^{(R)}} \geq \mathbf{v}^* - \epsilon \mathbf{1}.$$

Proof. By the property of an MDVSS, we have

$$\mathbf{v}^{\pi_{\min}^{(R)}} \leq \mathbf{v}^{(R)}.$$

Since $\mathbf{v}^{(R)} \leq \mathbf{v}^* + \epsilon \mathbf{1}$, we prove the first inequality. The second inequality follows similarly. \square

Next we consider when it is the case we achieve a good terminal value. The following proposition shows that an MDVSS(MIVSS) with an appropriate error parameters has a better terminal value than its initial value.

Proposition A.4. *Let $u \in (0, \beta), \beta = (1 - \gamma)^{-1}, R = \Theta[\beta \log(\beta/u)]$. Suppose an MDVSS (or MIVSS) $\{\mathbf{v}^{(i)}, \mathbf{Q}^{(i)}, \sigma^{(i)}, \boldsymbol{\epsilon}^{(i)}\}_{i=0}^R$ satisfies*

$$\|\mathbf{v}^{(0)} - \mathbf{v}^*\|_\infty \leq u \quad \text{and} \quad \boldsymbol{\epsilon}^{(i)} = \sqrt{L \cdot \text{var}(\mathbf{v}^{(0)})/m} + \beta \cdot (L/m)^{3/4} + u/(CR),$$

for some large constant $C > 1$ and $m \geq 1$. Then we have

$$\|\mathbf{v}^{(R)} - \mathbf{v}^*\|_\infty \leq u/2 \quad \text{for} \quad m = \tilde{\Omega}\left(\frac{1}{\min(1, u^2) \cdot (1 - \gamma)^3}\right).$$

Note that Proposition A.4 shows that in an MDVSS/MIVSS, the distance to the optimal value of the terminal value reduces by at least half of its initial value. Starting from some $\mathbf{v}^{(0)}$ with distance at most β to \mathbf{v}^* , by concatenating $O(\log(\beta/\epsilon))$ many MDVSS/MIVSS's, with the initial value of one sequence set as the terminal value of the last sequence, an ϵ -optimal value can be obtained. The remainder of this subsection devotes to proving the above proposition. Since MIVSS and MDVSS are symmetric, in the following analysis, we focus on MDVSS and the analysis follows similarly for MIVSS.

A.3.1 Auxiliary Markovian Strategy

Due to the lack of monotonicity we do not know how to use the optimal strategy σ^* to carefully account for the error accumulation of the MDVSS. To resolve this issue, we instead use the following auxiliary Markovian strategy as such a proxy.

Definition A.5 (Auxiliary Strategy). *Given a MDVSS, $\{\mathbf{v}^{(i)}, \mathbf{Q}^{(i)}, \sigma^{(i)}, \boldsymbol{\epsilon}^{(i)}\}_{i=0}^R$, we denote the Markovian auxiliary strategy for the max-player as*

$$\pi_{\text{aux max}}^{\infty(i)} = (\pi_{\text{aux max}}^{(i)}, \pi_{\text{aux max}}^{(i-1)}, \dots, \pi_{\text{aux max}}^{(1)}, \pi_{\text{max}}^*, \pi_{\text{max}}^*, \pi_{\text{max}}^* \dots),$$

where $\pi_{\text{aux max}}^{(i)}(s) = \arg \max_a \mathbf{Q}^{(i)}(s, a)$ for $s \in \mathcal{S}_{\text{max}}$. We denote the auxiliary strategy for the min-player as

$$\pi_{\text{aux min}}^{\infty(i)} = \sigma_{\text{min}}[\pi_{\text{aux max}}^{\infty(i)}] = (\pi_{\text{aux min}}^{(i)}, \pi_{\text{aux min}}^{(i-1)}, \dots, \pi_{\text{aux min}}^{(1)}, \pi_{\text{min}}^*, \pi_{\text{min}}^*, \pi_{\text{min}}^* \dots),$$

which is the optimal counter Markovian policy of $\pi_{\text{aux max}}^{\infty(i)}$, i.e.,

$$\forall s \in \mathcal{S}_{\text{min}} : \pi_{\text{aux min}}^{\infty(i)}(s) = \arg \min_a [\mathbf{r}(s, a) + \gamma \mathbf{P}(\cdot | s, a)^\top \mathbf{v}[\pi_{\text{aux}}^{\infty(i-1)}]].$$

We also denote

$$\sigma_{\text{aux}}^{\infty(i)} = \left[(\pi_{\text{aux min}}^{(i)}, \pi_{\text{aux max}}^{(i)}), (\pi_{\text{aux min}}^{(i-1)}, \pi_{\text{aux max}}^{(i-1)}), \dots, (\pi_{\text{aux min}}^{(1)}, \pi_{\text{aux max}}^{(1)}), \sigma^*, \sigma^*, \sigma^* \dots \right].$$

Furthermore, we denote $\sigma_{\text{aux}}^{(i)} = (\pi_{\text{aux min}}^{(i)}, \pi_{\text{aux max}}^{(i)})$ for $i \geq 1$ and $\sigma_{\text{aux}}^{(i)} = \sigma^*$ for $i \leq 0$.

For a Markovian strategy, we first show that the strategy has a value always smaller than the optimal value.

Lemma A.6. *For all $i \in [R]$, we have*

$$\mathbf{v}[\sigma_{\text{aux}}^{\infty(i)}] \leq \mathbf{v}^*.$$

Proof. Denote

$$\tilde{\sigma}_{\text{aux}}^{\infty(i)} = \left[(\pi_{\text{min}}^*, \pi_{\text{aux max}}^{(i)}), (\pi_{\text{min}}^*, \pi_{\text{aux max}}^{(i-1)}), \dots, (\pi_{\text{min}}^*, \pi_{\text{aux max}}^{(1)}), \sigma^*, \sigma^*, \sigma^* \dots \right].$$

Denote $\sigma_{\text{aux}}^{\infty(0)} = \tilde{\sigma}_{\text{aux}}^{\infty(0)} = (\sigma^*, \sigma^*, \dots)$. We first show that for all $i \in [R]$, $\mathbf{v}[\sigma_{\text{aux}}^{\infty(i)}] \leq \mathbf{v}[\tilde{\sigma}_{\text{aux}}^{\infty(i)}]$. Indeed it holds trivially for $i = 0$. Suppose it holds for some $i \geq 0$. Then, for each $s \in \mathcal{S}_{\text{min}}$, we have,

$$\begin{aligned} \mathbf{v}[\sigma_{\text{aux}}^{\infty(i)}](s) &= \min_a [\mathbf{r}(s, a) + \gamma \mathbf{P}(\cdot | s, a)^\top \mathbf{v}(\sigma_{\text{aux}}^{\infty(i-1)})] \\ &\leq [\mathbf{r}(s, \sigma^*(s)) + \gamma \mathbf{P}(\cdot | s, \sigma^*(s))^\top \mathbf{v}(\sigma_{\text{aux}}^{\infty(i-1)})] \\ &\leq [\mathbf{r}(s, \sigma^*(s)) + \gamma \mathbf{P}(\cdot | s, \sigma^*(s))^\top \mathbf{v}(\tilde{\sigma}_{\text{aux}}^{\infty(i-1)})] \quad (\text{due to } \mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}] \leq \mathbf{v}[\tilde{\sigma}_{\text{aux}}^{\infty(i-1)}]) \\ &= \mathbf{v}[\tilde{\sigma}_{\text{aux}}^{\infty(i)}](s). \end{aligned}$$

For each $s \in \mathcal{S}_{\text{max}}$, we have,

$$\begin{aligned} \mathbf{v}[\sigma_{\text{aux}}^{\infty(i)}](s) &= [\mathbf{r}(s, \sigma_{\text{aux}}^{(i)}(s)) + \gamma \mathbf{P}(\cdot | s, \sigma_{\text{aux}}^{(i)}(s))^\top \mathbf{v}(\sigma_{\text{aux}}^{\infty(i-1)})] \\ &\leq [\mathbf{r}(s, \sigma_{\text{aux}}^{(i)}(s)) + \gamma \mathbf{P}(\cdot | s, \sigma_{\text{aux}}^{(i)}(s))^\top \mathbf{v}(\tilde{\sigma}_{\text{aux}}^{\infty(i-1)})] \quad (\text{due to } \mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}] \leq \mathbf{v}[\tilde{\sigma}_{\text{aux}}^{\infty(i-1)}]) \\ &= \mathbf{v}[\tilde{\sigma}_{\text{aux}}^{\infty(i)}](s). \end{aligned}$$

Now, since $(\pi_{\text{max}}^*, \pi_{\text{max}}^*, \dots)$ is the optimal counter strategy of $(\pi_{\text{min}}^*, \pi_{\text{min}}^*, \dots)$, we have

$$\mathbf{v}[\tilde{\sigma}_{\text{aux}}^{\infty(i)}] \leq \mathbf{v}^*$$

holds similarly. This concludes the proof. \square

Consider the error vector $\boldsymbol{\epsilon}^{(i)}$. Recall that $\boldsymbol{\epsilon}_{\sigma_{\text{aux}}^{(i)}}^{(i)}$ denotes a vector in $\mathbb{R}^{\mathcal{S}}$ whose s -th entry is given by $\boldsymbol{\epsilon}^{(i)}(s, \sigma_{\text{aux}}^{(i)}(s))$. The next lemma shows a recursive relation between a Markovian strategy and the corresponding MDVSS values.

Lemma A.7. For all $i \in [R]$, we have

$$\mathbf{v}^{(i)} - \mathbf{v}[\sigma_{\text{aux}}^{\infty(i)}] \leq \gamma \mathbf{P}^{\sigma_{\text{aux}}^{(i)}} (\mathbf{v}^{(i-1)} - \mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}]) + \boldsymbol{\epsilon}_{\sigma_{\text{aux}}^{(i)}}^{(i)}$$

Proof. Note that $\mathbf{v}^{(i)} \geq \mathbf{v}^* \geq \mathbf{v}[\sigma_{\text{aux}}^{\infty(i)}]$. For each $s \in \mathcal{S}_{\min}$, we have

$$\begin{aligned} \mathbf{v}^{(i)}(s) &\leq \min_a \mathbf{Q}^{(i)}(s, a) \leq \mathbf{Q}^{(i)}(s, \sigma_{\text{aux}}^{(i)}(s)) \\ &\leq \mathbf{r}(s, \sigma_{\text{aux}}^{(i)}(s)) + \gamma \mathbf{P}(\cdot | s, \sigma_{\text{aux}}^{(i)}(s))^\top \mathbf{v}^{(i-1)} + \boldsymbol{\epsilon}^{(i)}(s, \sigma_{\text{aux}}^{(i)}(s)), \end{aligned}$$

and

$$\mathbf{v}[\sigma_{\text{aux}}^{\infty(i)}](s) = \mathbf{r}(s, \sigma_{\text{aux}}^{(i)}(s)) + \gamma \mathbf{P}(\cdot | s, \sigma_{\text{aux}}^{(i)}(s))^\top \mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}].$$

Thus

$$\mathbf{v}^{(i)}(s) - \mathbf{v}[\sigma_{\text{aux}}^{\infty(i)}](s) \leq \gamma \mathbf{P}(\cdot | s, \sigma_{\text{aux}}^{(i)}(s))^\top (\mathbf{v}^{(i-1)} - \mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}]) + \boldsymbol{\epsilon}^{(i)}(s, \sigma_{\text{aux}}^{(i)}(s)).$$

Similarly, for each $s \in \mathcal{S}_{\max}$, we have, $\mathbf{v}^{(i)}(s) \leq \max_a \mathbf{Q}^{(i)}(s, a) := \mathbf{Q}^{(i)}(s, \sigma_{\text{aux}}^{(i)}(s))$, thus

$$\begin{aligned} \mathbf{v}^{(i)}(s) - \mathbf{v}[\sigma_{\text{aux}}^{\infty(i)}](s) &\leq \mathbf{Q}^{(i)}(s, \sigma_{\text{aux}}^{(i)}(s)) - \mathbf{v}[\sigma_{\text{aux}}^{\infty(i)}](s) \\ &\leq \gamma \mathbf{P}(\cdot | s, \sigma_{\text{aux}}^{(i)}(s))^\top (\mathbf{v}^{(i-1)} - \mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}]) + \boldsymbol{\epsilon}^{(i)}(s, \sigma_{\text{aux}}^{(i)}(s)) \end{aligned}$$

as desired. \square

With an inductive application of the above lemma, we obtain the following corollary, which states an upper bound between the difference of $\mathbf{v}^{(R)}$ and $\mathbf{v}[\sigma_{\text{aux}}^{\infty(R)}]$. It connects the upper bound with a recursive propagation of the error.

Corollary A.8.

$$\begin{aligned} \mathbf{v}^{(R)} - \mathbf{v}[\sigma_{\text{aux}}^{\infty(R)}] &\leq \gamma^R \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdot \dots \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(1)}} (\mathbf{v}^{(0)} - \mathbf{v}^*) \\ &\quad + \sum_{i=1}^R \gamma^{R-i} \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdot \dots \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(i+1)}} \boldsymbol{\epsilon}_{\sigma_{\text{aux}}^{(i)}}^{(i)}. \end{aligned}$$

By this corollary, we know that the major error accumulation term is the second term.

A.3.2 Error Accumulation

We now consider the error accumulation in the sequence. As will show shortly, we relate $\boldsymbol{\epsilon}^{(i)}$ to the variance vector $\sqrt{\text{var}[\mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}]]} \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$, where $\text{var}(\mathbf{v})[s, a] := \text{var}_{s' \sim P(\cdot | s, a)}[\mathbf{v}(s')]$, $\forall (s, a), \mathbf{v}$. Therefore, it suffices to consider the following bound.

Lemma A.9.

$$\begin{aligned} &\sum_{i=1}^R \gamma^{R-i} \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdot \dots \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(i+1)}} \sqrt{\text{var}[\mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}]]_{\sigma_{\text{aux}}^{(i)}}} \\ &\leq \sqrt{R \sum_{i=1}^R \gamma^{2(R-i)} \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdot \dots \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(i+1)}} \text{var}[\mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}]]_{\sigma_{\text{aux}}^{(i)}}} \end{aligned}$$

Proof. Follows from Cauchy-Schwarz and that the \mathbf{P} matrices are non-negative with each row summing to 1. \square

The following lemma establishes a Bellman-like equation for the variance vector of a Markovian strategy.

Lemma A.10. For any Markovian strategy $\pi^\infty = (\pi^{(0)}, \pi^{(1)}, \dots)$, we have, for all $s \in \mathcal{S}$

$$\text{var} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s^t, \pi^{(t)}(s)) \middle| s^0 = s \right] = \left[\sum_{t=0}^{\infty} \gamma^{2(t+1)} \mathbf{P}^{\pi^{(0)}} \mathbf{P}^{\pi^{(1)}} \mathbf{P}^{\pi^{(2)}} \dots \mathbf{P}^{\pi^{(t-1)}} \text{var}[\mathbf{v}(\pi^{\infty(t+1)})]_{\pi^{(t)}} \right] (s) \quad (8)$$

Proof.

$$\text{var} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s^t, \pi^{(t)}(s^t)) \middle| s^0 = s \right] = \mathbb{E} \left[\left(\sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s^t, \pi^{(t)}(s^t)) \right)^2 \middle| s^0 = s \right] - \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s^t, \pi^{(t)}(s^t)) \middle| s^0 = s \right]^2.$$

For the second term, we have

$$\begin{aligned} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s^t, \pi^{(t)}(s^t)) \middle| s^0 = s \right]^2 &= \mathbf{v}[\pi^{\infty(0)}]^2(s) = \mathbf{r}(s, \pi^{(0)}(s))^2 \\ &\quad + \gamma^2 (\mathbf{P}^{\pi^{(0)}} \mathbf{v}[\pi^{\infty(1)}])^2(s) + 2\gamma \mathbf{r}(s, \pi^{(0)}(s)) (\mathbf{P}^{\pi^{(0)}} \mathbf{v}[\pi^{\infty(1)}])(s). \end{aligned}$$

For the first term, we have

$$\left(\sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s^t, \pi^{(t)}(s^t)) \right)^2 = \mathbf{r}(s, \pi^{(0)}(s))^2 + 2\mathbf{r}(s, \pi^{(0)}(s)) \left(\sum_{t=1}^{\infty} \gamma^t \mathbf{r}(s^t, \pi^{(t)}(s^t)) \right) + \left(\sum_{t=1}^{\infty} \gamma^t \mathbf{r}(s^t, \pi^{(t)}(s^t)) \right)^2$$

Note that

$$\begin{aligned} &\mathbb{E} \left[\left(\sum_{t=1}^{\infty} \gamma^t \mathbf{r}(s^t, \pi^{(t)}(s^t)) \right)^2 \middle| s^0 = s \right] \\ &= \mathbb{E} \left[\left(\sum_{t=1}^{\infty} \gamma^t \mathbf{r}(s^t, \pi^{(t)}(s^t)) \right)^2 \middle| s^0 = s \right] - \gamma^2 \sum_{s'} \mathbf{P}^{\pi^{(0)}}(s'|s) \mathbf{v}^2[\pi^{\infty(1)}](s') + \gamma^2 \sum_{s'} \mathbf{P}^{\pi^{(0)}}(s'|s) \mathbf{v}^2[\pi^{\infty(1)}](s') \\ &= \gamma^2 \sum_{s'} \mathbf{P}^{\pi^{(0)}}(s'|s) \text{var} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s^{t+1}, \pi^{(t+1)}(s)) \middle| s^1 = s' \right] + \gamma^2 \sum_{s'} \mathbf{P}^{\pi^{(0)}}(s'|s) \mathbf{v}^2[\pi^{\infty(1)}](s') \end{aligned}$$

Combining the above two equations, we have,

$$\begin{aligned} \mathbb{E} \left[\left(\sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s^t, \pi^{(t)}(s^t)) \right)^2 \middle| s^0 = s \right] &= \mathbf{r}(s, \pi^{(0)}(s))^2 + 2\gamma \mathbf{r}(s, \pi^{(0)}(s)) \sum_{s'} \mathbf{P}^{\pi^{(0)}}(s'|s) \mathbf{v}[\pi^{\infty(1)}](s') \\ &\quad + \gamma^2 \sum_{s'} \mathbf{P}^{\pi^{(0)}}(s'|s) \text{var} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s^{t+1}, \pi^{(t+1)}(s)) \middle| s^1 = s' \right] \\ &\quad + \gamma^2 \sum_{s'} \mathbf{P}^{\pi^{(0)}}(s'|s) \mathbf{v}^2[\pi^{\infty(1)}](s') \end{aligned}$$

We thus obtain

$$\begin{aligned} \text{var} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s^t, \pi^{(t)}(s^t)) \middle| s^0 = s \right] &= \gamma^2 \sum_{s'} \mathbf{P}^{\pi^{(0)}}(s'|s) \text{var} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s^{t+1}, \pi^{(t+1)}(s)) \middle| s^1 = s' \right] \\ &\quad + \gamma^2 \sum_{s'} \mathbf{P}^{\pi^{(0)}}(s'|s) \mathbf{v}^2[\pi^{\infty(1)}](s') - \gamma^2 (\mathbf{P}^{\pi^{(0)}} \mathbf{v}[\pi^{\infty(1)}])^2(s) \\ &= \gamma^2 \sum_{s'} \mathbf{P}^{\pi^{(0)}}(s'|s) \text{var} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s^{t+1}, \pi^{(t+1)}(s)) \middle| s^1 = s' \right] + \gamma^2 \text{var}((\mathbf{v}[\pi^{\infty(1)}]))_{\pi^{(0)}} \end{aligned}$$

Let LHS and RHS be the left hand side and right hand side of (8) respectively. Then we have,

$$\begin{aligned} LHS &= \gamma^2 \sum_{s'} \mathbf{P}^{\pi^{(0)}}(s'|s) \text{var} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s^{t+1}, \pi^{(t+1)}(s)) \middle| s^1 = s' \right] + \gamma^2 \text{var}((\mathbf{v}[\pi^{\infty(1)}]))_{\pi^{(0)}} \\ &= \gamma^4 \sum_{s', s''} \mathbf{P}^{\pi^{(0)}}(s'|s) \mathbf{P}^{\pi^{(1)}}(s''|s') \text{var} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}(s^{t+2}, \pi^{(t+2)}(s)) \middle| s^2 = s'' \right] + \gamma^4 \mathbf{P}^{\pi^{(0)}} \text{var}((\mathbf{v}[\pi^{\infty(2)}]))_{\pi^{(1)}} \\ &\quad + \gamma^2 \text{var}((\mathbf{v}[\pi^{\infty(1)}]))_{\pi^{(0)}}. \end{aligned}$$

Applying the above equality recursively for $\text{var}((\mathbf{v}[\pi^{\infty(i)}]))$ completes the proof. \square

Based on the above two lemmas, we immediately obtain the following worst-case bound for the error accumulation.

Corollary A.11.

$$\sqrt{R \sum_{i=1}^R \gamma^{2(R-i)} \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdot \dots \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(i+1)}} \text{var}(\mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}])_{\sigma_{\text{aux}}^{(i)}}} \leq \sqrt{\frac{R}{\gamma^2(1-\gamma)^2}}.$$

Proof. We use that

$$\begin{aligned} & \left[\sum_{i=1}^R \gamma^{2(R-i)} \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdot \dots \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(i+1)}} \text{var}(\mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}])_{\sigma_{\text{aux}}^{(R-i)}} \right] (s) \\ & \leq \frac{1}{\gamma^2} \cdot \text{var} \left[\sum_{i=0}^{\infty} \gamma^i \mathbf{r}(s^i, \sigma_{\text{aux}}^{(R-i)}(s)) \middle| s^0 = s \right]. \end{aligned}$$

Since $\sum_{i=0}^{\infty} \gamma^i \mathbf{r}(s^i, \sigma_{\text{aux}}^{(i)}(s)) \leq (1-\gamma)^{-1}$, we have $\text{var}(\mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}])_{\sigma_{\text{aux}}^{(R-i)}} \leq (1-\gamma)^{-2}$ as desired. \square

A.3.3 Putting Everything Together

Proof of Proposition A.4. By Corollary A.8, we have,

$$\begin{aligned} \mathbf{v}^{(R)} - \mathbf{v}^* & \leq \mathbf{v}^{(R)} - \mathbf{v}[\sigma_{\text{aux}}^{(R)}] \leq \gamma^R \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdot \dots \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(1)}} (\mathbf{v}^{(0)} - \mathbf{v}^*) \\ & \quad + \sum_{i=1}^R \gamma^{R-i} \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdot \dots \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(i+1)}} \boldsymbol{\epsilon}_{\sigma_{\text{aux}}^{(i)}}^{(i)} \\ & \leq u/4 + \underbrace{\sum_{i=1}^R \gamma^{R-i} \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdot \dots \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(i+1)}} \boldsymbol{\epsilon}_{\sigma_{\text{aux}}^{(i)}}^{(i)}}_{\textcircled{1}}, \end{aligned}$$

where the first inequality holds for sufficiently large R . Consider the second term. Since

$$\boldsymbol{\epsilon}^{(i)} = \sqrt{L \cdot \text{var}(\mathbf{v}^{(0)})/m} + \beta \cdot (L/m)^{3/4} + u/(CR).$$

We bound

$$\sum_{i=1}^R \gamma^{R-i} \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdot \dots \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(i+1)}} \cdot \beta \cdot (L/m)^{3/4} \leq R\beta \cdot (L/m)^{3/4}$$

and

$$\sum_{i=1}^R \gamma^{R-i} \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdot \dots \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(i+1)}} \cdot u/(CR) \leq Ru/(CR).$$

We thus have,

$$\textcircled{1} \leq \sum_{i=1}^R \gamma^{R-i} \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdot \dots \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(i+1)}} \sqrt{L \cdot \text{var}(\mathbf{v}^{(0)})_{\sigma_{\text{aux}}^{(i)}}/m} + R\beta \cdot (L/m)^{3/4} + R \cdot u/(CR).$$

Note that

$$\sqrt{\text{var}(\mathbf{v}^{(0)})_{\sigma_{\text{aux}}^{(i)}}} \leq \sqrt{\text{var}(\mathbf{v}^{\sigma_{\text{aux}}^{\infty(i-1)}})_{\sigma_{\text{aux}}^{(i)}}} + \|\mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}] - \mathbf{v}^{(0)}\|_{\infty}.$$

Now consider

$$\|\mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}] - \mathbf{v}^{(0)}\|_{\infty} \leq \|\mathbf{v}^{(0)} - \mathbf{v}^{(R)}\|_{\infty} + \|\mathbf{v}[\sigma_{\text{aux}}^{\infty(i-1)}] - \mathbf{v}^{(i-1)}\|_{\infty}.$$

We bound $\|\mathbf{v}^{(0)} - \mathbf{v}^{(R)}\|_\infty \leq u$. Applying Corollary A.8 again, we have

$$\begin{aligned} \|\mathbf{v}[\sigma_{\text{aux}}^{(i-1)}] - \mathbf{v}^{(i-1)}\|_\infty &\leq u/4 + \sum_{i=1}^R \gamma^{R-i} \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdots \mathbf{P}^{\sigma_{\text{aux}}^{(i+1)}} \boldsymbol{\epsilon}_{\sigma_{\text{aux}}^{(i)}}^{(i)} \\ &\leq u/4 + \sum_{i=1}^R \gamma^{R-i} \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdots \mathbf{P}^{\sigma_{\text{aux}}^{(i+1)}} \sqrt{L \cdot \text{var}(\mathbf{v}^{(0)})_{\sigma_{\text{aux}}^{(i)}}/m} \\ &\quad + R\beta \cdot (L/m)^{3/4} + R \cdot u/(CR) \end{aligned}$$

With a natural bound, $\text{var}[\mathbf{v}^{(0)}] \leq \beta^2 \cdot \mathbf{1}$, we have

$$\|\mathbf{v}[\sigma_{\text{aux}}^{(i-1)}] - \mathbf{v}^{(0)}\|_\infty \leq u + \|\mathbf{v}[\sigma_{\text{aux}}^{(i-1)}] - \mathbf{v}^{(i-1)}\|_\infty \leq R\beta\sqrt{L/m} + R\beta(L/m)^{3/4} + u/C'$$

for some constant $C' > 0$. Therefore,

$$\begin{aligned} \mathbf{v}^{(R)} - \mathbf{v}^* &\leq \sum_{i=1}^R \gamma^{R-i} \mathbf{P}^{\sigma_{\text{aux}}^{(R)}} \cdot \mathbf{P}^{\sigma_{\text{aux}}^{(R-1)}} \cdots \mathbf{P}^{\sigma_{\text{aux}}^{(i+1)}} \sqrt{\frac{L \text{var}(\mathbf{v}[\sigma_{\text{aux}}^{(i-1)}])_{\pi_{\text{aux}}^{(i)}}}{m}} \\ &\quad + R \left(R\beta\sqrt{\frac{L}{m}} + R\beta\left(\frac{L}{m}\right)^{3/4} + \frac{u}{C'} \right) \cdot \sqrt{\frac{L}{m}} + \frac{u}{4} + \frac{u}{C} + R\beta \cdot \left(\frac{L}{m}\right)^{3/4} \\ &\leq \sqrt{\frac{LR\beta^2}{\gamma^2 m}} + R \left(R\beta\sqrt{\frac{L}{m}} + R\beta\left(\frac{L}{m}\right)^{3/4} + \frac{u}{C'} \right) \cdot \sqrt{\frac{L}{m}} + \frac{u}{4} + \frac{u}{C} + R\beta \cdot \left(\frac{L}{m}\right)^{3/4} \\ &\leq u/2 \end{aligned}$$

for large enough constants, C , and that in m . □

A.4 Algorithm that Computes a Monotone Sequence

Here we show that Algorithm 1 or Algorithm 2 computes an MDVSS or MIVSS respectively.

Proposition A.12. *Let $u \in (0, \beta]$, $\beta = (1 - \gamma)^{-1}$, $\delta \in (0, 1)$, and $R = \Theta[\beta \log(\beta/u)]$. Further, let $L = \Theta(\log[\delta^{-1}\beta R|\mathcal{S}||\mathcal{A}|])$, and $m = \Omega(\beta^3 \cdot \max(u^{-2}, 1) \cdot \log(|\mathcal{S}||\mathcal{A}|\delta^{-1}))$. Then there exists an algorithm, on input a stochastic game with a generative model with a sampling oracle, $\mathcal{G} = (\mathcal{S} := \mathcal{S}_{\min} \cup \mathcal{S}_{\max}, \mathbf{P}, \mathbf{r}, \gamma)$, a value-strategy pair $(\pi^{(0)}, \mathbf{v}^{(0)})$ satisfying $\mathcal{T}_{\pi^{(0)}}[\mathbf{v}^{(0)}] \leq \mathbf{v}^{(0)}$, $\mathcal{T}[\mathbf{v}^{(0)}] \leq \mathbf{v}^{(0)}$ (or $\mathcal{T}_{\pi^{(0)}}[\mathbf{v}^{(0)}] \geq \mathbf{v}^{(0)}$, $\mathcal{T}[\mathbf{v}^{(0)}] \geq \mathbf{v}^{(0)}$), and $\|\mathbf{v}^{(0)} - \mathbf{v}^*\|_\infty \leq u$ for some $u > 0$, outputs, with probability at least $1 - \delta$, an MDVSS (or MIVSS) $\{\mathbf{v}^{(i)}, \mathbf{Q}^{(i)}, \pi^{(i)}, \boldsymbol{\epsilon}^{(i)}\}_{i=0}^R$ by querying*

$$Z = O[|\mathcal{S}||\mathcal{A}| \cdot (m + R\beta^2 \log[R|\mathcal{S}||\mathcal{A}|\delta^{-1}])]$$

samples, where

$$\boldsymbol{\epsilon}^{(i)} = \sqrt{L\boldsymbol{\sigma}_{\mathbf{v}^{(0)}}/m} + \beta \cdot (L/m)^{3/4} + u/(CR),$$

for some large constant $C > 1$ and $\boldsymbol{\sigma}_{\mathbf{v}^{(0)}} := \text{var}[\mathbf{v}^{(0)}]$ is the variance vector for vector $\mathbf{v}^{(0)}$. The algorithm uses space $O(|\mathcal{S}||\mathcal{A}|)$ and halts in time $O(Z)$.

This section is devoted to proving Proposition A.12. The algorithm of obtaining MDVSS and MIVSS is provided in Algorithm 1 and 2.

The Good Events Suppose we are given an arbitrary input vector $\mathbf{v}^{-(0)}, \mathbf{v}^{+(0)} \in [0, (1 - \gamma)^{-1}]^{\mathcal{S}}$ with $\mathbf{v}^* - u\mathbf{1} \leq \mathbf{v}^{-(0)} \leq \mathbf{v}^* \leq \mathbf{v}^{+(0)} \leq \mathbf{v}^* + u\mathbf{1}$, $\mathbf{v}^{-(0)} \leq \mathcal{T}[\mathbf{v}^{-(0)}]$, $\mathbf{v}^{-(0)} \leq \mathcal{T}_{\pi^-}[\mathbf{v}^{-(0)}]$, $\mathcal{T}[\mathbf{v}^{+(0)}] \leq \mathbf{v}^{+(0)}$ and $\mathcal{T}_{\pi^+}[\mathbf{v}^{+(0)}] \leq \mathbf{v}^{+(0)}$. Since the algorithm is randomized, to begin our analysis, we define a sequence of events for the iterates. We will show that these events happen with high probability via concentration inequalities.

Definition A.13. *Let $\tilde{\mathbf{w}}^-$ and $\tilde{\mathbf{w}}^+$ be the estimate defined in Line 16 (of Algorithm 1 or Algorithm 2 respectively). Denote $\alpha_1 \leftarrow L/m_1 \leq 1$. Let \mathcal{E}_0 be the event that*

$$|\tilde{\mathbf{w}}^\pm - \mathbf{P}\mathbf{v}^{\pm(0)}| \leq \sqrt{\alpha_1 \boldsymbol{\sigma}_{\mathbf{v}^{\pm(0)}}} + \alpha_1^{3/4} \cdot \|\mathbf{v}^{\pm(0)}\|_\infty \cdot \mathbf{1}. \quad (9)$$

For each $i > 0$, let $\tilde{\mathbf{g}}^{\pm(i)}$ be given in Line 32. Let \mathcal{E}_i be the event that

$$|\tilde{\mathbf{g}}^{\pm(i)} - \mathbf{P}[\mathbf{v}^{\pm(i)} - \mathbf{v}^{\pm(0)}]| \leq C(1 - \gamma)u \cdot \mathbf{1}. \quad (10)$$

for some sufficiently small constant $C > 0$.

Lemma A.14. For some sufficiently large constant c in L , $\Pr[\mathcal{E}_0] \geq 1 - O(\delta/R)$.

Proof. Note that $\|\mathbf{v}^{\pm(0)}\|_\infty \leq (1 - \gamma)^{-1}$. By a straightforward application of a Hoeffding bound and Bernstein inequality and a union bound over all (s, a) , we reach the desired inequality. More details can be find in the proof of Lemma 5.1 in Sidford et al. (2018a). \square

The Implications of the Good Events We now illustrate the consequences of these good events.

Lemma A.15 (Implications of \mathcal{E}_0). On \mathcal{E}_0 , we have, for all $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned} 0 \leq \mathbf{r}(s, a) + \gamma \mathbf{P}(\cdot | s, a)^\top \mathbf{v}^{- (0)} - \mathbf{Q}^{- (0)}(s, a) &\leq 2\sqrt{\alpha_1 \sigma_{\mathbf{v}^{- (0)}}} + 2\alpha_1^{3/4} \|\mathbf{v}^{- (0)}\|_\infty \quad \text{and} \\ 0 \leq \mathbf{Q}^{- (0)}(s, a) &\leq \mathbf{Q}^*(s, a), \end{aligned}$$

and

$$\begin{aligned} 0 \leq \mathbf{Q}^{+ (0)}(s, a) - \mathbf{r}(s, a) - \gamma \mathbf{P}(\cdot | s, a)^\top \mathbf{v}^{+ (0)} &\leq 2\sqrt{\alpha_1 \sigma_{\mathbf{v}^{+ (0)}}} + 2\alpha_1^{3/4} \|\mathbf{v}^{+ (0)}\|_\infty \quad \text{and} \\ \mathbf{Q}^*(s, a) \leq \mathbf{Q}^{+ (0)}(s, a) &\leq \beta, \end{aligned}$$

where $\alpha_1 = L/m_1$.

Proof. We prove the first inequality and the second inequality follows similarly. Condition on \mathcal{E}_0 , we have

$$|\mathbf{r} + \gamma \tilde{\mathbf{w}}^- - \mathbf{r} - \gamma \mathbf{P} \mathbf{v}^{- (0)}| \leq \sqrt{\alpha_1 \sigma_{\mathbf{v}^{- (0)}}} + \alpha_1^{3/4} \|\mathbf{v}^{- (0)}\|_\infty.$$

Since

$$\mathbf{Q}^{- (0)} = \max \left[\mathbf{r} + \gamma \tilde{\mathbf{w}}^- - \sqrt{\alpha_1 \sigma_{\mathbf{v}^{- (0)}}} - \alpha_1^{3/4} \|\mathbf{v}^{- (0)}\|_\infty, \mathbf{0} \right],$$

we have

$$0 \leq \mathbf{r} + \gamma \mathbf{P} \mathbf{v}^{- (0)} - \mathbf{Q}^{- (0)} \leq 2\sqrt{\alpha_1 \sigma_{\mathbf{v}^{- (0)}}} + 2\alpha_1^{3/4} \|\mathbf{v}^{- (0)}\|_\infty.$$

Moreover, since $\mathbf{v}^{(0)} \leq \mathbf{v}^*$, we have

$$\mathbf{Q}^{- (0)} \leq \mathbf{r}(s, a) + \gamma \mathbf{P}(\cdot | s, a)^\top \mathbf{v}^{- (0)} \leq \mathbf{Q}^*,$$

completing the proof. \square

Lemma A.16 (Implications of \mathcal{E}_i , (1)). Then for any $i > 0$, conditioning on $\mathcal{E}_0, \mathcal{E}_1, \dots, \mathcal{E}_R$, we have $\{\mathbf{v}^{- (i)}, \mathbf{Q}^{- (i)}, \pi^{- (i)}, \boldsymbol{\xi}^{- (i)}\}_{i=0}^R$ is an MIVSS and $\{\mathbf{v}^{+ (i)}, \mathbf{Q}^{+ (i)}, \pi^{+ (i)}, \boldsymbol{\xi}^{+ (i)}\}_{i=0}^R$ is an MDVSS where

$$\boldsymbol{\xi}^{\pm (i)} = 2\sqrt{\alpha_1 \sigma_{\mathbf{v}^{\pm (0)}}} + 2[\alpha_1^{3/4} \beta + C(1 - \gamma)u] \cdot \mathbf{1}$$

for some sufficiently small $C > 0$.

Proof. We prove the first part of the lemma, i.e., $\{\mathbf{v}^{- (i)}, \mathbf{Q}^{- (i)}, \pi^{- (i)}, \boldsymbol{\xi}^{- (i)}\}_{i=0}^R$ is an MIVSS. Then the second part follows similarly. It is clear from the definition (Line 32) that

$$\mathbf{v}^{- (0)} \leq \mathbf{v}^{- (1)} \dots \leq \mathbf{v}^{- (R)}.$$

To prove property 1 of MIVSS, we need additionally to show

$$\mathbf{v}^{- (R)} \leq \mathbf{v}^*.$$

This follows if property 2, i.e.,

$$\forall i \in [0, R]: \quad \mathbf{v}^{- (i)} \leq \mathcal{T}_{\pi^{- (i)}}[\mathbf{v}^{- (i)}], \quad \mathbf{v}^{- (i)} \leq \mathcal{T}[\mathbf{v}^{- (i)}].$$

Indeed,

$$\mathbf{v}^{-(i)} \leq \mathcal{T}[\mathbf{v}^{-(i)}] \leq \mathcal{T}^2[\mathbf{v}^{-(i)}] \dots \leq \mathcal{T}^\infty[\mathbf{v}^{-(i)}] = \mathbf{v}^*.$$

We now prove property 2 by induction on i . It immediately follows from the initial condition that

$$\mathbf{v}^{-(0)} \leq \mathcal{T}_{\pi^{-(0)}}[\mathbf{v}^{-(0)}], \quad \mathbf{v}^{-(0)} \leq \mathcal{T}[\mathbf{v}^{-(0)}].$$

Suppose this property holds for all $0, 1, \dots, i-1$ for some $i > 1$. We now consider the case i . Let $\tilde{\mathbf{Q}}^{-(i)} = \mathbf{r} + \gamma \mathbf{P} \mathbf{v}^{-(i-1)}$. Since $\|\tilde{\mathbf{g}}^{-(i)} - \mathbf{P}(\mathbf{v}^{-(i-1)} - \mathbf{v}^{-(0)})\|_\infty \leq C(1-\gamma)\epsilon$, we have

$$|\tilde{\mathbf{Q}}^{-(i)} - \mathbf{r} - \gamma \mathbf{w}^{-(i-1)} - \gamma \mathbf{g}^{-(i-1)}| \leq \boldsymbol{\xi}^{-(i)}/2. \quad (11)$$

Since

$$\begin{aligned} \mathbf{Q}^{-(i)}(s, a)y &= \max \left[\mathbf{r}(s, a) + \gamma \mathbf{w}^{-(i-1)}(s, a) + \gamma \mathbf{g}^{-(i-1)}(s, a), 0 \right] \leq \mathbf{r}(s, a) + \gamma \mathbf{P}(\cdot|s, a)^\top \mathbf{v}^{-(i-1)} \\ &= \tilde{\mathbf{Q}}^{-(i)}(s, a) \end{aligned}$$

we have, for any $s \in \mathcal{S}$

$$\min_a \mathbf{Q}^{-(i)}(s, a) \leq \min_a \tilde{\mathbf{Q}}^{-(i)}(s, a) \quad \text{and} \quad \max_a \mathbf{Q}^{-(i)}(s, a) \leq \max_a \tilde{\mathbf{Q}}^{-(i)}(s, a).$$

For each $s \in \mathcal{S}$, denote

$$\tilde{\mathbf{v}}^{-(i)}(s) = \mathbf{Q}^{-(i)}(s, \tilde{\pi}^{-(i)}(s)),$$

where $\tilde{\pi}^{-(i)}$ is given in Line 26 (that achieves $\max_a \mathbf{Q}^{-(i)}(s, a)$ or $\min_a \mathbf{Q}^{-(i)}(s, a)$). To show $\mathbf{v}^{-(i)} \leq \mathcal{T}[\mathbf{v}^{-(i)}]$ and $\mathbf{v}^{-(i)} \leq \mathcal{T}_{\pi^{-(i)}}[\mathbf{v}^{-(i)}]$, we do a case analysis for a state in \mathcal{S} . Firstly, we consider state $s \in \mathcal{S}_{\min}$. For each state $s \in \mathcal{S}_{\min}$, note that $\tilde{\pi}^{-(i)}(s) := \arg \min_a \mathbf{Q}^{-(i)}(s, a)$. By Line 32, $\mathbf{v}^{-(i)}(s)$ and $\pi^{-(i)}(s)$ have the following two possibilities,

1. $\mathbf{v}^{-(i-1)}(s) \leq \tilde{\mathbf{v}}^{-(i)}(s) \Rightarrow \mathbf{v}^{-(i)}(s) = \tilde{\mathbf{v}}^{-(i)}(s)$ and $\pi^{-(i)}(s) = \tilde{\pi}^{-(i)}(s)$;
2. $\mathbf{v}^{-(i-1)}(s) > \tilde{\mathbf{v}}^{-(i)}(s) \Rightarrow \mathbf{v}^{-(i)}(s) = \mathbf{v}^{-(i-1)}(s)$ and $\pi^{-(i)}(s) = \pi^{-(i-1)}(s)$.

Considering case 1., we have,

$$\begin{aligned} \mathbf{v}^{-(i)}(s) &= \mathbf{Q}^{-(i-1)}(s, \pi^{-(i)}(s)) \leq \min_a \tilde{\mathbf{Q}}^{-(i)}(s, a) = \mathcal{T}[\mathbf{v}^{-(i-1)}](s) \leq \mathcal{T}[\mathbf{v}^{-(i)}](s) \quad \text{and} \\ \mathbf{v}^{-(i)}(s) &= \mathbf{Q}^{-(i-1)}(s, \pi^{-(i)}(s)) \leq \tilde{\mathbf{Q}}^{-(i)}(s, \pi^{-(i)}(s)) = \mathcal{T}_{\pi^{-(i)}}[\mathbf{v}^{-(i-1)}](s) \leq \mathcal{T}_{\pi^{-(i)}}[\mathbf{v}^{-(i)}](s). \end{aligned}$$

Considering case 2., we have, by induction hypothesis $\mathbf{v}^{-(i-1)}(s) \leq \mathcal{T}[\mathbf{v}^{-(i-1)}](s)$, $\mathbf{v}^{-(i-1)}(s) \leq \mathcal{T}_{\pi^{-(i-1)}}[\mathbf{v}^{-(i-1)}](s)$. Thus

$$\begin{aligned} \mathbf{v}^{-(i)}(s) &\leq \mathcal{T}[\mathbf{v}^{-(i-1)}](s) \leq \mathcal{T}[\mathbf{v}^{-(i)}](s), \\ \mathbf{v}^{-(i)}(s) &\leq \mathcal{T}_{\pi^{-(i-1)}}[\mathbf{v}^{-(i-1)}](s) \leq \mathcal{T}_{\pi^{-(i-1)}}[\mathbf{v}^{-(i)}](s) = \mathcal{T}_{\pi^{-(i)}}[\mathbf{v}^{-(i)}](s). \end{aligned}$$

It follows similarly for the case of $s \in \mathcal{S}_{\max}$ (by just replacing the min by max in the above argument). This completes the induction step and hence the property 2.

We now prove property 3 and 4. By Equation 11, we immediately have,

$$\mathbf{Q}^{-(i)} \geq \mathbf{r} + \gamma \mathbf{P}(\mathbf{v}^{-(i)}) - \boldsymbol{\xi}^{-(i)}$$

proving property 3. Lastly, by Line 32, we have

$$\mathbf{v}^{-(i)} \geq \mathcal{T}[\mathbf{Q}^{-(i)}],$$

completing the proof of property 4 and the lemma. \square

The Probability of Good Events Note that the random samples in the successive improvement phase are independent with event \mathcal{E}_0 . We then have the following lemma.

Lemma A.17. *Suppose the algorithm does not halt at iteration $i \geq 1$, then,*

$$\Pr[\mathcal{E}_i | \mathcal{E}_0, \mathcal{E}_1, \dots, \mathcal{E}_{i-1}] \geq 1 - O(\delta/R).$$

Proof. On $\mathcal{E}_0, \mathcal{E}_1, \dots, \mathcal{E}_{i-1}$, we have $\mathbf{v}^{-(0)} \leq \mathbf{v}^{-(i)} \leq \mathbf{v}^* \leq \mathbf{v}^{+(i)} \leq \mathbf{v}^{+(0)}$, thus $\|\mathbf{v}^{\pm(i)} - \mathbf{v}^{\pm(0)}\| \leq u$. Applying Hoeffding bound, Bernstein's inequality and a union bound over all (s, a) , we have that with probability at least $1 - O(\delta/R)$, $\|g^{(i)} - P[v^{(i)} - \mathbf{v}^{(0)}]\|_\infty \leq (1 - \gamma)\epsilon/32$, completing the proof. \square

Therefore, we have the following lemma.

Lemma A.18. *Let $R = \lceil c_1 \beta \ln[\beta \epsilon^{-1}] \rceil$ be an integer for some constant c_1 . Then, with probability at least $1 - O(\delta)$, $\mathcal{E}_0, \{\mathcal{E}_i\}_{i=1}^R$ all happen.*

Proof. By a straightforward calculation, we have

$$\Pr[\cap_{i=0}^R \mathcal{E}_i] = \Pr[\mathcal{E}_0] \Pr[\mathcal{E}_1 | \mathcal{E}_0] \Pr[\mathcal{E}_2 | \mathcal{E}_0, \mathcal{E}_1] \dots \Pr[\mathcal{E}_R | \mathcal{E}_0, \mathcal{E}_1, \dots, \mathcal{E}_{R-1}] \geq 1 - O(\delta). \quad \square$$

Putting It Together

Proof of Proposition A.12. Let $\mathcal{E} = \cap_{i=0}^R \mathcal{E}_i$, we have shown, on \mathcal{E} , the outputs of Algorithm 1 and 2, $\{\mathbf{v}^{-(i)}, \mathbf{Q}^{-(i)}, \pi^{-(i)}, \boldsymbol{\xi}^{-(i)}\}_{i=0}^R$ is an MIVSS and $\{\mathbf{v}^{+(i)}, \mathbf{Q}^{+(i)}, \pi^{+(i)}, \boldsymbol{\xi}^{+(i)}\}_{i=0}^R$ is an MDVSS. By the above lemmas, \mathcal{E} happens with probability at least $1 - \Theta(\delta)$. This completes the proof of the proposition. \square

Proof of Theorem 1.1. The theorem is proved by combining Proposition A.3, A.4, and A.12. \square

Algorithm 2 QVI-MIVSS: algorithm for computing monotone increasing value-strategy sequences.

1: **Input:** A generative model for stochastic game $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathbf{r}, \mathbf{P}, \gamma)$;
 2: **Input:** Precision parameter $u \in [0, (1 - \gamma)^{-1}]$; and error probability $\delta \in (0, 1)$
 3: **Input:** Initial values $\mathbf{v}^{-(0)}, \pi^{-(0)}$ that satisfies monotonicity:

$$\mathbf{v}^* - u \mathbf{1} \leq \mathbf{v}^{-(0)} \leq \mathbf{v}^*, \quad \mathbf{v}^{-(0)} \leq \mathcal{T} \mathbf{v}^{-(0)}, \quad \text{and} \quad \mathbf{v}^{-(0)} \leq \mathcal{T}_{\pi^{-(0)}} \mathbf{v}^{-(0)};$$

4: **Output:** $\{\mathbf{v}^{-(i)}, \mathbf{Q}^{-(i)}, \pi^{-(i)}, \boldsymbol{\xi}^{-(i)}\}_{i=0}^R$ which is an MIVSS with high probability.
 5:
 6: **INITIALIZATION:**
 7: Let c_1, c_2, c_3, c be some tunable absolute constants;
 8: *Initialize constants:*
 9: $\beta \leftarrow (1 - \gamma)^{-1}$, and $R \leftarrow \lceil c_1 \beta \ln[\beta u^{-1}] \rceil$;
 10: $m_1 \leftarrow c_2 \beta^3 \cdot \min(1, u^{-2}) \cdot \log(8|\mathcal{S}||\mathcal{A}|\delta^{-1})$;
 11: $m_2 \leftarrow c_3 \beta^2 \log[2R|\mathcal{S}||\mathcal{A}|\delta^{-1}]$;
 12: $\alpha_1 \leftarrow L/m_1$ where $L = c \log(|\mathcal{S}||\mathcal{A}|\delta^{-1}(1 - \gamma)^{-1} \epsilon^{-1})$;
 13: *Obtain an initial batch of samples:*
 14: For each $(s, a) \in \mathcal{S} \times \mathcal{A}$: obtain independent samples $s_{s,a}^{(1)}, s_{s,a}^{(2)}, \dots, s_{s,a}^{(m_1)}$ from $\mathbf{P}(\cdot|s, a)$;
 15: Initialize: $\mathbf{w}^- = \tilde{\mathbf{w}}^- = \hat{\boldsymbol{\sigma}}^- = \mathbf{Q}^{-(0)} = \mathbf{Q}^{-(1)} \leftarrow \mathbf{0}_{\mathcal{S} \times \mathcal{A}}$ and $i \leftarrow 0$;
 16: **for** each $(s, a) \in \mathcal{S} \times \mathcal{A}$ **do**
 17: *Compute empirical estimates of $\mathbf{P}_{s,a}^\top \mathbf{v}^{-(0)}$ and $\text{var}(\mathbf{v}^{-(0)})(s, a)$:*
 18: $\tilde{\mathbf{w}}^-(s, a) \leftarrow \frac{1}{m_1} \sum_{j=1}^{m_1} \mathbf{v}^{-(0)}(s_{s,a}^{(j)})$;
 19: $\hat{\boldsymbol{\sigma}}^-(s, a) \leftarrow \frac{1}{m_1} \sum_{j=1}^{m_1} (\mathbf{v}^{-(0)}(s_{s,a}^{(j)}))^2 - (\tilde{\mathbf{w}}^-)^2(s, a)$;
 20: *Shift the empirical estimate to have one-sided error and guarantee monotonicity:*
 21: $\mathbf{w}^-(s, a) \leftarrow \tilde{\mathbf{w}}^-(s, a) - \sqrt{\alpha_1 \hat{\boldsymbol{\sigma}}^-(s, a)} - \alpha_1^{3/4} \beta$;
 22: *Compute coarse estimate of the Q-function and make sure its value is in $[0, \beta]$:*
 23: $\mathbf{Q}^{-(1)}(s, a) \leftarrow \text{clip}[\mathbf{r}(s, a) + \gamma \mathbf{w}^-(s, a), 0, \beta]$
 24: **end for**
 25:
 26: **REPEAT:** *successively improve*
 27: **for** $i = 1$ to R **do**
 28: *Compute the one-step dynamic programming:*
 29: Let $\mathbf{v}^{-(i)} \leftarrow \tilde{\mathbf{v}}^{-(i)} \leftarrow \mathcal{T} \mathbf{Q}^{-(i-1)}$, $\pi^{-(i)} \leftarrow \tilde{\pi}^{-(i)} \leftarrow \pi(\mathbf{Q}^{-(i-1)})$;
 30: *Compute strategy and value and maintain monotonicity:*
 31: For each $s \in \mathcal{S}$:
 32: if $\mathbf{v}^{-(i)}(s) \leq \mathbf{v}^{-(i-1)}(s)$, then $\mathbf{v}^{-(i)}(s) \leftarrow \mathbf{v}^{-(i-1)}(s)$ and $\pi^{-(i)}(s) \leftarrow \pi^{-(i-1)}(s)$;
 33: *Obtaining a small batch of samples:*
 34: For each $(s, a) \in \mathcal{S} \times \mathcal{A}$: draw independent samples $\tilde{s}_{s,a}^{(1)}, \tilde{s}_{s,a}^{(2)}, \dots, \tilde{s}_{s,a}^{(m_2)}$ from $\mathbf{P}(\cdot|s, a)$;
 35: *Compute the expected value, $\mathbf{g}^{\pm(i)}$, the estimate of $\mathbf{P}[\mathbf{v}^{\pm(i)} - \mathbf{v}^{\pm(0)}]$ with one-sided error:*
 36: Let $\tilde{\mathbf{g}}^{-(i)}(s, a) \leftarrow \frac{1}{m_2} \sum_{j=1}^{m_2} [\mathbf{v}^{-(i)}(\tilde{s}_{s,a}^{(j)}) - \mathbf{v}^{-(0)}(\tilde{s}_{s,a}^{(j)})]$;
 37: Let $\mathbf{g}^{-(i)}(s, a) \leftarrow \tilde{\mathbf{g}}^{-(i)}(s, a) - C(1 - \gamma)u$, where C is sufficiently small;
 38: *Estimate the approximation error:*
 39: $\boldsymbol{\xi}^{-(i)} \leftarrow 2\sqrt{\alpha_1 \boldsymbol{\sigma}_{\mathbf{v}^{-(0)}}} + 2[\alpha_1^{3/4} \beta + C(1 - \gamma)u] \cdot \mathbf{1}$
 40: *Improve $\mathbf{Q}^{-(i)}$ make sure its value is in $[0, \beta]$:*
 41: $\mathbf{Q}^{-(i+1)} \leftarrow \text{clip}[\mathbf{r} + \gamma \cdot [\mathbf{w}^- + \mathbf{g}^{-(i)}], 0, \beta]$;
 42: **end for**
 43: **return** $\{\mathbf{v}^{-(i)}, \mathbf{Q}^{-(i)}, \pi^{-(i)}, \boldsymbol{\xi}^{-(i)}\}_{i=0}^R$