

# Span-Agnostic Optimal Sample Complexity and Oracle Inequalities for Average-Reward RL

**Matthew Zurek** MATTHEW.ZUREK@WISC.EDU and **Yudong Chen** YUDONGCHEN@CS.WISC.EDU  
*Department of Computer Sciences, University of Wisconsin-Madison*

**Editors:** Nika Haghtalab and Ankur Moitra

## Abstract

We study the sample complexity of finding an  $\varepsilon$ -optimal policy in average-reward Markov Decision Processes (MDPs) with a generative model. The minimax optimal span-based complexity of  $\tilde{O}(SAH/\varepsilon^2)$ , where  $H$  is the span of the optimal bias function, has only been achievable with prior knowledge of the value of  $H$ . Prior-knowledge-free algorithms have been the objective of intensive research, but several natural approaches provably fail to achieve this goal. We resolve this problem, developing the first algorithms matching the optimal span-based complexity without  $H$  knowledge, both when the dataset size is fixed and when the suboptimality level  $\varepsilon$  is fixed. Our main technique combines the discounted reduction approach with a method for automatically tuning the effective horizon based on empirical confidence intervals or lower bounds on performance, which we term *horizon calibration*. We also develop an *empirical span penalization* approach, inspired by sample variance penalization, which satisfies an *oracle inequality* performance guarantee. In particular this algorithm can outperform the minimax complexity in benign settings such as when there exist near-optimal policies with span much smaller than  $H$ .

**Keywords:** Reinforcement learning, average-reward MDPs, sample complexity, adaptivity

## 1. Introduction

Reinforcement Learning (RL) has achieved significant empirical successes in various fields, demonstrating its potential to solve complex decision-making problems. RL is commonly modeled as to learn a policy which maximizes cumulative rewards within a Markov Decision Process (MDP), where the cumulative rewards can be measured in several different ways. We focus on the average-reward criterion, which involves the long-term average of collected rewards as the horizon goes to infinity, making it suitable for ongoing tasks without a natural endpoint.

A fundamental question in average-reward RL is the sample complexity for learning a near-optimal policy under a generative model of the MDP. This question has been the subject of intensive research. Recent work has established the minimax-optimal span-based complexity  $\tilde{O}(SA\|h^*\|_{\text{span}}/\varepsilon^2)$  for learning an  $\varepsilon$ -optimal policy (Zurek and Chen, 2025), where  $\|h^*\|_{\text{span}}$  denotes the span of the optimal bias  $h^*$  and is known to be a more refined complexity parameter than alternatives such as diameter or mixing times. However, this algorithm as well as earlier work all require prior knowledge of  $\|h^*\|_{\text{span}}$  (or other complexity parameters), which is generally unavailable, making the algorithms impractical. A flurry of subsequent research (Neu and Okolo, 2024; Tuynman et al., 2024; Jin et al., 2024; Zurek and Chen, 2024) has focused on removing the need for prior knowledge but failed to match the optimal span-based complexity. In fact, several natural approaches to knowledge-free optimal complexity, including span estimation and the average-reward plug-in method, are shown to provably fail (Tuynman et al., 2024; Zurek and Chen, 2024).

In this paper we resolve this problem, providing algorithms which obtain the optimal span-based complexity without knowing  $\|h^*\|_{\text{span}}$ , for both settings where we fix the dataset size  $n$  and where

we prescribe a target suboptimality level  $\varepsilon$ . Our algorithms are based upon reductions to discounted MDPs (DMDPs) combined with a novel technique of (effective-) *horizon calibration*, which chooses discount factors to maximize lower bounds or minimize confidence intervals on policy performance. This technique can be seen as related, but representing a simpler alternative, to the technique of sample variance penalization (SVP) from statistical learning (Maurer and Pontil, 2009; Duchi and Namkoong, 2019). We further develop an algorithm based more closely on a relaxed version of SVP, which we call *empirical span penalization*, which enjoys even stronger guarantees. In particular, this algorithm satisfies a complexity bound in terms of the minimum span  $\|h^\pi\|_{\text{span}}$  of any gain-optimal policy  $\pi$  in place of  $\|h^*\|_{\text{span}}$ . Moreover, it adapts to and competes with simpler (potentially suboptimal) policies with the best tradeoff between complexity and suboptimality. This bound is reminiscent of the *oracle inequalities* from the statistical learning literature (Deheuvels et al., 2007; Koltchinskii, 2011), but to the best of our knowledge is new to average-reward RL.

### 1.1. Related Work

The problem of learning optimal policies in average-reward MDPs (AMDPs) is studied in Jin and Sidford (2020, 2021); Li et al. (2024); Wang et al. (2022); Zhang and Xie (2023); Wang et al. (2023). We start with the recent work Zurek and Chen (2025), which was the first to obtain the optimal span-based sample complexity but required prior knowledge of  $\|h^*\|_{\text{span}}$  to do so. All earlier work also required knowledge of problem-dependent complexity parameters. More recent work, which we discuss below, has studied the setting without such prior knowledge; see Table 1 for a summary.

Tuynman et al. (2024) and Zurek and Chen (2025) show that it is generally impossible to obtain a multiplicative estimate of  $\|h^*\|_{\text{span}}$  with  $\text{poly}(SA\|h^*\|_{\text{span}})$  samples. See Appendix B for discussion of the relationship between our algorithms and estimating  $\|h^*\|_{\text{span}}$ . By estimating the MDP’s diameter  $D$ , which upper bounds  $\|h^*\|_{\text{span}}$  but can be arbitrarily larger (Bartlett and Tewari, 2009; Lattimore and Szepesvári, 2020), the work in Tuynman et al. (2024) removes the need for prior knowledge within the algorithm of Zurek and Chen (2025) but obtains a complexity involving  $D$  rather than  $\|h^*\|_{\text{span}}$ . The Q-learning-based algorithm in Jin et al. (2024) uses increasing discount factors and does not require prior knowledge. Their complexity bound however depends on the largest mixing time of all policies,  $\tau_{\text{unif}}$ , which satisfies  $3\tau_{\text{unif}} \geq \|h^*\|_{\text{span}}$  and can be infinite or arbitrarily larger than  $\|h^*\|_{\text{span}}$  (Wang et al., 2022; Zurek and Chen, 2024). (See Appendix A.2 for definitions of  $D$  and  $\tau_{\text{unif}}$ .)

Neu and Okolo (2024) and Zurek and Chen (2024) study, respectively, approaches based on stochastic saddle-point optimization and the average-reward plug-in method, both obtaining bounds involving the bias spans of certain policies output by the algorithm. These spans are not generally controlled by  $\|h^*\|_{\text{span}}$ ; in particular, Zurek and Chen (2024, Theorem 14) present an example where this is the case and show that the average-reward plug-in approach cannot achieve the optimal  $SA\|h^*\|_{\text{span}}/\varepsilon^2$  complexity. Zurek and Chen (2024) also analyzes a DMDP-reduction algorithm that uses a (relatively small) effective horizon independent of  $\|h^*\|_{\text{span}}$ , achieving a suboptimal complexity with  $\|h^*\|_{\text{span}}^2$  dependence.

While for consistency we present all results in Table 1 in terms of the sample complexity sufficient for  $\varepsilon$  suboptimality, we note that only the results of Zurek and Chen (2025), Tuynman et al. (2024), and our Theorem 2 are for the fixed- $\varepsilon$  setting (where a target suboptimality  $\varepsilon$  is provided to the algorithm); the others results are for the fixed- $n$  setting (where the dataset size is  $n$ ). We provide more discussion of the relationships between these two settings in Appendix A.3.

Also related to the present work are papers studying the sample complexity of the model-based/plug-in approach for discounted MDPs (Azar et al., 2012, 2013; Agarwal et al., 2020; Li et al., 2020; Zurek and Chen, 2024). We also note that Boone and Zhang (2024) recently developed an algorithm for the *online* setting achieving a  $\|h^*\|_{\text{span}}$ -based regret bound without requiring prior knowledge. This result does not imply any sample complexity bounds in our setting, because there is no general regret-to-PAC conversion for average-reward MDPs (Tuynman et al., 2024), and even if this were possible, their result appears to require  $\Omega(S^{40} A^{20} \|h^*\|_{\text{span}}^{10})$  interaction steps before achieving the optimal regret, which would imply a massive “burn-in” cost in our setting.

## 2. Problem Setup

A Markov decision process is a tuple  $(\mathcal{S}, \mathcal{A}, P, r)$ , where  $\mathcal{S}, \mathcal{A}$  are the state and action spaces, respectively, with finite cardinalities  $S := |\mathcal{S}|$  and  $A := |\mathcal{A}|$ ,  $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  is the transition kernel with  $\Delta(\mathcal{S})$  denoting the probability simplex on  $\mathcal{S}$ , and  $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  is the reward function. We only consider Markovian stationary policies of the form  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ . For initial state  $s_0 \in \mathcal{S}$  and policy  $\pi$ , let  $\mathbb{E}_{s_0}^\pi$  denote the expectation w.r.t. the distribution over trajectories  $(S_0, A_0, S_1, A_1, \dots)$  with  $S_0 = s_0$ ,  $A_t \sim \pi(S_t)$ , and  $S_{t+1} \sim P(\cdot | S_t, A_t)$ . Let  $P_\pi$  denote the transition probability matrix of the Markov chain induced by  $\pi$ , where  $(P_\pi)_{s,s'} := \sum_{a \in \mathcal{A}} \pi(a|s) P(s' | s, a)$ . Similarly let  $(r_\pi)_s := \sum_{a \in \mathcal{A}} \pi(a|s) r(s, a)$ . We also consider  $P$  as an  $(S \times A)$ -by- $S$  matrix with  $P_{sa,s'} = P(s' | s, a)$ , and  $r$  as an  $SA$ -dimensional vector. For a policy  $\pi$ , define the policy matrix  $M^\pi \in \mathbb{R}^{S \times SA}$  by  $M_{s,sa}^\pi = \pi(a|s)$  and  $M_{s,s'a}^\pi = 0$  if  $s \neq s'$ . Note that  $P_\pi = M^\pi P$  and  $r_\pi = M^\pi r$ . Also define the maximization operator  $M : \mathbb{R}^{SA} \rightarrow \mathbb{R}^S$  by  $M(x)_s = \max_a x_{sa}$ .

We assume  $P$  is unknown, but one has access to a generative model (a.k.a. simulator) (Kearns and Singh, 1998), which provides independent samples from  $P(\cdot | s, a)$  for each  $s \in \mathcal{S}, a \in \mathcal{A}$ . We assume  $r$  is known, which is standard (Agarwal et al., 2020; Li et al., 2020) as otherwise estimating  $r$  is relatively easy. Let  $\mathbf{0}, \mathbf{1} \in \mathbb{R}^S$  be the all-zero and all-one vectors, respectively.

**Discounted reward criterion** A discounted MDP is a tuple  $(\mathcal{S}, \mathcal{A}, P, r, \gamma)$ , where  $\gamma \in (0, 1)$  is the discount factor. For a policy  $\pi$ , the (discounted) value function  $V_\gamma^\pi : \mathcal{S} \rightarrow [0, \infty)$  is defined as  $V_\gamma^\pi(s) := \mathbb{E}_s^\pi [\sum_{t=0}^\infty \gamma^t R_t]$ , where  $R_t = r(S_t, A_t)$  is the reward received at time  $t$ . There always exists an optimal policy  $\pi_\gamma^*$  that satisfies  $V_{\gamma_\gamma^*}^\pi(s) = V_\gamma^*(s) := \sup_\pi V_\gamma^\pi(s), \forall s \in \mathcal{S}$  (Puterman, 1994). When using transition kernel  $\hat{P}$  we will accordingly write  $\hat{V}_\gamma^\pi$  for the associated value function. For reward functions  $r'$  other than  $r$ , we include the reward function in the subscript e.g.  $V_{\gamma,r'}^\pi$ .

**Average-reward criterion** In an MDP  $(\mathcal{S}, \mathcal{A}, P, r)$ , the average reward, a.k.a. the *gain*, of a policy  $\pi$  starting from state  $s$  is defined as  $\rho^\pi(s) := \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_s^\pi [\sum_{t=0}^{T-1} R_t]$ . The *bias function* of a stationary policy  $\pi$  is  $h^\pi(s) := \text{C-lim}_{T \rightarrow \infty} \mathbb{E}_s^\pi [\sum_{t=0}^{T-1} (R_t - \rho^\pi(S_t))]$ , where C-lim denotes the Cesaro limit. When the Markov chain induced by  $P_\pi$  is aperiodic, C-lim can be replaced with the usual limit. A policy  $\pi^*$  is *Blackwell-optimal* if there exists some discount factor  $\bar{\gamma} \in [0, 1)$  such that for all  $\gamma \geq \bar{\gamma}$  we have  $V_\gamma^{\pi^*} \geq V_\gamma^\pi, \forall \pi$ . When  $S$  and  $A$  are finite, there always exists some Blackwell-optimal policy, denoted by  $\pi^*$  (Puterman, 1994). Define the optimal gain  $\rho^* \in \mathbb{R}^S$  by  $\rho^*(s) = \sup_\pi \rho^\pi(s)$  and note that  $\rho^* := \rho^{\pi^*}$ . Define the optimal bias  $h^* := h^{\pi^*}$  (which is unique even when  $\pi^*$  is not). A policy  $\pi$  is *gain-optimal* if  $\rho^\pi = \rho^*$  and it is *bias-optimal* if in addition  $h^\pi = h^*$ . For  $x \in \mathbb{R}^S$ , define the span semi-norm  $\|x\|_{\text{span}} := \max_{s \in \mathcal{S}} x(s) - \min_{s \in \mathcal{S}} x(s)$ .

An MDP is communicating if for any states  $s$  and  $s'$ , some policy can reach  $s'$  from  $s$  with probability 1. An MDP is weakly communicating if the states can be partitioned into two subsets  $\mathcal{S} =$

Algorithm	Sample Complexity	Reference	Prior Knowledge
DMDP Reduction	$SA \frac{\ h^*\ _{\text{span}} + 1}{\varepsilon^2}$	Zurek and Chen (2025)	Yes
Diameter Estimation + DMDP Reduction	$SA \frac{D}{\varepsilon^2} + S^2 AD^2$	Tuynman et al. (2024)	No
Dynamic Horizon Q-Learning	$SA \frac{\tau_{\text{unif}}^8}{\varepsilon^8}$	Jin et al. (2024)	No
Stochastic Saddle-Point Optimization	$S^2 A^2 \frac{\ h^{\hat{\pi}}\ _{\text{span}}^4}{\varepsilon^2}$	Neu and Okolo (2024)	No
Plug-in Approach with Anchoring and Reward Perturbation	$SA \frac{\min\{D, \tau_{\text{unif}}\}}{\varepsilon^2} SA \frac{\ h^*\ _{\text{span}} + \min\{\ \hat{h}^*\ _{\text{span}}, \ \underline{h}^{\hat{\pi}}\ _{\text{span}}\}}{\varepsilon^2}$	Zurek and Chen (2024)	No
$\sqrt{n}$ -Horizon DMDP Reduction	$SA \frac{\ h^*\ _{\text{span}}^2 + 1}{\varepsilon^2}$	Zurek and Chen (2024)	No
DMDP Reduction + Horizon Calibration	$SA \frac{\ h^*\ _{\text{span}} + 1}{\varepsilon^2}$	Our Theorems 1 and 2	No
Span Penalization	$SA \inf_{\pi: \rho^\pi \text{ constant}} \left\{ \frac{\ h^\pi\ _{\text{span}}}{(\rho^\pi - \rho^* + \varepsilon)^2} \right\}$	Our Theorem 4	No

Table 1: **Algorithms and sample complexity bounds for average reward MDPs** for finding an  $\varepsilon$ -optimal policy under a generative model (up to log factors). See Appendix A.2 for the definitions of the complexity parameters  $D$ ,  $\tau_{\text{unif}}$ ,  $\|h^{\hat{\pi}}\|_{\text{span}}$ ,  $\|\hat{h}^*\|_{\text{span}}$ ,  $\|\underline{h}^{\hat{\pi}}\|_{\text{span}}$  used in prior work. Note that the diameter  $D$  and uniform mixing time  $3\tau_{\text{unif}}$  are both upper bounds of  $\|h^*\|_{\text{span}}$  and can be arbitrarily larger than  $\|h^*\|_{\text{span}}$ . The parameters  $\|h^{\hat{\pi}}\|_{\text{span}}$ ,  $\|\hat{h}^*\|_{\text{span}}$ ,  $\|\underline{h}^{\hat{\pi}}\|_{\text{span}}$  are not generally controlled by  $\|h^*\|_{\text{span}}$ . The guarantee for our Span Penalization algorithm involves the infimum over all policies  $\pi$  with constant (state-independent) gain  $\rho^\pi$ ; see Theorem 4 for an equivalent guarantee in terms of the dataset size.

$\mathcal{S}_1 \cup \mathcal{S}_2$  such that all states in  $\mathcal{S}_1$  are transient under all stationary policies and  $\mathcal{S}_2$  is communicating. In weakly communicating MDPs,  $\rho^*$  is a constant vector (all entries are equal). All results in this paper assume that  $P$  is weakly communicating. While not used in our results, the definitions of the MDP diameter  $D$  and uniform mixing time  $\tau_{\text{unif}}$  are given in Appendix A.2 for completeness. We use the standard  $\tilde{O}(\cdot)$  notation to hide logarithmic factors in  $S$ ,  $A$ ,  $\|h^*\|_{\text{span}}$ ,  $1/\varepsilon$ , and  $1/\delta$ .

### 3. Main Results

In this section, we present our algorithms and main results. Our algorithms involve the function  $\alpha(\delta, n) = 96\sqrt{\log(24SA n^5/\delta)} \log_2(\log_2(n+4))$ , which is  $\tilde{O}(1)$ ; see Remark 1 for its origin.

#### 3.1. Fixed- $n$ Setting

First we consider the setting where the number of samples per state-action pair,  $n$ , is fixed. Our objective is to learn a policy with the best possible rate of suboptimality. We refer to this as the *fixed- $n$  setting*. Our Algorithm 1 is based on using the dataset to form an empirical transition kernel  $\hat{P}$  and then computing a near-optimal policy in the DMDP  $(\hat{P}, r)$  for some discount factor  $\gamma$ . The key technique is a method for automatically calibrating  $\gamma$  (equivalently, the effective horizon  $\frac{1}{1-\gamma}$ ): we try multiple values of  $\gamma$ , and for each we compute a near-optimal policy  $\tilde{\pi}_\gamma$  for the DMDP  $(\hat{P}, r, \gamma)$  and a quantity  $\hat{L}(\gamma)$  that lower bounds its gain. We then use the discount factor  $\hat{\gamma}$  that optimizes this lower bound. For computational efficiency, we only need to try  $O(\log n)$  values of  $\gamma$ .

---

**Algorithm 1** Lower Bound Maximization via Horizon Calibration

---

**input:** Sample size per state-action pair  $n$

- 1: **for** each state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$  **do**
  - 2:     Collect  $n$  samples  $S_{s,a}^1, \dots, S_{s,a}^n$  from  $P(\cdot \mid s, a)$
  - 3:     Form the empirical transition kernel  $\hat{P}(s' \mid s, a) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{S_{s,a}^i = s'\}$ , for all  $s' \in \mathcal{S}$
  - 4: **end for**
  - 5: Form geometric discount factor range  $\mathcal{H} := \{\gamma : \text{there exists an integer } k \text{ such that } \sqrt{n} \leq \frac{1}{1-\gamma} = 2^k \leq n\}$
  - 6: **for** each discount factor  $\gamma \in \mathcal{H}$  **do**
  - 7:     Obtain policy  $\tilde{\pi}_\gamma$  and value function  $\tilde{V}_\gamma$  from  $\text{SOLVEDMDP}(\hat{P}, r, \gamma, \frac{1}{n})$
  - 8:     Compute objective value  $\hat{L}(\gamma) := (1-\gamma) \min_s \tilde{V}_\gamma(s) - 2\frac{1-\gamma}{n} - \alpha(\delta, n) \sqrt{\frac{\|\tilde{V}_\gamma\|_{\text{span}} + \frac{3}{n} + 1}{n}}$
  - 9: **end for**
  - 10: Find  $\hat{\gamma} = \operatorname{argmax}_{\gamma \in \mathcal{H}} \hat{L}(\gamma)$
  - 11: **return** policy  $\hat{\pi} := \tilde{\pi}_{\hat{\gamma}}$ , gain lower bound  $\hat{\rho} := \max\{\hat{L}(\hat{\gamma}), 0\} \mathbf{1}$
- 

**Theorem 1** Suppose  $P$  is weakly communicating. For some constant  $C_3$ , with probability at least  $1 - \delta$ , the policy  $\hat{\pi}$  and gain lower bound  $\hat{\rho}$  output by Algorithm 1 satisfy (elementwise)

$$\rho^{\hat{\pi}} \geq \hat{\rho} \geq \rho^* - C_3 \alpha(\delta, n)^2 \sqrt{\frac{\|h^*\|_{\text{span}} + 1}{n}} \mathbf{1}.$$

Theorem 1 shows that Algorithm 1 returns a policy  $\hat{\pi}$  with the minimax optimal rate of suboptimality without using any prior knowledge. The algorithm also returns a performance *certificate*  $\hat{\rho}$ , which lower-bounds  $\rho^{\hat{\pi}}$  (and  $\rho^*$ ), and this bound is tight up to an error of  $\tilde{O}(\sqrt{(\|h^*\|_{\text{span}} + 1)/n})$ .

We allow  $\text{SOLVEDMDP}$ , used in line 7 of Algorithm 1, to be any subroutine for approximately solving the empirical DMDP  $(\hat{P}, r, \gamma)$ , and simply require that it returns a deterministic policy  $\tilde{\pi}_\gamma$  and an approximate value function  $\tilde{V}_\gamma$  such that  $\hat{V}_\gamma^{\tilde{\pi}_\gamma} \geq \hat{V}_\gamma^* - \frac{1}{n} \mathbf{1}$  and  $\|\tilde{V}_\gamma - \hat{V}_\gamma^*\|_\infty \leq \frac{1}{n}$ . This can be done for instance using  $O(\frac{\log(n/(1-\gamma))}{1-\gamma})$  iterations of value iteration.

Algorithm 1 is reminiscent of sample variance penalization (SVP), a statistical learning algorithm which outputs a hypothesis minimizing the empirical risk plus an estimated variance term (Maurer and Pontil, 2009; Duchi and Namkoong, 2019). Here we clarify the connections and differences, which also provides intuitions for our algorithms. First, Algorithm 1 can be understood as controlling not only certain empirical variance (represented by the last term in the definition of  $\widehat{L}(\gamma)$  in line 8) but also a certain bias/approximation error due to discounted reduction. This is a significant difference since prior to our work it was not clear that such approximation error could be estimated/controlled without knowing  $\|h^*\|_{\text{span}}$ . See our proof sketch in Section 4 for more on this issue. Second, supposing that a lower bound like  $\rho^\pi \geq (1 - \gamma) \min_s \widehat{V}_\gamma^\pi(s) \mathbf{1} - \sqrt{\|\widehat{V}_\gamma^\pi\|_{\text{span}} + 1/n} \mathbf{1}$  holds for all policies  $\pi$  and all  $\gamma$  with high probability,<sup>1</sup> then an analogue of SVP would be (1) below, while Algorithm 1 can instead be seen as (approximately) solving (2).

$$\max_{\pi, \gamma} (1 - \gamma) \min_s \widehat{V}_\gamma^\pi(s) - \sqrt{\frac{\|\widehat{V}_\gamma^\pi\|_{\text{span}} + 1}{n}} \quad (1) \quad \max_\gamma (1 - \gamma) \min_s \widehat{V}_\gamma^*(s) - \sqrt{\frac{\|\widehat{V}_\gamma^*\|_{\text{span}} + 1}{n}} \quad (2)$$

Since  $\max_\pi \widehat{V}_\gamma^\pi = \widehat{V}_\gamma^*$ , solving (2) can be understood as only choosing  $\pi$  to optimize  $\widehat{V}_\gamma^\pi$  and then controlling the objective via  $\gamma$  tuning, whereas (1) optimizes all objective terms jointly. While (1) may appear more principled, it is not immediately clear how to solve such a problem, whereas (2) can simply utilize any DMDP solver for optimizing  $\widehat{V}_\gamma^\pi$  for fixed  $\gamma$ , then tune  $\gamma$  afterwards. However, this is not the final word on (1), as in Subsection 3.3 we develop our Algorithm 4 based on solving (1).

### 3.2. Fixed- $\varepsilon$ Setting

We next present our Algorithm 2 for the setting where one is given a target suboptimality  $\varepsilon$ , and the goal is to return a policy with suboptimality bounded by  $\varepsilon$  using as few samples as possible. We refer to this setting as the *fixed- $\varepsilon$  setting*. At a high level, we run our algorithm for the fixed- $n$  setting for a geometrically increasing sequence of dataset sizes  $\{n_i\}_i$ . However, the fixed- $\varepsilon$  setting is more challenging, because beyond lower-bounding the gain of some known policies, to obtain a termination condition we additionally need an (observable) upper bound on the optimal gain  $\rho^*$ , meaning we need to bound the gains of all policies. On iteration  $i$  with a dataset of size  $n_i$ , we compute both lower and upper bounds  $\widehat{L}_i(\gamma)$  and  $\widehat{U}_i(\gamma)$  for a range of  $\gamma$ , yielding different confidence intervals. The algorithm terminates once one such interval is sufficiently small and certifies the desired suboptimality level  $\varepsilon$ . We provide more details on how to compute such an upper bound  $\widehat{U}_i(\gamma)$  without knowing  $\|h^*\|_{\text{span}}$  in the proof sketches in Section 4.

**Theorem 2** *Suppose  $P$  is weakly communicating. There exist constants  $C_1, C_2$  such that for any  $\varepsilon > 0$ , with probability at least  $1 - \delta$ , Algorithm 2 uses at most*

$$N := 4C_1 \frac{\|h^*\|_{\text{span}} + 1}{\varepsilon^2} \log^3 \left( \frac{C_2 S A (\|h^*\|_{\text{span}} + 1)}{\delta \varepsilon} \right)$$

1. For a single policy  $\pi$  a similar bound follows from our techniques. Uniformity over all  $\pi$ 's would incur an additional  $\sqrt{S}$ .



**Algorithm 2** Confidence Interval Minimization via Horizon Calibration**input:** Target suboptimality  $\varepsilon > 0$ 

- 1: Set iteration number  $i = 0$
- 2: **repeat**
- 3:    $i \leftarrow i + 1$ ; set sample size per state-action pair  $n_i = 2^i$
- 4:   **for** each state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$  **do**
- 5:     Collect  $n_i$  samples  $S_{s,a}^1, \dots, S_{s,a}^{n_i}$  from  $P(\cdot \mid s, a)$
- 6:     Form the  $i$ th empirical transition kernel  $\hat{P}^{(i)}(s' \mid s, a) = \frac{1}{n} \sum_{j=1}^{n_i} \mathbb{I}\{S_{s,a}^j = s'\}, \forall s' \in \mathcal{S}$
- 7:   **end for**
- 8:   Form geometric discount factor range  $\mathcal{H}_i := \{\gamma : \text{there exists an integer } k \text{ such that } \sqrt{n_i} \leq \frac{1}{1-\gamma} = 2^k \leq n_i\}$
- 9:   **for** each discount factor  $\gamma \in \mathcal{H}_i$  **do**
- 10:     Obtain policy  $\tilde{\pi}_{\gamma,i}$  and value function  $\tilde{V}_{\gamma,i}$  from SOLVEDMDP( $\hat{P}^{(i)}, r, \gamma, \frac{1}{n_i}$ )
- 11:     Compute upper bound  $\hat{U}_i(\gamma) := (1 - \gamma) \max_s \tilde{V}_{\gamma,i}(s) + 5 \frac{1-\gamma}{n_i} + \frac{2\alpha(\delta, n_i)^2}{(1-\gamma)n_i} + 4\alpha(\delta, n_i) \sqrt{\frac{\|\tilde{V}_{\gamma,i}\|_{\text{span}} + 1 + \frac{3}{n_i}}{n_i}}$
- 12:     Compute lower bound  $\hat{L}_i(\gamma) := (1 - \gamma) \min_s \tilde{V}_{\gamma,i}(s) - 2 \frac{1-\gamma}{n_i} - \alpha(\delta, n_i) \sqrt{\frac{\|\tilde{V}_{\gamma,i}\|_{\text{span}} + \frac{3}{n_i} + 1}{n_i}}$
- 13:   **end for**
- 14:   Find discount factor with the smallest interval  $\hat{\gamma}_i := \operatorname{argmin}_{\gamma \in \mathcal{H}_i} \hat{U}_i(\gamma) - \hat{L}_i(\gamma)$
- 15: **until**  $\hat{U}_i(\hat{\gamma}_i) - \hat{L}_i(\hat{\gamma}_i) \leq \varepsilon$
- 16: **return** policy  $\hat{\pi} := \tilde{\pi}_{\hat{\gamma}_i,i}$ , optimal gain upper bound  $\hat{U} := \hat{U}_i(\hat{\gamma}_i)$  and lower bound  $\hat{L} := \hat{L}_i(\hat{\gamma}_i)$

samples per state-action pair and terminates after at most  $\log_2(N)$  outer iterations. Upon termination Algorithm 2 returns a policy  $\hat{\pi}$  and estimates  $\hat{U}, \hat{L}$  such that (elementwise)

$$\hat{L}\mathbf{1} \leq \rho^{\hat{\pi}} \leq \rho^* \leq \hat{U}\mathbf{1} \quad \text{and} \quad \hat{U} - \hat{L} \leq \varepsilon.$$

In particular, we have

$$\rho^{\hat{\pi}} \geq \rho^* - \varepsilon \mathbf{1}.$$

In fact, as shown in Lemma 9, the quantities  $\hat{U}_i(\gamma)$  are valid upper bounds of  $\rho^*$  for all  $i$  and  $\gamma \in \mathcal{H}_i$ , so the algorithm would still be correct if we instead used  $\min_{i, \gamma \in \mathcal{H}_i} \hat{U}_i(\gamma) - \max_{i', \gamma' \in \mathcal{H}_{i'}} \hat{L}_{i'}(\gamma') \leq \varepsilon$  as a termination condition; that is, we could use a different  $\gamma, i$  to compute the upper bound than the lower bound. The output policy should correspond to the best lower bound, that is, the  $\tilde{\pi}_{\gamma_i, i}$  such that  $i, \gamma \in \operatorname{argmax}_{i', \gamma' \in \mathcal{H}_{i'}} \hat{L}_{i'}(\gamma')$ . Also, within each outer-level iteration  $i$  of Algorithm 2 we consider  $O(\log n_i)$  values of  $\gamma$ , so by Theorem 2 the algorithm terminates after  $\tilde{O}(1)$  total iterations.

Compared to the algorithm of Zurek and Chen (2025), which also takes a discounted reduction approach but uses prior knowledge of  $\|h^*\|_{\text{span}}$  to set the discount factor, the  $\hat{\gamma}_i$  chosen by our algorithm should not be seen as implicitly estimating  $\|h^*\|_{\text{span}}$  and in fact cannot be used to do so (consistent with known hardness results on span estimation (Tuynman et al., 2024; Zurek and Chen, 2025)). Rather,  $\hat{\gamma}_i$  can be understood as balancing bounds on certain approximation and estimation error terms, potentially in a superior way than the  $\|h^*\|_{\text{span}}$ -knowledge-based choice

for non-worst-case instances. See Appendix B for further discussion of the relationship between  $\|h^*\|_{\text{span}}$  and the  $\hat{\gamma}_i$  computed in each iteration of our algorithm.

### 3.3. Span Penalization and Oracle Inequalities

Finally, we return to the fixed- $n$  setting and the goal of implementing the formulation (1) that resembles sample variance penalization. As we see momentarily, doing so allows us to obtain a stronger “oracle inequality” that optimally trades off suboptimality and complexity.

By superfluously introducing a span constraint of the form  $\|\hat{V}_\gamma^\pi\|_{\text{span}} \leq M$  to (1) and optimizing over  $M$ , we obtain the equivalent optimization problem

$$(1) \equiv \max_{\gamma, M} \max_{\pi: \|\hat{V}_\gamma^\pi\|_{\text{span}} \leq M} (1 - \gamma) \min_s \hat{V}_\gamma^\pi(s) - \sqrt{\frac{M+1}{n}}.$$

These manipulations are useful if, for fixed  $(\gamma, M)$ , we are able to solve the span-constrained optimization problem  $\max_{\pi: \|\hat{V}_\gamma^\pi\|_{\text{span}} \leq M} \hat{V}_\gamma^\pi$ . A natural approach is to attempt to apply value iteration with a span truncation step. This is inspired by Fruit et al. (2018), who first introduced this truncation operator and combined it with average-reward (undiscounted) value iteration to solve a certain bias-constrained gain optimization problem. Thus we define the span truncation operator  $\text{Clip}_M : \mathbb{R}^{\mathcal{S}} \rightarrow \mathbb{R}^{\mathcal{S}}$  where

$$\text{Clip}_M(V)(s) := \begin{cases} V(s) & \text{if } V(s) \leq M + \min_{s'} V(s') \\ M + \min_{s'} V(s') & \text{otherwise,} \end{cases} \quad (3)$$

or equivalently  $\text{Clip}_M(V) = \min\{V, (M + \min_{s'} V(s'))\mathbf{1}\}$ , where the outer min is elementwise. By combining  $\text{Clip}_M$ , which is  $\|\cdot\|_\infty$ -nonexpansive, with the discounted Bellman operator  $\mathcal{T}_\gamma(V) := M(r + \gamma PV)$ , we can now define our *Span-Constrained Planning* subroutine, given as Algorithm 3.

---

#### Algorithm 3 Span-Constrained Planning

---

**input:** Discounted MDP  $(P, r, \gamma)$ , span constraint bound  $M > 0$ , target error  $\varepsilon > 0$

- 1: Form clipped discounted Bellman operator  $\mathcal{L} := \text{Clip}_M \circ \mathcal{T}_\gamma$
  - 2: Set initial point  $V^0 = \mathbf{0} \in \mathbb{R}^{\mathcal{S}}$ , total iteration count  $T = \left\lceil \frac{\log(\frac{3}{(1-\gamma)^2\varepsilon})}{1-\gamma} \right\rceil$
  - 3: **for**  $t = 0, \dots, T-1$  **do**
  - 4:    $V^{t+1} = \mathcal{L}(V^t)$
  - 5: **end for**
  - 6: Set  $\hat{\pi}(s) \in \arg\max_{a \in \mathcal{A}} r(s, a) + \gamma P_{sa} V^T$  for all  $s \in \mathcal{S}$
  - 7: Compute minimum state value  $m = \min_{s \in \mathcal{S}} V^T(s)$
  - 8: Set truncated reward  $\tilde{r}(s, a) = \min\{m + M - \gamma P_{sa} V^T, r(s, a)\}$  for all  $s \in \mathcal{S}, a \in \mathcal{A}$
  - 9: **return** policy  $\hat{\pi}$ , approximate value function  $V^T$ , truncated reward  $\tilde{r}$
- 

Now we discuss why this subroutine returns a truncated reward  $\tilde{r}$  and whether it solves the aforementioned span-constrained planning problem  $\max_{\pi: \|\hat{V}_\gamma^\pi\|_{\text{span}} \leq M} \hat{V}_\gamma^\pi$ . (Here we discuss value functions based on a generic  $P$  rather than  $\hat{P}$  for clarity.) While the clipped Bellman operator  $\mathcal{L} = \text{Clip}_M \circ \mathcal{T}_\gamma$  has a unique fixed point  $V_{\gamma, M}^*$ , this fixed point generally may not satisfy any Bellman equation  $V_{\gamma, M}^* = M(r + \gamma PV_{\gamma, M}^*)$  for any policy  $\pi$ . This is to be expected, because such a



Bellman equation would imply that policy  $\pi$  has value function  $V_\gamma^\pi = V_{\gamma,M}^*$  and thus  $\|V_\gamma^\pi\|_{\text{span}} \leq M$ , but it is possible that no policies have value functions with  $\text{span} \leq M$ . (For example consider an MDP with  $A = 1$  where the only policy has a value function with large span.) The truncated reward  $\tilde{r}$  remedies both of these closely related issues: it is defined to (approximately) satisfy a Bellman equation  $V_{\gamma,M}^* = M^{\hat{\pi}}(\tilde{r} + \gamma P V_{\gamma,M}^*)$  for some  $\hat{\pi}$ , which would then imply that in the DMDP  $(P, \tilde{r}, \gamma)$  with the truncated reward, the value function  $V_{\gamma,\tilde{r}}^{\hat{\pi}}$  of policy  $\hat{\pi}$  does indeed have  $\|V_{\gamma,\tilde{r}}^{\hat{\pi}}\|_{\text{span}} \leq M$ . Thus, we do not actually solve  $\max_{\pi: \|V_\gamma^\pi\|_{\text{span}} \leq M} V_\gamma^\pi$ , but instead the relaxed problem  $\max_{\pi, \tilde{r}: \|V_{\gamma,\tilde{r}}^\pi\|_{\text{span}} \leq M, \tilde{r} \leq r} V_{\gamma,\tilde{r}}^\pi$ , which leads to a value function  $V_{\gamma,\tilde{r}}^{\hat{\pi}}$  that is at least as large as  $\max_{\pi: \|V_\gamma^\pi\|_{\text{span}} \leq M} V_\gamma^\pi$  yet still has span bounded by  $M$ . And, since elementwise  $\tilde{r} \leq r$ , the actual value function of  $\hat{\pi}$ ,  $V_\gamma^{\hat{\pi}}$ , is lower-bounded by  $V_{\gamma,\tilde{r}}^{\hat{\pi}}$ , which is still compatible with the lower-bound-based approach taken in Algorithm 4, which we are now prepared to develop.

The key properties of the Span-Constrained Planning Algorithm 3 are summarized in the following lemma.

**Lemma 3** *The operator  $\mathcal{L}$  defined in Algorithm 3 is  $\gamma$ -contractive and has unique fixed point  $V_{\gamma,M}^*$ . Moreover, the  $\hat{\pi}$ ,  $V^T$ , and  $\tilde{r}$  returned by Algorithm 3 satisfy*

1. (proximity to exact fixed point)  $\|V^T - V_{\gamma,M}^*\|_\infty \leq \varepsilon$ ;
2. (near-feasibility of  $\tilde{r}$  and  $\hat{\pi}$ )  $\tilde{r} \leq r$ ,  $\|V_{\gamma,\tilde{r}}^{\hat{\pi}} - V_{\gamma,M}^*\|_\infty \leq \varepsilon$ , and  $\|V_{\gamma,\tilde{r}}^{\hat{\pi}}\|_{\text{span}} \leq M + 2\varepsilon$ ;
3. (near-optimality of  $\hat{\pi}$ ) for any policy  $\pi'$  and reward function  $r' \leq r$  such that  $\|V_{\gamma,r'}^{\pi'}\|_{\text{span}} \leq M$ , we have  $V_{\gamma,M}^* \geq V_{\gamma,r'}^{\pi'}$  and  $V_\gamma^{\hat{\pi}} \geq V_{\gamma,\tilde{r}}^{\hat{\pi}} \geq V_{\gamma,M}^* - \varepsilon \mathbf{1} \geq V_{\gamma,r'}^{\pi'} - \varepsilon \mathbf{1}$ .

We can now present our final Algorithm 4, titled *Empirical Span Penalization*. In summary, it can be understood as solving the problem

$$\max_{\gamma, M} \max_{\pi, \tilde{r}: \|V_{\gamma,\tilde{r}}^\pi\|_{\text{span}} \leq M, \tilde{r} \leq r} \min_s \hat{V}_{\gamma,\tilde{r}}^\pi(s) - \sqrt{\frac{M+1}{n}}.$$

For reasons apparent in the proof sketches in Section 4, it is unclear whether there is any provable benefit of optimizing over all  $\gamma$  as opposed to choosing  $\gamma$  such that  $\frac{1}{1-\gamma} \approx \sqrt{nM}$ , so we opt to only optimize  $M$  and simply set  $\gamma$  in this way. Like for our previous algorithms, we only need to try  $O(\log n)$  values of  $M$  to approximately solve the problem.

**Theorem 4** *Suppose  $P$  is weakly communicating. For some constant  $C_4$ , with probability at least  $1 - \delta$ , the policy  $\hat{\pi}$  and gain lower bound  $\hat{\rho}$  output by Algorithm 4 satisfy (elementwise)*

$$\rho^{\hat{\pi}} \geq \hat{\rho} \geq \sup_{\pi: \rho^\pi \text{ constant}} \left\{ \rho^\pi - C_4 \alpha(\delta, n) \sqrt{\frac{\|h^\pi\|_{\text{span}} + 1}{n}} \mathbf{1} \right\}.$$

Theorem 4 shows that the output policy  $\hat{\pi}$  satisfies an “oracle inequality”:  $\hat{\pi}$  competes with the constant-gain policy  $\pi$  that has the best tradeoff between suboptimality and complexity, as measured by  $\rho^* - \rho^\pi$  and  $\|h^\pi\|_{\text{span}}$ , respectively. This bound, when written in the equivalent form

$$\rho^* - \rho^{\hat{\pi}} \leq \inf_{\pi: \rho^\pi \text{ constant}} \left\{ (\rho^* - \rho^\pi) + C_4 \alpha(\delta, n) \sqrt{\frac{\|h^\pi\|_{\text{span}} + 1}{n}} \mathbf{1} \right\},$$

**Algorithm 4** Empirical Span Penalization**input:** Sample size per state-action pair  $n$ 

- 
- 1: **for** each state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$  **do**
  - 2:     Collect  $n$  samples  $S_{s,a}^1, \dots, S_{s,a}^n$  from  $P(\cdot \mid s, a)$
  - 3:     Form the empirical transition kernel  $\hat{P}(s' \mid s, a) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{S_{s,a}^i = s'\}$ , for all  $s' \in \mathcal{S}$
  - 4: **end for**
  - 5: **for**  $i = 2, \dots, \lceil \log_2 n \rceil$  **do**
  - 6:     Set span constraint bound  $M_i = 2^i$
  - 7:     Set discount factor  $\gamma_i$  such that  $\frac{1}{1-\gamma_i} = \max \left\{ \frac{\sqrt{nM_i}}{\alpha(\delta, n)2\sqrt{2}}, 1 \right\}$
  - 8:     Obtain policy  $\tilde{\pi}_i$ , value function  $\tilde{V}_i$ , truncated reward  $\tilde{r}_i$  from the Span-Constrained Planning Algorithm 3 with input DMDP  $(\hat{P}, r, \gamma_i)$ , span constraint bound  $M_i$ , and target error  $\frac{1}{n}$
  - 9:     Compute objective value  $\hat{L}(i) := (1 - \gamma_i) \min_s \tilde{V}_i(s) - \alpha(\delta, n) \sqrt{\frac{M_i+1}{n}}$
  - 10: **end for**
  - 11: Find  $\hat{i} = \operatorname{argmax}_i \hat{L}(i)$
  - 12: **return** policy  $\hat{\pi} := \tilde{\pi}_{\hat{i}}$ , gain lower bound  $\hat{\rho} := \max\{\hat{L}(\hat{i}), 0\} \mathbf{1}$
- 

resembles the oracle inequalities (Deheuvels et al., 2007; Koltchinskii, 2011) from the statistics literature that feature a similar tradeoff.<sup>2</sup> It is immediate that Theorem 4 is at least as strong as Theorem 1 by setting  $\pi$  to be the Blackwell-optimal policy  $\pi^*$ , since  $\rho^{\pi^*} = \rho^*$  is constant and  $\|h^{\pi^*}\|_{\text{span}} = \|h^*\|_{\text{span}}$  by definition.

In fact, since there may generally be many gain-optimal  $\pi$ 's with  $h^\pi \neq h^*$ , Theorem 4 implies

$$\rho^{\hat{\pi}} \geq \rho^* - C_4 \alpha(\delta, n) \sqrt{\frac{\inf_{\pi: \rho^\pi = \rho^*} \|h^\pi\|_{\text{span}} + 1}{n}} \mathbf{1},$$

thus replacing  $\|h^*\|_{\text{span}}$  with the smallest bias of any gain-optimal policy. The above bound does not contradict the minimax lower bound rate of  $\sqrt{\|h^*\|_{\text{span}}/n}$  (Wang et al., 2022), which is based on worst-case instances, but we present an example in Appendix D where  $\inf_{\pi: \rho^\pi = \rho^*} \|h^\pi\|_{\text{span}}$  can be arbitrarily smaller than  $\|h^*\|_{\text{span}}$ . To the best of our knowledge, no other algorithm can achieve such a bound. Another situation where Theorem 4 outperforms the minimax rate is when all gain-optimal policies have large span relative to  $n$  but there is some near-optimal policy with much smaller span, such that learning a near-optimal policy is still possible. We provide such an example in Appendix D. In practice, RL is commonly applied to problems for which exact optimal policies are extremely complicated, and yet such problems may be solved to a reasonable degree of optimality by following relatively simple heuristics. Our result shows that it is possible to be automatically adaptive to the policy with the best tradeoff of complexity and suboptimality.

## 4. Proof Sketches

In this section, we outline the key ideas in the proofs of our main theorems.

---

2. An oracle inequality controls the error of a statistical estimator in terms of that of an oracle that selects the optimal model by trading off approximation error and estimation error.

#### 4.1. Proof sketch for Theorem 1

We start by outlining the ideas behind Algorithm 1 and Theorem 1 for the fixed- $n$  setting. First we review the approach of Zurek and Chen (2025) and simplifications due to Zurek and Chen (2024), which achieve the same rate as Theorem 1 but require knowing  $\|h^*\|_{\text{span}}$ . The algorithm of Zurek and Chen (2025) chooses a discount factor  $\gamma^*$  in a way that depends crucially on  $\|h^*\|_{\text{span}}$ , and then solves the  $\gamma^*$ -discounted empirical MDP  $(\hat{P}, r, \gamma^*)$ , where  $\hat{P}$  is constructed as in Algorithm 1. The resulting policy  $\hat{\pi}_{\gamma^*}$  is shown to be near-optimal for the AMDP  $(P, r)$  using the reduction from Wang et al. (2022). We now illustrate how this optimally-chosen  $\gamma^*$  trades off certain approximation and learning error terms. Letting  $\hat{\pi}_\gamma$  be the policy output by the above-described procedure but with a general discount factor  $\gamma$ , it is shown that

$$\begin{aligned} \left\| \rho^{\hat{\pi}_\gamma} - \rho^* \right\|_\infty &\stackrel{(i)}{\lesssim} (1 - \gamma) \left( \|h^*\|_{\text{span}} + \|V_\gamma^{\hat{\pi}_\gamma} - V_\gamma^*\|_\infty \right) \\ &\stackrel{(ii)}{\lesssim} (1 - \gamma) \left( \|h^*\|_{\text{span}} + \frac{1}{1 - \gamma} \sqrt{\frac{\|V_\gamma^*\|_{\text{span}} + \|\hat{V}_\gamma^*\|_{\text{span}} + 1}{n}} \right) \\ &\stackrel{(iii)}{\lesssim} (1 - \gamma) \|h^*\|_{\text{span}} + \sqrt{\frac{\|h^*\|_{\text{span}} + \|\hat{V}_\gamma^*\|_{\text{span}} + 1}{n}}, \end{aligned} \quad (4)$$

where the notation  $\lesssim$  ignores constant/log factors, (i) is the AMDP-to-DMDP reduction of Wang et al. (2022), (ii) uses Zurek and Chen (2024, Theorem 9) to bound  $\|V_\gamma^{\hat{\pi}_\gamma} - V_\gamma^*\|_\infty$  in terms of the variance parameters  $\|V_\gamma^*\|_{\text{span}}$ ,  $\|\hat{V}_\gamma^*\|_{\text{span}}$ , and (iii) uses  $\|V_\gamma^*\|_{\text{span}} \lesssim \|h^*\|_{\text{span}}$  due to Wei et al. (2020, Lemma 2). To proceed, one can invoke Zurek and Chen (2024, Lemma 12) (itself summarizing a key step due to Zurek and Chen 2025 similar to our Lemma 7), which controls the variance parameter  $\|\hat{V}_\gamma^*\|_{\text{span}}$  associated with the empirically optimal policy  $\hat{\pi}$  like

$$\|\hat{V}_\gamma^*\|_{\text{span}} \lesssim \|V_\gamma^*\|_{\text{span}} + \left\| \hat{V}_\gamma^* - V_\gamma^* \right\|_\infty \lesssim \|h^*\|_{\text{span}} + \frac{1}{1 - \gamma} \sqrt{\frac{\|h^*\|_{\text{span}} + \|\hat{V}_\gamma^*\|_{\text{span}} + 1}{n}}. \quad (5)$$

Solving this recursion yields  $\|\hat{V}_\gamma^*\|_{\text{span}} \lesssim \|h^*\|_{\text{span}} + \frac{1}{(1 - \gamma)^2 n}$ . Plugging back into (4) gives the following bound and in particular the term  $T_3$ :

$$\left\| \rho^{\hat{\pi}_\gamma} - \rho^* \right\|_\infty \lesssim \underbrace{(1 - \gamma) \|h^*\|_{\text{span}}}_{T_1} + \underbrace{\sqrt{\frac{\|h^*\|_{\text{span}} + 1}{n}}}_{T_2} + \underbrace{\frac{1}{(1 - \gamma)n}}_{T_3} \quad (6)$$

Intuitively, this argument is reminiscent of localization techniques in statistical learning (e.g., Bartlett et al. 2005) since it controls the variance parameters associated with near-optimal policies in terms of the variance of an optimal policy. For this reason we call  $T_3$  the *localization error*. Term  $T_1$  is the approximation error due to the DMDP reduction, and term  $T_2$  is independent of  $\gamma$  and is exactly the desired minimax rate. From (6) it is immediate that choosing  $\gamma^*$  so that  $\frac{1}{1 - \gamma^*} \approx \sqrt{n(\|h^*\|_{\text{span}} + 1)}$  will balance all three terms and achieve the minimax rate.

Now we derive an algorithm which does not require prior knowledge of  $\|h^*\|_{\text{span}}$ . The simplest idea is to attempt to estimate  $\|h^*\|_{\text{span}}$  and plug the estimate into the formula for  $\gamma^*$ . However,

estimating  $\|h^*\|_{\text{span}}$  up to a constant factor is shown to be impossible in the worst case (Zurek and Chen, 2025; Tuynman et al., 2024). Moreover, Zurek and Chen (2024, Theorem 14) provides an example where choosing  $\gamma$  too large provably fails to obtain the minimax-optimal rate, suggesting that the localization error term  $T_3$  could not be removed with a different analysis. A partially successful approach is to attempt to replace the terms in (6) with *observable/estimable* upper bounds that depend on  $\gamma$ , because observability would enable us to minimize these upper bounds over  $\gamma$ . Doing a similar “localization” approach as in (5), but instead upper-bounding  $\|V_\gamma^*\|_{\text{span}}$  in terms of  $\|\hat{V}_\gamma^*\|_{\text{span}}$ , one obtains  $\|V_\gamma^*\|_{\text{span}} \lesssim \|\hat{V}_\gamma^*\|_{\text{span}} + \frac{1}{(1-\gamma)^2 n}$ . Using this bound but otherwise following the derivation of (6), we can get

$$\|\rho^{\hat{\pi}_\gamma} - \rho^*\|_\infty \lesssim \underbrace{(1-\gamma)\|h^*\|_{\text{span}}}_{T_1} + \underbrace{\sqrt{\frac{\|\hat{V}_\gamma^*\|_{\text{span}} + 1}{n}}}_{T'_2} + \underbrace{\frac{1}{(1-\gamma)n}}_{T_3}, \quad (7)$$

where now both  $T_3$  and  $T'_2$  are observable, and based on (5), there are some choices of  $\gamma$  (including  $\gamma^*$ ) for which  $T'_2$  will not be too much larger than  $T_2$ . However,  $T_1$  still depends on the unknown  $\|h^*\|_{\text{span}}$ . This approach can actually be salvaged with one key modification and forms the basis of our result (Theorem 2), but here we take a slightly different approach to address term  $T_1$ . Since  $T_1$  is decreasing monotonically in  $\gamma$ , we believe *Lepski’s trick* (Lepski and Spokoiny, 1997) could actually be applied to estimate  $\rho^*$  with the minimax rate, but obtaining a policy with a matching order of suboptimality appears to require more work.

In particular, instead of optimizing a bound on the suboptimality  $\|\rho^{\hat{\pi}_\gamma} - \rho^*\|_\infty$  of the policy  $\hat{\pi}_\gamma$  that is optimal for the empirical DMDP  $(\hat{P}, r, \gamma)$ , we optimize a lower bound on the performance of  $\hat{\pi}_\gamma$  since we are considering the fixed- $n$  setting. For any  $\pi$  and any  $\gamma$ , it is known that  $\rho^\pi \geq (1-\gamma)V_\gamma^\pi - (1-\gamma)\|V_\gamma^\pi\|_{\text{span}}\mathbf{1}$ .<sup>3</sup> The significance of this statement applied to our problem is that if we only desire a lower bound for  $\rho^{\hat{\pi}_\gamma}$ , then the approximation error term  $T_1$  can be replaced with  $(1-\gamma)\|V_\gamma^{\hat{\pi}_\gamma}\|_{\text{span}}$ , which can be estimated. Specifically, using the bound  $\|V_\gamma^{\hat{\pi}_\gamma} - \hat{V}_\gamma^*\|_\infty \lesssim \frac{1}{1-\gamma} \sqrt{\frac{\|\hat{V}_\gamma^*\|_{\text{span}} + 1}{n}}$ , we have

$$\begin{aligned} \rho^{\hat{\pi}_\gamma} &\geq (1-\gamma)V_\gamma^{\hat{\pi}_\gamma} - (1-\gamma)\|V_\gamma^{\hat{\pi}_\gamma}\|_{\text{span}}\mathbf{1} \geq (1-\gamma)\hat{V}_\gamma^* - (1-\gamma)\left(\|\hat{V}_\gamma^*\|_{\text{span}} - C\|V_\gamma^{\hat{\pi}_\gamma} - \hat{V}_\gamma^*\|_\infty\right)\mathbf{1} \\ &\geq (1-\gamma)\hat{V}_\gamma^* - (1-\gamma)\|\hat{V}_\gamma^*\|_{\text{span}}\mathbf{1} - C'\sqrt{\frac{\|\hat{V}_\gamma^*\|_{\text{span}} + 1}{n}}\mathbf{1}, \end{aligned} \quad (8)$$

where  $C, C' \leq \tilde{O}(1)$ . Let  $L(\gamma)$  be the minimum entry of the RHS of (8). Then  $L(\gamma)$  is observable and we can show  $L(\gamma^*) \geq \rho^* - \sqrt{\frac{\|h^*\|_{\text{span}} + 1}{n}}$ , so returning the policy  $\pi_{\hat{\gamma}}$  where  $\hat{\gamma} = \arg\max_\gamma L(\gamma)$ , we obtain the minimax rate. (No explicit localization error term appears in (8), but lower-bounding  $L(\gamma)$  requires the “localization bound”  $\|\hat{V}_\gamma^*\|_{\text{span}} \lesssim \|h^*\|_{\text{span}} + \frac{1}{(1-\gamma)^2 n}$ , so it still appears implicitly in the size of  $\|\hat{V}_\gamma^*\|_{\text{span}}$ .) This is essentially our Theorem 1.

## 4.2. Proof sketch for Theorem 2

Next we describe the ideas behind Algorithm 2 and Theorem 2 for the fixed- $\varepsilon$  setting. It suffices to develop a method of bounding the suboptimality  $\|\rho^{\hat{\pi}} - \rho^*\|_\infty$  within the fixed- $n$  setting, as we

3. We actually refine this slightly to  $\rho^\pi \geq (1-\gamma) \min_s V_\gamma^\pi(s)\mathbf{1}$ , which is what appears in our algorithms.

can then double the dataset size until the suboptimality bound is  $\leq \varepsilon$ . We return to (7) and the problematic term  $T_1$ , which originates from the AMDP-to-DMDP reduction of Wang et al. (2022). We show that actually  $\|h^*\|_{\text{span}}$  can be replaced with  $\|V_\gamma^*\|_{\text{span}}$ , that is,

$$\|\rho^{\hat{\pi}_\gamma} - \rho^*\|_\infty \lesssim (1 - \gamma) \left( \|V_\gamma^*\|_{\text{span}} + \|V_\gamma^{\hat{\pi}_\gamma} - V_\gamma^*\|_\infty \right). \quad (9)$$

See Lemma 5 for the precise statement. We could essentially recover the reduction result of Wang et al. (2022) (step (i) of (6)) by using  $\|V_\gamma^*\|_{\text{span}} \lesssim \|h^*\|_{\text{span}}$ , but the subtle improvement of (9) gives a new approximation error term  $(1 - \gamma)\|V_\gamma^*\|_{\text{span}}$  which, with aforementioned arguments, can be bounded by fully observable quantities involving  $\hat{V}_\gamma^*$ . Specifically, (7) can be replaced with

$$\|\rho^{\hat{\pi}_\gamma} - \rho^*\|_\infty \lesssim (1 - \gamma)\|\hat{V}_\gamma^*\|_{\text{span}} + \sqrt{\frac{\|\hat{V}_\gamma^*\|_{\text{span}} + 1}{n_i}} + \frac{1}{(1 - \gamma)n_i}, \quad (10)$$

where we now use a variable dataset size  $n_i$ , since we intend to double the size until a termination condition is satisfied. Comparing with our Algorithm 2, the RHS of (10) is essentially the difference between the upper and lower bounds defined in the algorithm, that is the quantity  $\hat{U}_i(\gamma) - \hat{L}_i(\gamma)$ . Thus Algorithm 2 is slightly more useful than the procedure sketched here, since it provides valid upper and lower bounds rather than just the difference between such bounds, but this does not require much more work. Similarly to the previous algorithm, when  $\gamma = \gamma^*$ , the bound (10) will be  $\lesssim \sqrt{(\|h^*\|_{\text{span}} + 1)/n}$  thanks to the “localization bound”  $\|\hat{V}_\gamma^*\|_{\text{span}} \lesssim \|h^*\|_{\text{span}} + \frac{1}{(1 - \gamma)^2 n}$  arising from (5). This concludes our motivation and proof sketch of Theorem 2. Also see Appendix B for more discussion on the relationship between  $\|h^*\|_{\text{span}}$  and the  $\hat{\gamma}_i$  that minimizes the right hand side of (10) over  $\gamma$ .

### 4.3. Proof sketch of Theorem 4

For Theorem 4, our starting point is a lower bound similar to (8) used within Algorithm 1, but with a few key differences. Suppose we apply the Span-Constrained Planning Algorithm 3 to the empirical DMDP  $(\hat{P}, r, \gamma)$  with span constraint  $M$  and, for simplicity, an arbitrarily small target error, to get the policy  $\hat{\pi}_M$ , truncated reward  $\tilde{r}$ , and the exact fixed point  $\hat{V}_{\gamma, M}^* = \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}_M}$ . We also fix a comparator policy  $\pi$  such that  $\|h^\pi\|_{\text{span}} + 1 \leq \frac{M}{4}$  and  $\rho^\pi$  is constant, which enable us to show that  $\|V_\gamma^\pi\|_{\text{span}} \lesssim \|h^\pi\|_{\text{span}}$  for any  $\gamma$  and that  $\|\hat{V}_\gamma^\pi - V_\gamma^\pi\|_\infty \lesssim \sqrt{M/(1 - \gamma)^2 n}$  using standard techniques (Zurek and Chen, 2024). By setting  $\gamma$  so that  $\frac{1}{1 - \gamma} \lesssim \sqrt{nM}$ , we can ensure that  $\|\hat{V}_\gamma^\pi\|_{\text{span}} \leq M$  (by bounding it in terms of  $\|V_\gamma^\pi\|_{\text{span}} \lesssim \|h^\pi\|_{\text{span}} \lesssim M$  and  $\|\hat{V}_\gamma^\pi - V_\gamma^\pi\|_\infty$ ), which is essential to ensure that  $\pi$  is “feasible” for the empirical span-constrained problem, guaranteeing that  $\hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}_M} \geq \hat{V}_\gamma^\pi$  by Lemma 3. These assumptions and condition on  $\gamma$  also imply that  $\hat{V}_\gamma^\pi \geq \frac{1}{1 - \gamma}\rho^\pi - O(M)$ . Combining

all these steps we can show that

$$\begin{aligned}
\rho^{\hat{\pi}_M} &\geq (1 - \gamma) \min_s V_{\gamma}^{\hat{\pi}_M}(s) \mathbf{1} \stackrel{\text{(I)}}{\geq} (1 - \gamma) \min_s V_{\gamma, \tilde{r}}^{\hat{\pi}_M}(s) \mathbf{1} \\
&\geq (1 - \gamma) \min_s \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}_M}(s) \mathbf{1} - (1 - \gamma) \left\| \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}_M} - V_{\gamma, \tilde{r}}^{\hat{\pi}_M} \right\|_{\infty} \mathbf{1} \\
&\stackrel{\text{(II)}}{\geq} (1 - \gamma) \min_s \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}_M}(s) \mathbf{1} - C \sqrt{\frac{M+1}{n}} \stackrel{\text{(III)}}{\geq} (1 - \gamma) \min_s \hat{V}_{\gamma}^{\pi}(s) \mathbf{1} - C \sqrt{\frac{M+1}{n}} \\
&\stackrel{\text{(IV)}}{\geq} \rho^{\pi} - (1 - \gamma) C' M \mathbf{1} - C \sqrt{\frac{M+1}{n}},
\end{aligned}$$

where  $C, C' \leq \tilde{O}(1)$ . Here (I) follows from  $\tilde{r} \leq r$ , and (II) follows from the error bound

$$\left\| \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}_M} - V_{\gamma, \tilde{r}}^{\hat{\pi}_M} \right\|_{\infty} \lesssim \frac{1}{1 - \gamma} \sqrt{\frac{\|\hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}_M}\|_{\text{span}} + 1}{n}} = \frac{1}{1 - \gamma} \sqrt{\frac{M+1}{n}}, \quad (11)$$

and (III) and (IV) follow from the aforementioned “feasibility” of  $\pi$  and the lower bound on  $\hat{V}_{\gamma}^{\pi}$ . The key new challenge is to develop the error bound (11), which is complicated due to the statistical dependence between  $\hat{P}$  and  $\hat{\pi}_M$ . Without the span constraint this is addressed by Agarwal et al. (2020); Li et al. (2020) through the use of “absorbing MDP” constructions, which enable concentration inequalities for  $|(\hat{P}_{sa} - P_{sa})\hat{V}_{\gamma}^{\star}|$ . Instead, we desire such bounds for  $|(\hat{P}_{sa} - P_{sa})\hat{V}_{\gamma, M}^{\star}|$ , that is, involving the empirical span-constrained optimal value functions. Based on the contractivity of the span-constrained Bellman operator  $\mathcal{L}$ , we develop a new absorbing MDP construction for span-constrained value functions, ultimately leading to the desired bound (11).

## 5. Conclusion

We resolve the problem of achieving the minimax optimal span-based complexity in average-reward RL without prior knowledge of the span or other unknown MDP complexity parameters. Our algorithms apply to both the fixed-sample-size and fixed-suboptimality settings, and moreover achieve optimal tradeoff between complexity and suboptimality, surpassing the minimax lower bound in benign settings. Future directions of interest include generalizing to general/multichain MDPs and extending beyond the tabular simulator setting.

## Acknowledgments

Y. Chen and M. Zurek acknowledge support by National Science Foundation grants CCF-2233152 and DMS-2023239.

## References

Alekh Agarwal, Sham Kakade, and Lin F Yang. Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pages 67–83. PMLR, 2020.



- Mohammad Gheshlaghi Azar, Rémi Munos, and Hilbert J Kappen. On the sample complexity of reinforcement learning with a generative model. In *Proceedings of the 29th International Conference on Machine Learning*, pages 1707–1714, 2012.
- Mohammad Gheshlaghi Azar, Rémi Munos, and Hilbert J. Kappen. Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model. *Machine Learning*, 91(3):325–349, June 2013. ISSN 1573-0565. doi: 10.1007/s10994-013-5368-1. URL <https://doi.org/10.1007/s10994-013-5368-1>.
- Peter Bartlett and Ambuj Tewari. REGAL: a regularization based algorithm for reinforcement learning in weakly communicating mdps. In *Uncertainty in Artificial Intelligence: Proceedings of the 25th Conference*, pages 35–42. AUAI Press, 2009.
- Peter L. Bartlett, Olivier Bousquet, and Shahar Mendelson. Local Rademacher complexities. *The Annals of Statistics*, 33(4):1497–1537, August 2005. ISSN 0090-5364, 2168-8966. doi: 10.1214/009053605000000282. URL <https://projecteuclid.org/journals/annals-of-statistics/volume-33/issue-4/Local-Rademacher-complexities/10.1214/009053605000000282.full>. Publisher: Institute of Mathematical Statistics.
- Victor Boone and Zihan Zhang. Achieving Tractable Minimax Optimal Regret in Average Reward MDPs, June 2024. URL <http://arxiv.org/abs/2406.01234>. arXiv:2406.01234 [cs].
- Paul Deheuvels, Eustasio Del Barrio, and Sara Van De Geer. *Lectures on Empirical Processes: Theory and Statistical Applications*, volume 6 of *EMS Series of Lectures in Mathematics*. EMS Press, 1 edition, January 2007. ISBN 978-3-03719-027-2 978-3-03719-527-7. doi: 10.4171/027. URL <https://ems.press/doi/10.4171/027>.
- John Duchi and Hongseok Namkoong. Variance-based regularization with convex objectives. *Journal of Machine Learning Research*, 20(1):2450–2504, January 2019. ISSN 1532-4435.
- Ronan Fruit, Matteo Pirota, Alessandro Lazaric, and Ronald Ortner. Efficient bias-span-constrained exploration-exploitation in reinforcement learning. In *International Conference on Machine Learning*, pages 1578–1586. PMLR, 2018.
- Ying Jin, Ramki Gummadi, Zhengyuan Zhou, and Jose Blanchet. Feasible \$Q\$-Learning for Average Reward Reinforcement Learning. In *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, pages 1630–1638. PMLR, April 2024. URL <https://proceedings.mlr.press/v238/jin24b.html>. ISSN: 2640-3498.
- Yujia Jin and Aaron Sidford. Efficiently solving MDPs with stochastic mirror descent. In *International Conference on Machine Learning*, pages 4890–4900. PMLR, 2020.
- Yujia Jin and Aaron Sidford. Towards tight bounds on the sample complexity of average-reward MDPs. In *International Conference on Machine Learning*, pages 5055–5064. PMLR, 2021.
- Michael Kearns and Satinder Singh. Finite-Sample Convergence Rates for Q-Learning and Indirect Algorithms. In *Advances in Neural Information Processing Systems*, volume 11. MIT Press, 1998. URL <https://proceedings.neurips.cc/paper/1998/hash/99adff456950dd9629a5260c4de21858-Abstract.html>.

- Vladimir Koltchinskii. *Oracle Inequalities in Empirical Risk Minimization and Sparse Recovery Problems: École d'Été de Probabilités de Saint-Flour XXXVIII-2008*, volume 2033 of *Lecture Notes in Mathematics*. Springer Berlin Heidelberg, Berlin, Heidelberg, 2011. ISBN 978-3-642-22146-0 978-3-642-22147-7. doi: 10.1007/978-3-642-22147-7. URL <https://link.springer.com/10.1007/978-3-642-22147-7>.
- Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, Cambridge ; New York, NY, 2020. ISBN 978-1-108-57140-1.
- Oleg V. Lepski and Vladimir G. Spokoiny. Optimal pointwise adaptive methods in nonparametric estimation. *The Annals of Statistics*, 25(6):2512–2546, 1997.
- Gen Li, Yuting Wei, Yuejie Chi, Yuantao Gu, and Yuxin Chen. Breaking the Sample Size Barrier in Model-Based Reinforcement Learning with a Generative Model. In *Advances in Neural Information Processing Systems*, volume 33, pages 12861–12872. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/96ea64f3a1aa2fd00c72faacf0cb8ac9-Abstract.html>.
- Tianjiao Li, Feiyang Wu, and Guanghui Lan. Stochastic first-order methods for average-reward markov decision processes. *Mathematics of Operations Research*, 2024.
- Andreas Maurer and Massimiliano Pontil. Empirical Bernstein Bounds and Sample Variance Penalization, July 2009. URL <http://arxiv.org/abs/0907.3740>. Proc. Computational Learning Theory Conference (COLT 2009). arXiv:0907.3740 [stat] version: 1.
- Gergely Neu and Nneka Okolo. Dealing with unbounded gradients in stochastic saddle-point optimization, June 2024. URL <http://arxiv.org/abs/2402.13903>. Forty-first International Conference on Machine Learning. arXiv:2402.13903 [cs, math, stat] version: 2.
- Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley Series in Probability and Statistics. Wiley, 1 edition, April 1994. ISBN 978-0-471-61977-2 978-0-470-31688-7. doi: 10.1002/9780470316887. URL <https://onlinelibrary.wiley.com/doi/book/10.1002/9780470316887>.
- Adrienne Tuynman, Rémy Degenne, and Emilie Kaufmann. Finding good policies in average-reward Markov Decision Processes without prior knowledge, May 2024. URL <http://arxiv.org/abs/2405.17108>. arXiv:2405.17108 [cs].
- Jinghan Wang, Mengdi Wang, and Lin F. Yang. Near Sample-Optimal Reduction-based Policy Learning for Average Reward MDP, December 2022. URL <http://arxiv.org/abs/2212.00603>. arXiv:2212.00603 [cs].
- Shengbo Wang, Jose Blanchet, and Peter Glynn. Optimal Sample Complexity of Reinforcement Learning for Mixing Discounted Markov Decision Processes, September 2023. URL <http://arxiv.org/abs/2302.07477>. arXiv:2302.07477.
- Chen-Yu Wei, Mehdi Jafarnia Jahromi, Haipeng Luo, Hiteshi Sharma, and Rahul Jain. Model-free reinforcement learning in infinite-horizon average-reward markov decision processes. In *International conference on machine learning*, pages 10170–10180. PMLR, 2020.

Zihan Zhang and Qiaomin Xie. Sharper Model-free Reinforcement Learning for Average-reward Markov Decision Processes, June 2023. URL <http://arxiv.org/abs/2306.16394>. The Thirty Sixth Annual Conference on Learning Theory. arXiv:2306.16394 [cs].

Matthew Zurek and Yudong Chen. The Plug-in Approach for Average-Reward and Discounted MDPs: Optimal Sample Complexity Analysis, October 2024. URL <http://arxiv.org/abs/2410.07616>. International Conference on Algorithmic Learning Theory (ALT 2025). arXiv:2410.07616 [cs].

Matthew Zurek and Yudong Chen. Span-Based Optimal Sample Complexity for Weakly Communicating and General Average Reward MDPs. *Advances in Neural Information Processing Systems*, 37:33455–33504, January 2025. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/hash/3acbe9dc3ale8d48a57b16e9aef91879-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2024/hash/3acbe9dc3ale8d48a57b16e9aef91879-Abstract-Conference.html).

## Appendix A. Additional Notation and Guide to Appendices

In this section, we provide additional notations and definitions, and outline the organization of the remainder of the appendices.

### A.1. Additional Notation

in the proofs of our main theorems, we make use of some additional notations and standard facts about average-reward MDPs. For any policy  $\pi$ , its gain and bias  $\rho^\pi$  and  $h^\pi$  satisfy the optimality equations  $\rho^\pi = P_\pi \rho^\pi$  and  $\rho^\pi + h^\pi = r_\pi + P_\pi h^\pi$ , where the second equation is sometimes called the (average-reward) Bellman equation or Poisson equation. We let  $P_\pi^\infty = \text{C-lim}_{T \rightarrow \infty} (P_\pi)^T$  denote the limiting matrix of the Markov chain induced by the policy  $\pi$ . When this Markov chain is aperiodic, the Cesaro limit, C-lim, can be replaced with the usual limit. Note that  $P_\pi^\infty P_\pi = P_\pi P_\pi^\infty = P_\pi^\infty$  and  $\rho^\pi = P_\pi^\infty r_\pi$ . For any policy  $\pi$  we define the Bellman consistency operator  $\mathcal{T}_\gamma^\pi : \mathbb{R}^S \rightarrow \mathbb{R}^S$  by  $\mathcal{T}_\gamma^\pi(x) = r_\pi + \gamma P_\pi x$ .

For any  $x \in \mathbb{R}^S$  and any policy  $\pi$ , we define a policy-specific next-state transition variance vector  $\mathbb{V}_{P_\pi}[x] \in \mathbb{R}^S$  as  $(\mathbb{V}_{P_\pi}[x])_s := \sum_{s' \in S} (P_\pi)_{s,s'} [x(s') - \sum_{s''} (P_\pi)_{s,s''} x(s'')]^2$ . We use  $\|B\|_{\infty \rightarrow \infty}$  to denote the  $\|\cdot\|_\infty$  operator norm of matrix  $B$ , which is equal to the maximum  $\ell^1$ -norms of any row of  $B$ .

### A.2. Complexity Parameters

In this section, we provide definitions for the diameter  $D$  and the uniform mixing time  $\tau_{\text{unif}}$ , which are standard complexity parameters (along with optimal bias span  $\|h^*\|_{\text{span}}$ ) used in prior work on average-reward MDPs. We also give the definitions for the quantities  $\|\hat{h}^\pi\|_{\text{span}}$ ,  $\|\hat{\underline{h}}^\star\|_{\text{span}}$ , and  $\|\underline{h}^\pi\|_{\text{span}}$  appearing in Table 1, which arise in the results of [Neu and Okolo \(2024\)](#) and [Zurek and Chen \(2024\)](#). The results in the present paper do not involve these parameters; their definitions are included here for completeness and for comparison with our results.

First we define the diameter. For any state  $s \in \mathcal{S}$ , let  $\eta_s = \inf\{t \geq 1 : S_t = s\}$  denote the hitting time of state  $s$ , which is a random variable (in the probability space of trajectories in the MDP  $P$ ). Then the diameter  $D$  is defined as

$$D := \max_{s_1 \neq s_2} \inf_{\pi \in \Pi_{\text{MD}}} \mathbb{E}_{s_1}^\pi [\eta_{s_2}],$$

where  $\Pi_{\text{MD}}$  is the set of all Markovian deterministic policies.  $D < \infty$  if and only if the MDP is communicating, so in particular it is generally infinite in weakly communicating MDPs.

Now we define the uniform mixing time  $\tau_{\text{unif}}$ . This parameter is only defined if we first assume that for all  $\pi \in \Pi_{\text{MD}}$ , the Markov chain induced by  $P_\pi$  has a unique stationary distribution, which we call  $\nu_\pi$  (considered as a row vector in  $\mathbb{R}^S$ ). Then for any policy  $\pi \in \Pi_{\text{MD}}$ , we can define its mixing time  $\tau_\pi := \inf\{t \geq 1 : \max_{s \in S} \|e_s^\top (P_\pi)^t - \nu_\pi^\top\|_1 \leq \frac{1}{2}\}$ . Finally, we define the uniform mixing time  $\tau_{\text{unif}} := \sup_{\pi \in \Pi_{\text{MD}}} \tau_\pi$ .

It always holds that  $\|h^*\|_{\text{span}} \leq D$  ([Bartlett and Tewari, 2009](#)) and  $\|h^*\|_{\text{span}} \leq 3\tau_{\text{unif}}$  ([Wang et al., 2022](#); [Zurek and Chen, 2024](#)). In general,  $D$  and  $\tau_{\text{unif}}$  are not comparable ([Wang et al., 2022](#)), and they both can be arbitrarily larger than  $\|h^*\|_{\text{span}}$ .

The prior work listed in Table 1 also involves the quantities  $\|\hat{h}^\pi\|_{\text{span}}$ ,  $\|\hat{\underline{h}}^\star\|_{\text{span}}$  and  $\|\underline{h}^\pi\|_{\text{span}}$ . These are all dependent on the outputs of certain (randomized) algorithms which introduce them, and

they are generally not controlled in terms of  $\|h^*\|_{\text{span}}$ . In particular, the algorithm described in (Neu and Okolo, 2024, Theorem 4.1) returns a policy  $\hat{\pi}$ , and their complexity bounds depend on  $\|h^{\hat{\pi}}\|_{\text{span}}$ , the span of the bias of  $\hat{\pi}$  under the true MDP  $(P, r)$ . Zurek and Chen (2024) define an *anchored* empirical AMDP  $(\hat{P}, r)$  with the transition matrix  $\hat{P} := (1 - \eta)\hat{P} + \eta\mathbf{1}e_{s_0}^\top$ , where  $\eta \in [0, 1]$  is an algorithmic parameter,  $s_0$  is an arbitrary state and  $e_{s_0} \in \mathbb{R}^S$  is the vector with all zeros except for the entry  $s_0$  being equal to 1. Letting  $\hat{h}^*$  denote the optimal bias function in this anchored AMDP, and letting  $\hat{h}^{\hat{\pi}}$  denote the bias function of the policy  $\hat{\pi}$  output by (the anchored and unperturbed version of) Algorithm 1 in Zurek and Chen (2024), their sample complexity guarantees are given in terms of the quantities  $\|\hat{h}^*\|_{\text{span}}$  and  $\|\hat{h}^{\hat{\pi}}\|_{\text{span}}$ . Note that Zurek and Chen (2024, Theorem 14) gives an example where both of these terms are much larger than  $\|h^*\|_{\text{span}}$ .

We refer to Neu and Okolo (2024) and Zurek and Chen (2024) for the complete definitions of  $\|h^{\hat{\pi}}\|_{\text{span}}$ ,  $\|\hat{h}^*\|_{\text{span}}$  and  $\|\hat{h}^{\hat{\pi}}\|_{\text{span}}$  as well as their relationship with other complexity parameters.

### A.3. Fixed- $n$ vs. Fixed- $\varepsilon$ Settings

Here we compare the fixed- $n$  and fixed- $\varepsilon$  settings. Suppose we have some algorithm for the fixed- $n$  setting which outputs a policy  $\hat{\pi}$ . If the complexity parameters appearing in the suboptimality guarantee for this algorithm (e.g.,  $C\sqrt{\|h^*\|_{\text{span}} \log(1/\delta)/n}$ ) are known, then the algorithm can be converted to the fixed- $\varepsilon$  setting by “inverting” this bound to find  $n$  large enough to make it  $\leq \varepsilon$ . For the minimax rate of solving discounted MDPs of  $\tilde{O}(\sqrt{\frac{1}{(1-\gamma)^3 n}})$  (Azar et al., 2013), since  $\gamma$  is known, this is always possible. However, this is not generally possible in the average-reward setting, since the bounds involve unknown complexity parameters such as the optimal bias span  $\|h^*\|_{\text{span}}$ , the uniform mixing time  $\tau_{\text{unif}}$ , or the diameter  $D$ . A natural remedy is to attempt to estimate these parameters and then use the estimate to invert the suboptimality bound, which is the approach taken by Tuynman et al. (2024) to obtain an algorithm for the fixed- $\varepsilon$  setting which has a  $D$ -based sample complexity, which is possible since  $D$  is estimable. However, the optimal bias span is not estimable (Tuynman et al., 2024; Zurek and Chen, 2025) and the uniform mixing time does not seem to be efficiently computable, so this approach does not seem to work. But even if the “inversion” cannot be performed algorithmically, we can still use it from an analytical perspective to compare algorithms, as we do in Table 1. We note then that some algorithms from Table 1 that are prior-knowledge-free for the fixed- $n$  setting cannot be converted to the fixed- $\varepsilon$  setting without complexity parameter knowledge. (The goal of removing prior knowledge is still nontrivial for the fixed- $n$  setting; the algorithm of Zurek and Chen (2025) requires prior knowledge in the fixed- $n$  setting in order to tune the discount factor used within a discounted reduction.) Finally, we briefly note that while lower bounds for average-reward RL have been developed for the fixed- $\varepsilon$  setting (e.g., Wang et al. 2022), it is straightforward to convert them to the fixed- $n$  setting.

### A.4. Remark on $\alpha$

We discuss the origin of the function  $\alpha$  that appears in our algorithms and bounds.

**Remark 1** *The quantity  $\alpha(\delta, n) = 96\sqrt{\log(24SA n^5/\delta)} \log_2(\log_2(n+4))$  appearing in our bounds arises from our choice to make use of concentration bounds developed by Zurek and Chen (2024). This quantity is used within our algorithms, so we make it explicit but did not attempt to optimize it. We note Zurek and Chen (2024) also did not optimize constants/log factors, and thus we*

conjecture that a smaller function  $\alpha$  may be sufficient for their inequalities to hold, in which case the improvement could be carried over to our work in a black-box manner.

### A.5. Guide to Appendices

In Appendix B, we discuss the relationship between the  $\hat{\gamma}$  selected by our algorithm and the true optimal bias span  $\|h^*\|_{\text{span}}$ . Appendix C contains the proofs for all of our main theorems, and Appendix D contains examples mentioned in Subsection 3.3 of situations where the guarantee of Theorem 4 could be much better than the minimax rate. Within Appendix C, we first prove Theorem 2 in Appendix C.1, then Theorem 1 in Appendix C.2. In Appendix C.3 we prove Lemma 3 regarding the properties of the span-constrained planning subroutine Algorithm 3, and finally in Appendix C.4 we prove Theorem 4.

### Appendix B. Relationship Between $\hat{\gamma}_i$ and $\|h^*\|_{\text{span}}$

Recall equation (10) from the proof sketch of Theorem 2. There, we see that the discount factor  $\hat{\gamma}_i$  selected by Algorithm 2 in a given iteration  $i$  is (approximately) the minimizer of a function  $\hat{B}(\gamma)$  that upper-bounds the suboptimality of a policy  $\hat{\pi}_\gamma$  in the following manner:

$$\|\rho^{\hat{\pi}_\gamma} - \rho^*\|_\infty \lesssim \hat{B}(\gamma) := (1 - \gamma)\|\hat{V}_\gamma^*\|_{\text{span}} + \sqrt{\frac{\|\hat{V}_\gamma^*\|_{\text{span}} + 1}{n_i}} + \frac{1}{(1 - \gamma)n_i}.$$

We can understand  $\hat{B}(\gamma)$  by relating it to the following (deterministic) quantity

$$B(\gamma) := (1 - \gamma)\|V_\gamma^*\|_{\text{span}} + \sqrt{\frac{\|V_\gamma^*\|_{\text{span}} + 1}{n_i}} + \frac{1}{(1 - \gamma)n_i}. \quad (12)$$

Using an error bound of the form  $\|V_\gamma^* - \hat{V}_\gamma^*\|_\infty \lesssim \frac{1}{1 - \gamma} \sqrt{\frac{\|V_\gamma^*\|_{\text{span}} + \|\hat{V}_\gamma^*\|_{\text{span}} + 1}{n_i}}$  (Zurek and Chen, 2024, Theorem 9), as well as the “localization” bound  $\|V_\gamma^*\|_{\text{span}} \lesssim \|\hat{V}_\gamma^*\|_{\text{span}} + \frac{1}{(1 - \gamma)^2 n_i}$  which appears in the proof sketch, we can show that (with high probability)  $B(\gamma) \lesssim \hat{B}(\gamma)$ , and likewise repeating the same steps but with the roles of  $\hat{V}_\gamma^*$  and  $V_\gamma^*$  reversed we can also show that  $\hat{B}(\gamma) \lesssim B(\gamma)$ . Therefore, since  $\hat{B}(\gamma)$  and  $B(\gamma)$  are equivalent up to logarithmic factors with high probability, we can understand  $\hat{\gamma}_i$  by considering the  $\gamma$  which minimizes  $B(\gamma)$ .

We now turn to the “oracle choice” of  $\gamma^*$  made in Zurek and Chen (2024) (also discussed in our proof sketches) such that  $\frac{1}{1 - \gamma^*} \approx \sqrt{n_i(\|h^*\|_{\text{span}} + 1)}$ . Using the fact that  $\|V_\gamma^*\|_{\text{span}} \lesssim \|h^*\|_{\text{span}}$  (Wei et al., 2020), we are guaranteed the final term in the definition (12) of  $B(\gamma^*)$  is the largest and equal to  $\sqrt{\frac{\|h^*\|_{\text{span}} + 1}{n_i}}$ . This implies that the  $\hat{\gamma}_i$  selected by minimizing  $\hat{B}(\gamma)$  (and also the minimizing  $\gamma$  of  $B(\gamma)$ ) must satisfy  $\frac{1}{1 - \hat{\gamma}_i} \lesssim \frac{1}{1 - \gamma^*} \approx \sqrt{n_i(\|h^*\|_{\text{span}} + 1)}$ , since for all  $\gamma \geq \gamma^*$  the final term of  $B(\gamma)$  is  $\gtrsim \frac{1}{(1 - \gamma^*)n_i} \gtrsim B(\gamma^*)$ . However,  $\frac{1}{1 - \hat{\gamma}_i}$  may potentially be much smaller than  $\sqrt{n_i(\|h^*\|_{\text{span}} + 1)}$ , so in particular it is not generally possible to convert  $\hat{\gamma}_i$  into some constant-factor estimate of  $\|h^*\|_{\text{span}}$  (which would contradict the hardness of estimating  $\|h^*\|_{\text{span}}$  as established in Tuynman et al. 2024; Zurek and Chen 2025). Indeed, it is possible to compute the minimizer  $\tilde{\gamma}^*$  of  $B(\gamma)$  for the MDP



instances that were used to show the hardness of estimating  $\|h^*\|_{\text{span}}$  in Zurek and Chen (2025, Proof of Theorem 3). We find that such instances (which have  $n_i$  as an input parameter) have  $\|V_\gamma^*\|_{\text{span}} \lesssim 1 + \frac{1}{1-\gamma} \frac{1}{\sqrt{n_i}}$ , and for relatively small  $\gamma$  this bound can be much better than the bound  $\|V_\gamma^*\|_{\text{span}} \lesssim \|h^*\|_{\text{span}}$ ; consequently, our sharper DMDP reduction (14), which replaces  $\|h^*\|_{\text{span}}$  with  $\|V_\gamma^*\|_{\text{span}}$ , gives a much better bound (also see (9) in the proof sketches). This fact about  $\|V_\gamma^*\|_{\text{span}}$  for these instances also implies that the minimizing  $\tilde{\gamma}^*$  has  $\frac{1}{1-\tilde{\gamma}^*} \approx \sqrt{n_i} \ll \sqrt{n_i(\|h^*\|_{\text{span}} + 1)}$ . Since one of the instances used in Zurek and Chen (2025, Proof of Theorem 3) has an arbitrarily large bias span  $\|h^*\|_{\text{span}}$ , this in particular implies that  $\frac{1}{1-\tilde{\gamma}^*}$  (and thus also  $\frac{1}{1-\tilde{\gamma}_i}$ ) may be arbitrarily smaller than  $\sqrt{n_i(\|h^*\|_{\text{span}} + 1)}$ . This minimizing choice  $\tilde{\gamma}^*$  gives  $B(\tilde{\gamma}^*) \lesssim \frac{1}{\sqrt{n_i}} \ll \sqrt{\frac{\|h^*\|_{\text{span}}}{n_i}}$ , which implies that these instances, while being worst-case for estimating the optimal bias span, are *not* worst-case for the task of finding a gain-optimal policy or estimating the optimal gain.

In summary, although  $\hat{\gamma}_i$  is chosen in our algorithm to calibrate certain approximation and estimation error terms at least as well as the “oracle choice”  $\gamma^*$  (which depends on  $\|h^*\|_{\text{span}}$ ), this  $\hat{\gamma}_i$  should not be seen as estimating  $\|h^*\|_{\text{span}}$ , but rather as calibrating tighter bounds on approximation/estimation errors; in particular, as discussed in the proof sketches, these tighter bounds replace  $\|h^*\|_{\text{span}}$  with  $\|V_\gamma^*\|_{\text{span}}$ . In fact, the “oracle choice”  $\frac{1}{1-\gamma^*} \approx \sqrt{n_i(\|h^*\|_{\text{span}} + 1)}$  is not only more difficult to compute, but can generally achieve a worse tradeoff than the  $\tilde{\gamma}^*$  which minimizes  $B(\gamma)$ . Still, for worst-case instances such as those appearing in the minimax lower bounds for learning a gain-optimal policy or estimating the optimal gain, the oracle  $\gamma^*$  must indeed approximately minimize  $B(\gamma)$ , since we know for such instances that  $B(\gamma) \gtrsim \sqrt{\|h^*\|_{\text{span}}/n_i}$  for all  $\gamma$  (since they are hard instances), and also that  $B(\gamma^*) \lesssim \sqrt{\|h^*\|_{\text{span}}/n_i}$ . This observation suggests that for instances that are hard for optimizing/estimating the gain, we may actually be able to estimate  $\|h^*\|_{\text{span}}$  to a constant factor using  $\hat{\gamma}_i$ , which would imply such instances are *not* hard for the task of span estimation.

## Appendix C. Proofs of Main Theorems

In this section, we prove our main theorems.

### C.1. Proof of Theorem 2

The proof has two main steps. First we will show that all upper and lower bounds computed within Algorithm 2 are valid upper/lower bounds on  $\rho^*$  and the gain of the output policy. This implies that whenever the algorithm terminates, it will have output a policy which is  $\varepsilon$ -optimal. The second step is to show that once a sufficient number of iterations have occurred (and thus a sufficiently large sample size is chosen), the confidence interval corresponding to some discount factor  $\gamma^*$  will be sufficiently small, thus ensuring termination of the algorithm at or before this point.

The following lemma is very similar to results of Wei et al. (2020) and Wang et al. (2022). While the bound  $\|\rho^* - (1-\gamma)V_\gamma^*\|_\infty \leq (1-\gamma)\|h^*\|_{\text{span}}$  has appeared previously (Wei et al., 2020, Lemma 2), to the best of our knowledge a bound of the form  $\|\rho^* - (1-\gamma)V_\gamma^*\|_\infty \leq (1-\gamma)\|V_\gamma^*\|_{\text{span}}$  has not previously appeared in the literature. As explained in the proof sketches in Section 4, the difference is extremely significant in our situation, since  $\|V_\gamma^*\|_{\text{span}}$  can be estimated while  $\|h^*\|_{\text{span}}$

generally cannot be (Zurek and Chen, 2025; Tuynman et al., 2024). This new result thus provides a computable upper bound for the “approximation error” due to reducing the AMDP to a DMDP with a certain discount factor, which can be used algorithmically. We also note that (16) below is already known and appears in (Wang et al., 2022, Lemma 6). (Their proof has a minor issue because it is not the case that all policies  $\pi$  have  $\rho^\pi$  which is state-independent/constant in a weakly communicating MDP, but their proof does not actually need this fact and the result is still true.)

**Lemma 5** *Fix a discount factor  $\gamma \in [0, 1)$ . Then*

$$(1 - \gamma) \left( \min_s V_\gamma^*(s) \right) \mathbf{1} \leq \rho^* \leq (1 - \gamma) \left( \max_s V_\gamma^*(s) \right) \mathbf{1} \quad (13)$$

*which implies*

$$\|\rho^* - (1 - \gamma)V_\gamma^*\|_\infty \leq (1 - \gamma) \|V_\gamma^*\|_{\text{span}}. \quad (14)$$

*Additionally, for any fixed policy  $\pi$ , we have that*

$$(1 - \gamma) \left( \min_s V_\gamma^\pi(s) \right) \mathbf{1} \leq \rho^\pi \leq (1 - \gamma) \left( \max_s V_\gamma^\pi(s) \right) \mathbf{1}, \quad (15)$$

*which implies*

$$\|\rho^\pi - (1 - \gamma)V_\gamma^\pi\|_\infty \leq (1 - \gamma) \|V_\gamma^\pi\|_{\text{span}}. \quad (16)$$

*Consequently, for any policy  $\pi$ , we have that*

$$\rho^\pi \geq \rho^* - (1 - \gamma) \left( \|V_\gamma^\pi - V_\gamma^*\|_\infty + \|V_\gamma^*\|_{\text{span}} \right) \mathbf{1}.$$

**Proof** First, for the statement with a fixed policy  $\pi$ , as shown in Wang et al. (2022) we have that

$$\rho^\pi = P_\pi^\infty r_\pi = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P_\pi^\infty r_\pi = (1 - \gamma) P_\pi^\infty \sum_{t=0}^{\infty} \gamma^t P_\pi^t r_\pi = (1 - \gamma) P_\pi^\infty V_\gamma^\pi$$

which implies that, for each  $s \in \mathcal{S}$ ,  $\rho^\pi(s) = e_s^\top P_\pi^\infty ((1 - \gamma)V_\gamma^\pi)$ , and thus  $\rho^\pi(s)$  is a convex combination of entries of  $(1 - \gamma)V_\gamma^\pi$ , implying that  $(1 - \gamma) \min_{s'} V_\gamma^\pi(s') \leq \rho^\pi(s) \leq (1 - \gamma) \max_{s'} V_\gamma^\pi(s')$ . By subtracting  $(1 - \gamma)V_\gamma^\pi(s)$  we obtain that

$$(1 - \gamma) \left( \min_{s'} V_\gamma^\pi(s') - V_\gamma^\pi(s) \right) \leq \rho^\pi(s) - (1 - \gamma)V_\gamma^\pi(s) \leq (1 - \gamma) \left( \max_{s'} V_\gamma^\pi(s') - V_\gamma^\pi(s) \right)$$

which can be further lower and upper bounded as

$$-(1 - \gamma) \|V_\gamma^\pi\|_{\text{span}} \leq \rho^\pi(s) - (1 - \gamma)V_\gamma^\pi(s) \leq (1 - \gamma) \|V_\gamma^\pi\|_{\text{span}}$$

implying (16).

Now we show the statement (13). For any  $s \in \mathcal{S}$ , we similarly have

$$\rho^*(s) = (1 - \gamma) e_s^\top P_{\pi^*}^\infty V_\gamma^* \leq (1 - \gamma) e_s^\top P_{\pi^*}^\infty V_\gamma^* \leq (1 - \gamma) \max_s V_\gamma^*(s)$$

and also for any  $s \in \mathcal{S}$

$$\rho^*(s) \geq \rho^{\pi^*}(s) = (1 - \gamma) e_s^\top P_{\pi^*}^\infty V_\gamma^* \geq (1 - \gamma) \min_s V_\gamma^*(s).$$

Statement (14) follows in an identical manner as to how (16) followed from (15).

For the final desired statement, using the previous results we have

$$\begin{aligned} \rho^\pi &\geq (1 - \gamma) \min_s V_\gamma^\pi(s) \mathbf{1} \geq (1 - \gamma) \min_s V_\gamma^*(s) \mathbf{1} - (1 - \gamma) \|V_\gamma^\pi - V_\gamma^*\|_\infty \mathbf{1} \\ &\geq (1 - \gamma) \max_s V_\gamma^*(s) \mathbf{1} - (1 - \gamma) \|V_\gamma^*\|_{\text{span}} \mathbf{1} - (1 - \gamma) \|V_\gamma^\pi - V_\gamma^*\|_\infty \mathbf{1} \\ &\geq \rho^* - (1 - \gamma) \|V_\gamma^*\|_{\text{span}} \mathbf{1} - (1 - \gamma) \|V_\gamma^\pi - V_\gamma^*\|_\infty \mathbf{1}. \end{aligned}$$

■

Here we repeat the definitions of some quantities which are defined in Algorithm 2, and additionally we extend their definitions to hold for all integers  $i \geq 1$  (rather than just the iterations which appear in the course of the algorithm before its termination). For all  $i \geq 1$ , we define  $\hat{P}^{(i)}$  to be constructed in the manner shown in Algorithm 2 via taking  $n_i = 2^i$  samples from all state-action pairs. We also let  $\mathcal{H}_i := \{\gamma : \text{there exists an integer } k \text{ such that } \sqrt{n_i} \leq \frac{1}{1-\gamma} = 2^k \leq n_i\}$ . We can then define, for all  $\gamma \in \mathcal{H}_i$ , the policy  $\tilde{\pi}_{\gamma,i}$  and value function  $\tilde{V}_{\gamma,i}$  as the outputs from SOLVEDMDP( $\hat{P}^{(i)}, r, \gamma, \frac{1}{n_i}$ ). With such quantities we can then define  $\hat{U}_i(\gamma), \hat{L}_i(\gamma)$  following the definitions given in 2, and finally we let  $\hat{\gamma}_i := \operatorname{argmin}_{\gamma \in \mathcal{H}_i} \hat{U}_i(\gamma) - \hat{L}_i(\gamma)$ . We also use  $\hat{V}_{\gamma,i}^\pi$  to denote the value function of policy  $\pi$  in the DMDP ( $\hat{P}^{(i)}, r, \gamma$ ) and  $\hat{V}_{\gamma,i}^*$  to denote the optimal value function of this DMDP.

**Lemma 6** Define the function  $\alpha(\delta, \tilde{n}) = 24\sqrt{16 \log\left(\frac{24SA\tilde{n}^5}{\delta}\right) \log_2(\log_2(\tilde{n} + 4))}$ . Fix  $\delta > 0$ . Then with probability at least  $1 - \delta$ , we have for all integers  $i \geq 1$  and all  $\gamma \in \mathcal{H}_i$  that

$$\|V_\gamma^{\pi^*} - \hat{V}_{\gamma,i}^{\pi^*}\|_\infty \leq \frac{\alpha(\delta, n_i)}{1 - \gamma} \sqrt{\frac{\|V_\gamma^{\pi^*}\|_{\text{span}} + 1}{n_i}} \quad (17)$$

and also the subroutine on line 10 outputs a policy  $\tilde{\pi}_{\gamma,i}$  such that

$$\hat{V}_{\gamma,i}^{\tilde{\pi}_{\gamma,i}} \geq \hat{V}_{\gamma,i}^* - \frac{1}{n_i} \mathbf{1} \quad (18)$$

$$\|\tilde{V}_{\gamma,i} - \hat{V}_{\gamma,i}^*\|_\infty \leq \frac{1}{n_i} \quad (19)$$

$$\|\hat{V}_{\gamma,i}^{\tilde{\pi}_{\gamma,i}} - V_\gamma^{\tilde{\pi}_{\gamma,i}}\|_\infty \leq \frac{\alpha(\delta, n_i)}{1 - \gamma} \sqrt{\frac{\|\hat{V}_{\gamma,i}^{\tilde{\pi}_{\gamma,i}}\|_{\text{span}} + 1}{n_i}}. \quad (20)$$

**Proof** We fix a pair  $i$  and  $\gamma \in \mathcal{H}_i$  and establish the desired bounds, and then we will take a union bound. From (Zurek and Chen, 2024, Proof of Theorem 9), in particular (Zurek and Chen, 2024, Equations (31) and (32)), we have with probability at least  $1 - 2\delta$  that

$$\|V_\gamma^{\pi^*} - \hat{V}_{\gamma,i}^{\pi^*}\|_\infty \leq \frac{24 \log_2 \log_2\left(\frac{1}{1-\gamma} + 4\right)}{1 - \gamma} \sqrt{16 \frac{\|V_\gamma^*\|_{\text{span}} + 1}{n_i} \log\left(\frac{12SA n_i}{(1 - \gamma)^2 \delta}\right)}$$

and

$$\left\| V_{\gamma}^{\tilde{\pi}_{\gamma}} - \widehat{V}_{\gamma}^{\tilde{\pi}_{\gamma}} \right\|_{\infty} \leq \frac{24 \log_2 \log_2 \left( \frac{1}{1-\gamma} + 4 \right)}{1-\gamma} \sqrt{16 \frac{\left\| \widehat{V}_{\gamma}^{\tilde{\pi}_{\gamma}} \right\|_{\text{span}} + 1}{n_i} \log \left( \frac{12 S A n_i}{(1-\gamma)^2 \delta} \right)}.$$

By definition of  $\mathcal{H}_i$  we have that  $\frac{1}{1-\gamma} \leq n_i$ , which we use to simplify the bounds slightly to

$$\left\| V_{\gamma}^{\pi_{\gamma}^*} - \widehat{V}_{\gamma}^{\pi_{\gamma}^*} \right\|_{\infty} \leq \frac{24 \log_2 \log_2 (n_i + 4)}{1-\gamma} \sqrt{16 \frac{\left\| V_{\gamma}^* \right\|_{\text{span}} + 1}{n_i} \log \left( \frac{12 S A n_i^3}{\delta} \right)} \quad (21)$$

and

$$\left\| V_{\gamma}^{\tilde{\pi}_{\gamma}} - \widehat{V}_{\gamma}^{\tilde{\pi}_{\gamma}} \right\|_{\infty} \leq \frac{24 \log_2 \log_2 (n_i + 4)}{1-\gamma} \sqrt{16 \frac{\left\| \widehat{V}_{\gamma}^{\tilde{\pi}_{\gamma}} \right\|_{\text{span}} + 1}{n_i} \log \left( \frac{12 S A n_i^3}{\delta} \right)}. \quad (22)$$

Now keeping  $i$  fixed and taking a union bound over all  $\gamma \in \mathcal{H}_i$ , we have that inequalities (21) and (22) hold for all  $\gamma \in \mathcal{H}_i$  with probability at least  $1 - 2\delta|\mathcal{H}_i|$ . Also

$$|\mathcal{H}_i| \leq 1 + \log_2 \frac{n_i}{\sqrt{n_i}} \leq 1 + \frac{1}{2} \log_2 n_i \leq 1 + \frac{1}{2} n_i \leq n_i$$

using the facts that  $\log_2 x \leq x$  and that  $n_i \geq 2$  (so  $1 \leq \frac{1}{2} n_i$ ). Now adjusting the failure probability, we have (for any  $\delta_i > 0$ ) that with probability at least  $1 - \delta_i$ , for all  $\gamma \in \mathcal{H}_i$ , both

$$\left\| V_{\gamma}^{\pi_{\gamma}^*} - \widehat{V}_{\gamma}^{\pi_{\gamma}^*} \right\|_{\infty} \leq \frac{24 \log_2 \log_2 (n_i + 4)}{1-\gamma} \sqrt{16 \frac{\left\| V_{\gamma}^* \right\|_{\text{span}} + 1}{n_i} \log \left( \frac{24 S A n_i^4}{\delta_i} \right)} \quad (23)$$

and

$$\left\| V_{\gamma}^{\tilde{\pi}_{\gamma}} - \widehat{V}_{\gamma}^{\tilde{\pi}_{\gamma}} \right\|_{\infty} \leq \frac{24 \log_2 \log_2 (n_i + 4)}{1-\gamma} \sqrt{16 \frac{\left\| \widehat{V}_{\gamma}^{\tilde{\pi}_{\gamma}} \right\|_{\text{span}} + 1}{n_i} \log \left( \frac{24 S A n_i^4}{\delta_i} \right)} \quad (24)$$

are true. Now we union bound over all natural numbers  $i \geq 1$ . We set  $\delta_i = \frac{\delta}{2^i} = \frac{\delta}{n_i}$  and thus obtain that with probability at least  $1 - \sum_{i \geq 1} \delta_i = 1 - \delta$ , for all  $i \geq 1$  and all  $\gamma \in \mathcal{H}_i$  we have both (23) and (24). By our choice of  $\delta_i$ 's these can be written as

$$\begin{aligned} \left\| V_{\gamma}^{\pi_{\gamma}^*} - \widehat{V}_{\gamma}^{\pi_{\gamma}^*} \right\|_{\infty} &\leq \frac{24 \log_2 \log_2 (n_i + 4)}{1-\gamma} \sqrt{16 \frac{\left\| V_{\gamma}^* \right\|_{\text{span}} + 1}{n_i} \log \left( \frac{24 S A n_i^5}{\delta} \right)} \\ &= \frac{\alpha(\delta, n_i)}{1-\gamma} \sqrt{\frac{\left\| V_{\gamma}^* \right\|_{\text{span}} + 1}{n_i}} \end{aligned}$$

and

$$\begin{aligned} \left\| V_{\gamma}^{\tilde{\pi}_{\gamma}} - \widehat{V}_{\gamma}^{\tilde{\pi}_{\gamma}} \right\|_{\infty} &\leq \frac{24 \log_2 \log_2 (n_i + 4)}{1-\gamma} \sqrt{16 \frac{\left\| \widehat{V}_{\gamma}^{\tilde{\pi}_{\gamma}} \right\|_{\text{span}} + 1}{n_i} \log \left( \frac{24 S A n_i^5}{\delta} \right)} \\ &= \frac{\alpha(\delta, n_i)}{1-\gamma} \sqrt{\frac{\left\| \widehat{V}_{\gamma}^{\tilde{\pi}_{\gamma}} \right\|_{\text{span}} + 1}{n_i}} \end{aligned}$$

as desired. ■

**Lemma 7** *Suppose that there exists some  $m, \beta > 1$ ,  $\gamma < 1$ , policy  $\pi$ , MDPs  $(P, r)$ ,  $(\bar{P}, r)$ , such that*

$$\|V_\gamma^{\pi_\gamma^*} - \bar{V}_\gamma^{\pi_\gamma^*}\|_\infty \leq \frac{\beta}{1-\gamma} \sqrt{\frac{\|V_\gamma^{\pi_\gamma^*}\|_{\text{span}} + 1}{m}} \quad (25)$$

$$\bar{V}_\gamma^\pi \geq \bar{V}_\gamma^* - \frac{1}{m} \mathbf{1} \quad (26)$$

$$\|\bar{V}_\gamma^\pi - V_\gamma^\pi\|_\infty \leq \frac{\beta}{1-\gamma} \sqrt{\frac{\|\bar{V}_\gamma^\pi\|_{\text{span}} + 1}{m}} \quad (27)$$

where for any  $\pi'$  we use  $\bar{V}_\gamma^{\pi'}$  to denote the value of policy  $\pi'$  in the DMDP  $(\bar{P}, r, \gamma)$ , where  $\bar{V}_\gamma^*$  denotes the optimal value function in the DMDP  $(\bar{P}, r, \gamma)$ , and where  $\pi_\gamma^*$  is the optimal policy for the DMDP  $(P, r, \gamma)$ . Then

$$\|\bar{V}_\gamma^* - V_\gamma^*\|_\infty \leq \frac{2\beta^2}{(1-\gamma)^2 m} + \frac{4\beta}{1-\gamma} \sqrt{\frac{\|V_\gamma^*\|_{\text{span}} + 1 + \frac{1}{m}}{m}} + \frac{4}{m} \quad (28)$$

and

$$\|\bar{V}_\gamma^* - V_\gamma^*\|_\infty \leq \frac{2\beta^2}{(1-\gamma)^2 m} + \frac{4\beta}{1-\gamma} \sqrt{\frac{\|\bar{V}_\gamma^*\|_{\text{span}} + 1 + \frac{1}{m}}{m}} + \frac{4}{m}. \quad (29)$$

**Proof** Using the triangle inequality several times and (26), we have

$$\begin{aligned} V_\gamma^* &\geq V_\gamma^\pi \\ &\geq \bar{V}_\gamma^\pi - \|\bar{V}_\gamma^\pi - V_\gamma^\pi\|_\infty \mathbf{1} \\ &\geq \bar{V}_\gamma^* - \frac{1}{m} \mathbf{1} - \|\bar{V}_\gamma^\pi - V_\gamma^\pi\|_\infty \mathbf{1} \\ &\geq \bar{V}_\gamma^{\pi_\gamma^*} - \frac{1}{m} \mathbf{1} - \|\bar{V}_\gamma^\pi - V_\gamma^\pi\|_\infty \mathbf{1} \\ &\geq V_\gamma^{\pi_\gamma^*} - \left\| V_\gamma^{\pi_\gamma^*} - \bar{V}_\gamma^{\pi_\gamma^*} \right\|_\infty \mathbf{1} - \frac{1}{m} \mathbf{1} - \|\bar{V}_\gamma^\pi - V_\gamma^\pi\|_\infty \mathbf{1}. \end{aligned}$$

Since by definition  $V_\gamma^{\pi_\gamma^*} = V_\gamma^*$ , subtracting this term from all expressions we obtain

$$\mathbf{0} \geq \bar{V}_\gamma^* - V_\gamma^* - \frac{1}{m} \mathbf{1} - \|\bar{V}_\gamma^\pi - V_\gamma^\pi\|_\infty \mathbf{1} \geq -\left\| V_\gamma^{\pi_\gamma^*} - \bar{V}_\gamma^{\pi_\gamma^*} \right\|_\infty \mathbf{1} - \frac{1}{m} \mathbf{1} - \|\bar{V}_\gamma^\pi - V_\gamma^\pi\|_\infty \mathbf{1}$$

which after rearranging implies that

$$\frac{1}{m} \mathbf{1} + \|\bar{V}_\gamma^\pi - V_\gamma^\pi\|_\infty \mathbf{1} \geq \bar{V}_\gamma^* - V_\gamma^* \geq -\left\| V_\gamma^{\pi_\gamma^*} - \bar{V}_\gamma^{\pi_\gamma^*} \right\|_\infty \mathbf{1}$$

which further implies

$$\left\| \bar{V}_\gamma^\star - V_\gamma^\star \right\|_\infty \leq \left\| \bar{V}_\gamma^\pi - V_\gamma^\pi \right\|_\infty + \left\| V_\gamma^{\pi^\star} - \bar{V}_\gamma^{\pi^\star} \right\|_\infty + \frac{1}{m}. \quad (30)$$

Now we will focus on establishing the first inequality (28) in the lemma statement. For the first term on the RHS of (30), we can use condition (27) and then triangle inequality to bound

$$\begin{aligned} \left\| \bar{V}_\gamma^\pi - V_\gamma^\pi \right\|_\infty &\leq \frac{\beta}{1-\gamma} \sqrt{\frac{\left\| \bar{V}_\gamma^\pi \right\|_{\text{span}} + 1}{m}} \\ &\leq \frac{\beta}{1-\gamma} \sqrt{\frac{\left\| V_\gamma^\star \right\|_{\text{span}} + \left\| \bar{V}_\gamma^\pi - V_\gamma^\star \right\|_{\text{span}} + 1}{m}} \\ &\leq \frac{\beta}{1-\gamma} \sqrt{\frac{\left\| V_\gamma^\star \right\|_{\text{span}} + 2 \left\| \bar{V}_\gamma^\star - V_\gamma^\star \right\|_\infty + \frac{1}{m} + 1}{m}} \end{aligned} \quad (31)$$

using that

$$\left\| \bar{V}_\gamma^\pi - V_\gamma^\star \right\|_{\text{span}} \leq \left\| \bar{V}_\gamma^\pi - \bar{V}_\gamma^\star \right\|_{\text{span}} + \left\| \bar{V}_\gamma^\star - V_\gamma^\star \right\|_{\text{span}} \leq \frac{1}{m} + 2 \left\| \bar{V}_\gamma^\star - V_\gamma^\star \right\|_\infty$$

since  $\|\cdot\|_{\text{span}} \leq 2 \|\cdot\|_\infty$  and  $\mathbf{0} \leq \bar{V}_\gamma^\pi - \bar{V}_\gamma^\star \leq \frac{1}{m} \mathbf{1}$ . For the second term on the RHS of (30), we can use the condition (25). Thus combining (30), (31), and (25), we obtain

$$\begin{aligned} \left\| \bar{V}_\gamma^\star - V_\gamma^\star \right\|_\infty &\leq \frac{\beta}{1-\gamma} \sqrt{\frac{\left\| V_\gamma^\star \right\|_{\text{span}} + 2 \left\| \bar{V}_\gamma^\star - V_\gamma^\star \right\|_\infty + 1 + \frac{1}{m}}{m}} + \frac{\beta}{1-\gamma} \sqrt{\frac{\left\| V_\gamma^{\pi^\star} \right\|_{\text{span}} + 1}{m}} + \frac{1}{m} \\ &\leq \frac{2\beta}{1-\gamma} \sqrt{\frac{\left\| V_\gamma^\star \right\|_{\text{span}} + 1 + \frac{1}{m}}{m}} + \frac{\beta}{1-\gamma} \sqrt{\frac{2 \left\| \bar{V}_\gamma^\star - V_\gamma^\star \right\|_\infty + \frac{1}{m}}{m}} + \frac{1}{m} \end{aligned} \quad (32)$$

where we used that  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  (and  $V_\gamma^\star = V_\gamma^{\pi^\star}$ ). Rearranging (32) as

$$\left\| \bar{V}_\gamma^\star - V_\gamma^\star \right\|_\infty - \frac{\sqrt{2}\beta}{(1-\gamma)\sqrt{m}} \sqrt{\left\| \bar{V}_\gamma^\star - V_\gamma^\star \right\|_\infty} - \frac{2\beta}{1-\gamma} \sqrt{\frac{\left\| V_\gamma^\star \right\|_{\text{span}} + 1 + \frac{1}{m}}{m}} + \frac{1}{m} \leq 0$$

and then using the quadratic formula to find the largest possible root of this quadratic polynomial in

$\sqrt{\left\| \bar{V}_\gamma^\star - V_\gamma^\star \right\|_\infty}$ , we obtain that

$$\sqrt{\left\| \bar{V}_\gamma^\star - V_\gamma^\star \right\|_\infty} \leq \frac{1}{2} \frac{\sqrt{2}\beta}{(1-\gamma)\sqrt{m}} + \frac{1}{2} \sqrt{\left( \frac{\sqrt{2}\beta}{(1-\gamma)\sqrt{m}} \right)^2 + 4 \frac{2\beta}{1-\gamma} \sqrt{\frac{\left\| V_\gamma^\star \right\|_{\text{span}} + 1 + \frac{1}{m}}{m}} + \frac{4}{m}}. \quad (33)$$



Now squaring both sides and using that  $(a + b)^2 \leq 2a^2 + 2b^2$  which implies  $(\frac{a}{2} + \frac{b}{2})^2 \leq \frac{a^2}{2} + \frac{b^2}{2}$ , we obtain

$$\begin{aligned} \|\bar{V}_\gamma^\star - V_\gamma^\star\|_\infty &\leq \frac{1}{2} \frac{2\beta^2}{(1-\gamma)^2 m} + \frac{1}{2} \left( \left( \frac{\sqrt{2}\beta}{(1-\gamma)\sqrt{m}} \right)^2 + 4 \frac{2\beta}{1-\gamma} \sqrt{\frac{\|V_\gamma^\star\|_{\text{span}} + 1 + \frac{1}{m}}{m}} + \frac{4}{m} \right) \\ &= \frac{2\beta^2}{(1-\gamma)^2 m} + \frac{4\beta}{1-\gamma} \sqrt{\frac{\|V_\gamma^\star\|_{\text{span}} + 1 + \frac{1}{m}}{m}} + \frac{4}{m} \end{aligned} \quad (34)$$

as desired.

We next focus on establishing the second inequality (29) from the lemma statement. We now will bound the first term on the RHS of (30) in terms of  $\|\bar{V}_\gamma^\star\|_{\text{span}}$  with an analogous argument but instead using only that  $\|\bar{V}_\gamma^\pi\|_{\text{span}} \leq \|\bar{V}_\gamma^\star\|_{\text{span}} + \|\bar{V}_\gamma^\pi - \bar{V}_\gamma^\star\|_{\text{span}} \leq \|\bar{V}_\gamma^\star\|_{\text{span}} + \frac{1}{m}$ , to obtain

$$\|\bar{V}_\gamma^\pi - V_\gamma^\pi\|_\infty \leq \frac{\beta}{1-\gamma} \sqrt{\frac{\|\bar{V}_\gamma^\star\|_{\text{span}} + \frac{1}{m} + 1}{m}}. \quad (35)$$

We bound the second term on the RHS of (30) in terms of  $\|\bar{V}_\gamma^\star\|_{\text{span}}$  by using (25) and then that  $\|V_\gamma^{\pi^\star}\|_{\text{span}} \leq \|\bar{V}_\gamma^\star\|_{\text{span}} + \|\bar{V}_\gamma^\star - V_\gamma^\star\|_{\text{span}} \leq \|\bar{V}_\gamma^\star\|_{\text{span}} + 2\|\bar{V}_\gamma^\star - V_\gamma^\star\|_\infty$  to obtain

$$\|V_\gamma^{\pi^\star} - \bar{V}_\gamma^{\pi^\star}\|_\infty \leq \frac{\beta}{1-\gamma} \sqrt{\frac{\|V_\gamma^{\pi^\star}\|_{\text{span}} + 1}{m}} \leq \frac{\beta}{1-\gamma} \sqrt{\frac{\|\bar{V}_\gamma^\star\|_{\text{span}} + 2\|\bar{V}_\gamma^\star - V_\gamma^\star\|_\infty + 1}{m}}. \quad (36)$$

Combining (30), (35), and (36), we obtain

$$\begin{aligned} \|\bar{V}_\gamma^\star - V_\gamma^\star\|_\infty &\leq \frac{\beta}{1-\gamma} \sqrt{\frac{\|\bar{V}_\gamma^\star\|_{\text{span}} + \frac{1}{m} + 1}{m}} + \frac{\beta}{1-\gamma} \sqrt{\frac{\|\bar{V}_\gamma^\star\|_{\text{span}} + 2\|\bar{V}_\gamma^\star - V_\gamma^\star\|_\infty + 1}{m}} + \frac{1}{m} \\ &\leq \frac{2\beta}{1-\gamma} \sqrt{\frac{\|\bar{V}_\gamma^\star\|_{\text{span}} + \frac{1}{m} + 1}{m}} + \frac{\beta}{1-\gamma} \sqrt{\frac{2\|\bar{V}_\gamma^\star - V_\gamma^\star\|_\infty}{m}} + \frac{1}{m} \end{aligned} \quad (37)$$

using  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  for the second inequality. Since (37) is identical to (32) except for the presence of  $\|\bar{V}_\gamma^\star\|_{\text{span}}$  instead of  $\|V_\gamma^\star\|_{\text{span}}$ , we can take completely analogous steps to obtain (29). ■

The following lemma is not necessary to show the validity of the lower bounds, but it is necessary for the upper bounds.

**Lemma 8** *Under the event in Lemma 6, for all integers  $i \geq 1$  and all  $\gamma \in \mathcal{H}_i$ , we have*

$$\|\hat{V}_{\gamma,i}^\star - V_\gamma^\star\|_\infty \leq \frac{2\alpha(\delta, n_i)^2}{(1-\gamma)^2 n_i} + \frac{4\alpha(\delta, n_i)}{1-\gamma} \sqrt{\frac{\|V_\gamma^\star\|_{\text{span}} + 1 + \frac{1}{n_i}}{n_i}} + \frac{4}{n_i} \quad (38)$$

and

$$\left\| \widehat{V}_{\gamma,i}^* - V_{\gamma}^* \right\|_{\infty} \leq \frac{2\alpha(\delta, n_i)^2}{(1-\gamma)^2 n_i} + \frac{4\alpha(\delta, n_i)}{1-\gamma} \sqrt{\frac{\left\| \widehat{V}_{\gamma,i}^* \right\|_{\text{span}} + 1 + \frac{1}{n_i}}{n_i}} + \frac{4}{n_i} \quad (39)$$

**Proof** This follows immediately from Lemma 7 since the event described in Lemma 6 exactly meets the conditions of Lemma 7 (by setting  $m = n_i$  and  $\beta = \alpha(\delta, n_i)$ , for all  $i$  and  $\gamma \in \mathcal{H}_i$ ).  $\blacksquare$

Now we show that under the event in Lemma 6, all lower and upper bounds constructed within Algorithm 2 are valid lower/upper bounds of  $\rho^*$ .

**Lemma 9** *Under the event in Lemma 6, for all integers  $i \geq 1$  and all  $\gamma \in \mathcal{H}_i$ , we have*

$$\widehat{L}_i(\gamma) \mathbf{1} \leq \rho^{\pi_{\gamma,i}} \leq \rho^* \leq \widehat{U}_i(\gamma) \mathbf{1}. \quad (40)$$

**Proof** Fix an arbitrary integer  $i \geq 1$  and  $\gamma \in \mathcal{H}_i$ . By (15), optimality of  $\rho^*$ , and (13) we have

$$\begin{aligned} (1-\gamma) \left( \min_s V_{\gamma}^{\pi_{\gamma,i}}(s) \right) \mathbf{1} &\leq \rho^{\pi_{\gamma,i}} \\ &\leq \rho^* \\ &\leq (1-\gamma) \left( \max_s V_{\gamma}^*(s) \right) \mathbf{1}. \end{aligned} \quad (41)$$

We start by lower-bounding the LHS of (41). Using (18) and (20),

$$\begin{aligned} (1-\gamma) \left( \min_s V_{\gamma}^{\pi_{\gamma,i}}(s) \right) \mathbf{1} &\geq (1-\gamma) \left( \min_s \widehat{V}_{\gamma,i}^{\pi_{\gamma,i}}(s) \right) \mathbf{1} - (1-\gamma) \left\| \widehat{V}_{\gamma,i}^{\pi_{\gamma,i}} - V_{\gamma}^{\pi_{\gamma,i}} \right\|_{\infty} \mathbf{1} \\ &\geq (1-\gamma) \left( \min_s \widehat{V}_{\gamma,i}^*(s) \right) \mathbf{1} - \frac{1-\gamma}{n_i} \mathbf{1} - (1-\gamma) \left\| \widehat{V}_{\gamma,i}^{\pi_{\gamma,i}} - \widehat{V}_{\gamma,i}^* \right\|_{\infty} \mathbf{1} \\ &\geq (1-\gamma) \left( \min_s \widehat{V}_{\gamma,i}^*(s) \right) \mathbf{1} - \frac{1-\gamma}{n_i} \mathbf{1} - \alpha(\delta, n_i) \sqrt{\frac{\left\| \widehat{V}_{\gamma,i}^{\pi_{\gamma,i}} \right\|_{\text{span}} + 1}{n_i}} \mathbf{1}. \end{aligned} \quad (42)$$

Now we replace the quantities in (42) with observable quantities in terms of  $\widetilde{V}_{\gamma,i}$ . Using the requirement (18) to relate  $\widehat{V}_{\gamma,i}^{\pi_{\gamma,i}}$  and  $\widehat{V}_{\gamma,i}^*$ , and then (19) which bounds  $\left\| \widetilde{V}_{\gamma,i} - \widehat{V}_{\gamma,i}^* \right\|_{\infty}$ , we have

$$\begin{aligned} \left\| \widehat{V}_{\gamma,i}^{\pi_{\gamma,i}} \right\|_{\text{span}} &\leq \left\| \widehat{V}_{\gamma,i}^* \right\|_{\text{span}} + \left\| \widehat{V}_{\gamma,i}^* - \widehat{V}_{\gamma,i}^{\pi_{\gamma,i}} \right\|_{\text{span}} \\ &\leq \left\| \widetilde{V}_{\gamma,i} \right\|_{\text{span}} + \left\| \widetilde{V}_{\gamma,i} - \widehat{V}_{\gamma,i}^* \right\|_{\text{span}} + \left\| \widehat{V}_{\gamma,i}^* - \widehat{V}_{\gamma,i}^{\pi_{\gamma,i}} \right\|_{\text{span}} \\ &\leq \left\| \widetilde{V}_{\gamma,i} \right\|_{\text{span}} + 2 \left\| \widetilde{V}_{\gamma,i} - \widehat{V}_{\gamma,i}^* \right\|_{\infty} + \frac{1}{n_i} \\ &\leq \left\| \widetilde{V}_{\gamma,i} \right\|_{\text{span}} + \frac{2}{n_i} + \frac{1}{n_i} \end{aligned} \quad (43)$$

(where  $\|\widehat{V}_{\gamma,i}^* - \widehat{V}_{\gamma,i}^{\pi_{\gamma,i}}\|_{\text{span}} \leq \frac{1}{n_i}$  since  $\mathbf{0} \leq \widehat{V}_{\gamma,i}^* - \widehat{V}_{\gamma,i}^{\pi_{\gamma,i}} \leq \frac{1}{n_i} \mathbf{1}$ ). Using this bound (43), as well as (19) again, we can further bound (42) as

$$\begin{aligned}
& (1 - \gamma) \left( \min_s \widehat{V}_{\gamma,i}^*(s) \right) \mathbf{1} - \frac{1 - \gamma}{n_i} \mathbf{1} - \alpha(\delta, n_i) \sqrt{\frac{\|\widehat{V}_{\gamma,i}^{\pi_{\gamma,i}}\|_{\text{span}} + 1}{n_i}} \mathbf{1} \\
& \geq (1 - \gamma) \left( \min_s \widetilde{V}_{\gamma,i}(s) \right) \mathbf{1} - (1 - \gamma) \|\widetilde{V}_{\gamma,i} - \widehat{V}_{\gamma,i}^*\|_{\infty} \mathbf{1} - \frac{1 - \gamma}{n_i} \mathbf{1} - \alpha(\delta, n_i) \sqrt{\frac{\|\widetilde{V}_{\gamma,i}\|_{\text{span}} + \frac{3}{n_i} + 1}{n_i}} \mathbf{1} \\
& \geq (1 - \gamma) \left( \min_s \widetilde{V}_{\gamma,i}(s) \right) \mathbf{1} - 2 \frac{1 - \gamma}{n_i} \mathbf{1} - \alpha(\delta, n_i) \sqrt{\frac{\|\widetilde{V}_{\gamma,i}\|_{\text{span}} + \frac{3}{n_i} + 1}{n_i}} \mathbf{1} \\
& = \widehat{L}_i(\gamma) \mathbf{1}
\end{aligned} \tag{44}$$

Now we upper-bound the RHS of (41). The inequality (39) plays a key role. Using this bound, as well as (19) to bound  $\|\widetilde{V}_{\gamma,i} - \widehat{V}_{\gamma,i}^*\|_{\infty}$ , and the fact that  $\|\widehat{V}_{\gamma,i}^*\|_{\text{span}} \leq \|\widetilde{V}_{\gamma,i}\|_{\text{span}} + \|\widehat{V}_{\gamma,i}^* - \widetilde{V}_{\gamma,i}\|_{\text{span}} \leq \|\widetilde{V}_{\gamma,i}\|_{\text{span}} + \frac{2}{n_i}$  (since  $\|\cdot\|_{\text{span}} \leq 2 \|\cdot\|_{\infty}$  and using (19)), we obtain

$$\begin{aligned}
& (1 - \gamma) \left( \max_s V_{\gamma}^*(s) \right) \mathbf{1} \\
& \leq (1 - \gamma) \left( \max_s \widehat{V}_{\gamma,i}^*(s) \right) \mathbf{1} + (1 - \gamma) \|\widehat{V}_{\gamma,i}^* - V_{\gamma}^*\|_{\infty} \mathbf{1} \\
& \leq (1 - \gamma) \left( \max_s \widetilde{V}_{\gamma,i}(s) \right) \mathbf{1} + (1 - \gamma) \|\widetilde{V}_{\gamma,i} - \widehat{V}_{\gamma,i}^*\|_{\infty} \mathbf{1} + (1 - \gamma) \|\widehat{V}_{\gamma,i}^* - V_{\gamma}^*\|_{\infty} \mathbf{1} \\
& \leq (1 - \gamma) \left( \max_s \widetilde{V}_{\gamma,i}(s) \right) \mathbf{1} + \frac{1 - \gamma}{n_i} \mathbf{1} + \frac{2\alpha(\delta, n_i)^2}{(1 - \gamma)n_i} \mathbf{1} + 4\alpha(\delta, n_i) \sqrt{\frac{\|\widehat{V}_{\gamma,i}^*\|_{\text{span}} + 1 + \frac{1}{n_i}}{n_i}} \mathbf{1} + (1 - \gamma) \frac{4}{n_i} \mathbf{1} \\
& \leq (1 - \gamma) \left( \max_s \widetilde{V}_{\gamma,i}(s) \right) \mathbf{1} + 5 \frac{1 - \gamma}{n_i} \mathbf{1} + \frac{2\alpha(\delta, n_i)^2}{(1 - \gamma)n_i} \mathbf{1} + 4\alpha(\delta, n_i) \sqrt{\frac{\|\widetilde{V}_{\gamma,i}\|_{\text{span}} + 1 + \frac{3}{n_i}}{n_i}} \mathbf{1} \\
& = \widehat{U}_i(\gamma) \mathbf{1}.
\end{aligned} \tag{45}$$

Combining (41), (44), and (45), and unfixing  $i$  and  $\gamma$ , we obtain the desired conclusion.  $\blacksquare$

Lemma 9 implies that whenever the algorithm terminates, the resulting policy will be  $\varepsilon$ -optimal and the resulting confidence interval will be valid and of size  $\leq \varepsilon$ . Next we show that the algorithm will terminate by a certain iteration, by showing that on this iteration the confidence interval corresponding to some discount factor will be small.

First we define, for all integers  $i \geq 1$ , the discount factor  $\gamma_i^*$  as

$$\gamma_i^* = \inf \left\{ \gamma \in \mathcal{H}_i : \frac{1}{1 - \gamma} \geq \sqrt{n_i \left( \|h^*\|_{\text{span}} + 1 \right)} \right\}. \tag{46}$$

(If  $\|h^*\|_{\text{span}}$  is large relative to  $n_i$  then, since the largest value of  $\frac{1}{1 - \gamma}$  for  $\gamma \in \mathcal{H}_i$  is  $n_i$ , the above set might be empty, in which case by usual convention we would have  $\gamma_i^* = \inf \emptyset = \infty$ .)

**Lemma 10** For all integers  $i \geq 1$ , if

$$\|h^*\|_{\text{span}} + 1 \leq n_i \quad (47)$$

then  $\gamma_i^*$  is finite, and furthermore

$$\sqrt{n_i (\|h^*\|_{\text{span}} + 1)} \leq \frac{1}{1 - \gamma_i^*} \leq 2\sqrt{n_i (\|h^*\|_{\text{span}} + 1)}. \quad (48)$$

**Proof** First, by condition (47) we have that  $\sqrt{n_i (\|h^*\|_{\text{span}} + 1)} \leq n_i$ , and since the largest element of  $\mathcal{H}_i$  is  $n_i$ , the set in the definition (46) of  $\gamma_i^*$  will be nonempty and thus  $\gamma_i^*$  will be finite.

Furthermore since  $\|h^*\|_{\text{span}} + 1 \geq 1$ , we have  $\sqrt{n_i (\|h^*\|_{\text{span}} + 1)} \geq \sqrt{n_i}$ .  $\sqrt{n_i}$  may not be a member of  $\mathcal{H}_i$  but the smallest power of 2 which is  $\geq \sqrt{n_i}$  will be a member of  $\mathcal{H}_i$ . Therefore  $\min \mathcal{H}_i \leq \frac{1}{1 - \gamma_i^*} \leq \max \mathcal{H}_i$ , so by the construction of  $\mathcal{H}_i$  in line (8) of Algorithm 2, condition (48) must hold as  $\mathcal{H}_i$  contains all powers of 2 within  $[\min \mathcal{H}_i, \max \mathcal{H}_i]$ . ■

**Lemma 11** Under the event in Lemma 6, for all integers  $i \geq 1$  such that  $\|h^*\|_{\text{span}} + 1 \leq n_i$ , it holds that

$$\widehat{U}_i(\gamma_i^*) - \widehat{L}_i(\gamma_i^*) \leq 90\alpha(\delta, n_i)^2 \sqrt{\frac{\|h^*\|_{\text{span}} + 1}{n_i}}.$$

In particular there exist some absolute constants  $C_1, C_2$  such that letting

$$B = \left\lceil \log_2 \left( C_1 \frac{\|h^*\|_{\text{span}} + 1}{\varepsilon^2} \log^3 \left( \frac{C_2 S A (\|h^*\|_{\text{span}} + 1)}{\delta \varepsilon} \right) \right) \right\rceil$$

we have

$$\widehat{U}_B(\gamma_B^*) - \widehat{L}_B(\gamma_B^*) \leq \varepsilon.$$

**Proof** For any integer  $i \geq 1$  and any  $\gamma \in \mathcal{H}_i$  we have

$$\begin{aligned} \widehat{U}_i(\gamma) - \widehat{L}_i(\gamma) &= (1 - \gamma) \left( \max_s \widetilde{V}_{\gamma,i}(s) - \min_s \widetilde{V}_{\gamma,i}(s) \right) + 7 \frac{1 - \gamma}{n_i} + \frac{2\alpha(\delta, n_i)^2}{(1 - \gamma)n_i} \\ &\quad + 5\alpha(\delta, n_i) \sqrt{\frac{\|\widetilde{V}_{\gamma,i}\|_{\text{span}} + 1 + \frac{3}{n_i}}{n_i}} \\ &= (1 - \gamma) \|\widetilde{V}_{\gamma,i}\|_{\text{span}} + 7 \frac{1 - \gamma}{n_i} + \frac{2\alpha(\delta, n_i)^2}{(1 - \gamma)n_i} + 5\alpha(\delta, n_i) \sqrt{\frac{\|\widetilde{V}_{\gamma,i}\|_{\text{span}} + 1 + \frac{3}{n_i}}{n_i}}. \end{aligned} \quad (49)$$

Now we will set  $\gamma = \gamma_i^*$  and relate all terms in (49) to  $\|h^*\|_{\text{span}}$ , but first we bound  $\|\tilde{V}_{\gamma,i}\|_{\text{span}}$ . Under the event of Lemma 6 we have from (19), and from (38) in Lemma 8, that

$$\begin{aligned}
\|\tilde{V}_{\gamma,i}\|_{\text{span}} &\leq \|V_\gamma^*\|_{\text{span}} + \|\tilde{V}_{\gamma,i} - V_\gamma^*\|_{\text{span}} \\
&\leq \|V_\gamma^*\|_{\text{span}} + 2 \|\tilde{V}_{\gamma,i} - V_\gamma^*\|_\infty \\
&\leq \|V_\gamma^*\|_{\text{span}} + 2 \|\tilde{V}_{\gamma,i} - \hat{V}_{\gamma,i}^*\|_\infty + 2 \|\hat{V}_{\gamma,i}^* - V_\gamma^*\|_\infty \\
&\leq \|V_\gamma^*\|_{\text{span}} + 2 \frac{1}{n_i} + 2 \left( \frac{2\alpha(\delta, n_i)^2}{(1-\gamma)^2 n_i} + \frac{4\alpha(\delta, n_i)}{1-\gamma} \sqrt{\frac{\|V_\gamma^*\|_{\text{span}} + 1 + \frac{1}{n_i}}{n_i}} + \frac{4}{n_i} \right).
\end{aligned} \tag{50}$$

Bounding (50) by substituting  $\gamma = \gamma_i^*$  and using Lemma 10 to bound  $\frac{1}{1-\gamma_i^*}$ , and using that  $\|V_\gamma^*\|_{\text{span}} \leq 2 \|h^*\|_{\text{span}}$  (Wei et al., 2020, Lemma 2), we have

$$\begin{aligned}
\|\tilde{V}_{\gamma_i^*,i}\|_{\text{span}} &\leq 2 \|h^*\|_{\text{span}} + \frac{2}{n_i} + \frac{4\alpha(\delta, n_i)^2}{(1-\gamma_i^*)^2 n_i} + \frac{8\alpha(\delta, n_i)}{1-\gamma_i^*} \sqrt{\frac{2 \|h^*\|_{\text{span}} + 1 + \frac{1}{n_i}}{n_i}} + \frac{8}{n_i} \\
&\leq 2 \|h^*\|_{\text{span}} + \frac{10}{n_i} + \frac{4\alpha(\delta, n_i)^2}{n_i} 4n_i (\|h^*\|_{\text{span}} + 1) \\
&\quad + 8\alpha(\delta, n_i) 2 \sqrt{n_i (\|h^*\|_{\text{span}} + 1)} \sqrt{\frac{2 \|h^*\|_{\text{span}} + 1 + \frac{1}{n_i}}{n_i}} \\
&\leq 2 \|h^*\|_{\text{span}} + 5 + 16\alpha(\delta, n_i)^2 (\|h^*\|_{\text{span}} + 1) \\
&\quad + 16\alpha(\delta, n_i) \sqrt{n_i (\|h^*\|_{\text{span}} + 1)} \sqrt{\frac{2 \|h^*\|_{\text{span}} + 2}{n_i}} \\
&= 2 \|h^*\|_{\text{span}} + 5 + 16\alpha(\delta, n_i)^2 (\|h^*\|_{\text{span}} + 1) + 16\sqrt{2}\alpha(\delta, n_i) (\|h^*\|_{\text{span}} + 1) \\
&\leq 44\alpha(\delta, n_i)^2 (\|h^*\|_{\text{span}} + 1)
\end{aligned} \tag{51}$$

where in the third inequality we used that  $n_i \geq 2$  to bound  $\frac{1}{n_i} \leq \frac{1}{2} \leq 1$ , and in the final inequality we used that  $\alpha(\delta, n_i) \geq 1$  (which is immediate from the form of  $\alpha$  and the facts that  $\delta \leq 1, n_i \geq 1$ ) and that  $5 + 16 + 16\sqrt{2} \leq 44$ . Substituting the inequality (51) into (49) and simplifying in the same

ways (including setting  $\gamma = \gamma_i^*$ ), we obtain that

$$\begin{aligned}
& \widehat{U}_i(\gamma_i^*) - \widehat{L}_i(\gamma_i^*) \\
& \leq (1 - \gamma_i^*) \left\| \widetilde{V}_{\gamma_i^*, i} \right\|_{\text{span}} + 7 \frac{1 - \gamma_i^*}{n_i} + \frac{2\alpha(\delta, n_i)^2}{(1 - \gamma_i^*)n_i} + 5\alpha(\delta, n_i) \sqrt{\frac{\left\| \widetilde{V}_{\gamma_i^*, i} \right\|_{\text{span}} + 1 + \frac{3}{n_i}}{n_i}} \\
& \leq (1 - \gamma_i^*) 44\alpha(\delta, n_i)^2 \left( \|h^*\|_{\text{span}} + 1 \right) + 7 \frac{1 - \gamma_i^*}{n_i} + \frac{2\alpha(\delta, n_i)^2}{(1 - \gamma_i^*)n_i} \\
& \quad + 5\alpha(\delta, n_i) \sqrt{\frac{44\alpha(\delta, n_i)^2 \left( \|h^*\|_{\text{span}} + 1 \right) + 1 + \frac{3}{n_i}}{n_i}} \\
& \leq \frac{1}{\sqrt{n_i \left( \|h^*\|_{\text{span}} + 1 \right)}} 44\alpha(\delta, n_i)^2 \left( \|h^*\|_{\text{span}} + 1 \right) + 7 \frac{1}{n_i} + \frac{2\alpha(\delta, n_i)^2}{n_i} 2\sqrt{n_i \left( \|h^*\|_{\text{span}} + 1 \right)} \\
& \quad + 5\alpha(\delta, n_i) \sqrt{\frac{47\alpha(\delta, n_i)^2 \left( \|h^*\|_{\text{span}} + 1 \right)}{n_i}} \\
& = \left( 48 + 5\sqrt{47} \right) \alpha(\delta, n_i)^2 \sqrt{\frac{\|h^*\|_{\text{span}} + 1}{n_i}} + \frac{7}{n_i} \\
& \leq 83\alpha(\delta, n_i)^2 \sqrt{\frac{\|h^*\|_{\text{span}} + 1}{n_i}} + 7\sqrt{\frac{\|h^*\|_{\text{span}} + 1}{n_i}} \\
& \leq 90\alpha(\delta, n_i)^2 \sqrt{\frac{\|h^*\|_{\text{span}} + 1}{n_i}}.
\end{aligned}$$

We have thus shown the first part of the lemma statement.

For the second part of the lemma statement, we need to find some  $i$  such that  $90\alpha(\delta, n_i)^2 \sqrt{\frac{\|h^*\|_{\text{span}} + 1}{n_i}} \leq \varepsilon$ . First we compute

$$\begin{aligned}
\alpha(\delta, n_i)^4 &= 24^4 16^2 \log^2 \left( \frac{24SAn_i^5}{\delta} \right) \log^4(\log_2(n_i + 4)) \\
&\leq 24^4 16^2 117 \log^2 \left( \frac{24SAn_i^5}{\delta} \right) \log(n_i) \\
&\leq 24^4 16^2 5^2 117 \log^2 \left( \frac{24SAn_i}{\delta} \right) \log(n_i) \\
&\leq 24^4 16^2 5^2 117 \log^3 \left( \frac{24SAn_i}{\delta} \right)
\end{aligned}$$



so

$$\begin{aligned}
& 90\alpha(\delta, n_i)^2 \sqrt{\frac{\|h^*\|_{\text{span}} + 1}{n_i}} \leq \varepsilon \\
& \iff n_i \geq 90^2 \alpha(\delta, n_i)^4 \frac{\|h^*\|_{\text{span}} + 1}{\varepsilon^2} \\
& \stackrel{\text{(I)}}{\iff} n_i \geq 90^2 24^4 16^2 5^2 117 \frac{\|h^*\|_{\text{span}} + 1}{\varepsilon^2} \log^3 \left( \frac{24SA n_i}{\delta} \right) \\
& \stackrel{\text{(II)}}{\iff} n_i \geq 10 \cdot 90^2 24^4 16^2 5^2 117 \frac{\|h^*\|_{\text{span}} + 1}{\varepsilon^2} \log^3 \left( \frac{240SA}{\delta} 90^2 24^4 16^2 5^2 117 \frac{\|h^*\|_{\text{span}} + 1}{\varepsilon^2} \right)
\end{aligned} \tag{52}$$

where (I) is due to Lemma 12 and (II) is due to Lemma 13. Therefore if we set

$$B = \left\lceil \log_2 \left( C_1 \frac{\|h^*\|_{\text{span}} + 1}{\varepsilon^2} \log^3 \left( \frac{C_2 SA (\|h^*\|_{\text{span}} + 1)}{\delta \varepsilon} \right) \right) \right\rceil$$

where  $C_1 = 10 \cdot 90^2 24^4 16^2 5^2 117$  and  $C_2 = 240 \cdot 90^2 24^4 16^2 5^2 117$ , then (52) will be satisfied for  $i = B$ .  $\blacksquare$

**Lemma 12**  $(\log_2 \log_2(x + 4))^4 \leq 117 \log x$  for all  $x \geq 2$ .

**Proof** Since  $z \mapsto 3 \cdot 2^{z/4}$  is convex,  $3 \cdot 2^{z/4}$  is lower-bounded by its tangent line at  $z = 4$ , so for all  $z$  we have

$$3 \cdot 2^{z/4} \geq 3 \cdot 2^{4/4} + 3 \cdot 2^{4/4} \cdot \ln(2^{1/4}) \cdot (z - 4) = 6 + \frac{3 \ln(2)}{2} (z - 4).$$

Since  $\frac{3 \ln(2)}{2} \in (1, 1.1)$ , we have for all  $z \geq 4$  that

$$6 + \frac{3 \ln(2)}{2} (z - 4) \geq 6 + (z - 4) = z + 2$$

and for all  $0 \leq z \leq 4$  that

$$6 + \frac{3 \ln(2)}{2} (z - 4) \geq 6 + 1.1(z - 4) > 1.1z + 1.5 \geq z + 1.5$$

so we have that  $3 \cdot 2^{z/4} \geq z + 1.5$  for all  $z \geq 0$ . Also  $\log_2(1 + \log_2 3) < 1.5$ , so we have  $3 \cdot 2^{z/4} \geq z + \log_2(1 + \log_2 3)$ . Now letting  $y = 2^z$  or equivalently  $z = \log_2 y$  (and we must have  $y \geq 1$  since  $z \geq 0$ ), we have that

$$\begin{aligned}
\log_2(y + \log_2 3) & \leq \log_2(y + y \log_2 3) = \log_2 y + \log_2(1 + \log_2 3) \\
& = z + \log_2(1 + \log_2 3) \leq 3 \cdot 2^{z/4} = 3y^{1/4}.
\end{aligned}$$

This implies

$$(\log_2(y + \log_2 3))^4 \leq 3^4 y.$$

Now letting  $x = 2^y$  or equivalently  $y = \log_2(x)$  (and we must have  $x \geq 2$  since  $y \geq 1$ ), we have that

$$(\log_2(\log_2(x+4)))^4 \leq (\log_2(\log_2(3x)))^4 = (\log_2(\log_2 x + \log_2 3))^4 = (\log_2(y + \log_2 3))^4 \leq 3^4 y,$$

where the first inequality is because  $x+4 \leq 3x$  for  $x \geq 2$ . We also have

$$3^4 y = 3^4 \log_2(x) = 3^4 \log_2(e) \ln x < 117 \ln(x).$$

Thus, we can conclude the desired result for any  $x \geq 2$  by making the appropriate choice of  $z$  (that is,  $z = \log_2 \log_2 x$ ).  $\blacksquare$

**Lemma 13** *Suppose that  $x, y \geq 1$ . Then if  $n \geq 10x \log^3(10xy)$ , we have that  $n \geq x \log^3(yn)$ .*

**Proof** The desired conclusion  $n \geq x \log^3(yn)$  is equivalent to  $\frac{n}{\log^3(yn)} \geq x$ . The derivative of  $\frac{n}{\log^3(yn)}$  with respect to  $n$  is  $\frac{\log(yn)-3}{\log^4(yn)}$  which is  $\geq 0$  if  $n \geq \frac{e^3}{y}$ , so the function  $\frac{n}{\log^3(yn)}$  is monotone non-decreasing for  $n \geq \frac{e^3}{y}$ . Thus it suffices to find some  $m$  such that  $m \geq \frac{e^3}{y}$  and  $\frac{m}{\log^3(ym)} \geq x$ , because then by monotonicity we have that  $n \geq m$  implies  $\frac{n}{\log^3(yn)} \geq \frac{m}{\log^3(ym)} \geq x$ .

Now we claim that  $m = 10x \log^3(10xy)$  satisfies these conditions. First, since  $x, y \geq 1$ , we have that  $\log(10xy) > \log(e^2) = 2$ , so  $m = 10x \log^3(10xy) > 10 \cdot 2^3 = 80 \geq e^3 \geq \frac{e^3}{y}$ , meeting the first required condition. Next, we have that

$$\begin{aligned} \log^3(ym) &= \log^3(y10x \log^3(10xy)) = (\log(10xy) + 3 \log \log(10xy))^3 \\ &\leq \left( \log(10xy) + \frac{3}{e} \log(10xy) \right)^3 \leq 10 \log^3(10xy) \end{aligned}$$

where we used that  $\log x \leq \frac{x}{e}$  and then in the last step we used that  $(1 + \frac{3}{e})^3 \leq 10$ . (To obtain the inequality  $\log x \leq \frac{x}{e}$ , first we can show the inequality  $x \leq e^{x/e}$  by noting that  $e^{x/e}$  is convex, using the first-order convexity condition, and taking the tangent line to  $e^{x/e}$  at  $x = e$ . Then we can take log of both sides.) Thus

$$\frac{m}{\log^3(ym)} \geq \frac{10x \log^3(10xy)}{10 \log^3(10xy)} \geq x$$

as desired.  $\blacksquare$

Combining the consequences of Lemmas 9 and 11, we can prove Theorem 2.

**Proof of Theorem 2** For this proof we assume the event in Lemma 6 holds, which occurs with probability at least  $1 - \delta$ . First we argue that the algorithm terminates and uses at most the claimed number of samples. By Lemma 11, we have that  $\widehat{U}_B(\gamma_B^*) - \widehat{L}_B(\gamma_B^*) \leq \varepsilon$ . Since this would trigger the termination condition of Algorithm 2, this implies that the algorithm must terminate on or before

iteration  $B$ . Therefore the total number of samples used per state-action pair is at most

$$\begin{aligned} \sum_{i=1}^B n_i &= \sum_{i=1}^B 2^i \leq 2^{B+1} \\ &= 2 \cdot 2^{\left\lceil \log_2 \left( C_1 \frac{\|h^*\|_{\text{span}} + 1}{\varepsilon^2} \log^3 \left( \frac{C_2 SA(\|h^*\|_{\text{span}} + 1)}{\delta \varepsilon} \right) \right) \right\rceil} \\ &\leq 4C_1 \frac{\|h^*\|_{\text{span}} + 1}{\varepsilon^2} \log^3 \left( \frac{C_2 SA(\|h^*\|_{\text{span}} + 1)}{\delta \varepsilon} \right) =: N. \end{aligned}$$

Also we can further bound  $B$ , which is an upper bound on the number of iterations, by

$$\begin{aligned} B &= \left\lceil \log_2 \left( C_1 \frac{\|h^*\|_{\text{span}} + 1}{\varepsilon^2} \log^3 \left( \frac{C_2 SA(\|h^*\|_{\text{span}} + 1)}{\delta \varepsilon} \right) \right) \right\rceil \\ &\leq \log_2 \left( C_1 \frac{\|h^*\|_{\text{span}} + 1}{\varepsilon^2} \log^3 \left( \frac{C_2 SA(\|h^*\|_{\text{span}} + 1)}{\delta \varepsilon} \right) \right) + 1 \\ &= \log_2 \left( 2C_1 \frac{\|h^*\|_{\text{span}} + 1}{\varepsilon^2} \log^3 \left( \frac{C_2 SA(\|h^*\|_{\text{span}} + 1)}{\delta \varepsilon} \right) \right) \\ &\leq \log_2(N). \end{aligned}$$

Since the algorithm terminates and by definition of  $\hat{L}$  and  $\hat{U}$ , it is immediate that we have  $\hat{U} - \hat{L} \leq \varepsilon$ . By Lemma 9, we have for all  $i \geq 1$  and all  $\gamma \in \mathcal{H}_i$  that

$$\hat{L}_i(\gamma)\mathbf{1} \leq \rho^{\tilde{\pi}_{\gamma,i}} \leq \rho^* \leq \hat{U}_i(\gamma)\mathbf{1}$$

so in particular we have, for the (random) final iteration  $I$ , that

$$\hat{L}\mathbf{1} = \hat{L}_I(\hat{\gamma}_I)\mathbf{1} \leq \rho^{\tilde{\pi}_{\gamma,I}} = \rho^{\hat{\pi}} \leq \rho^* \leq \hat{U}_I(\hat{\gamma}_I)\mathbf{1} = \hat{U}\mathbf{1}.$$

Finally, combining this with the fact that  $\hat{U} - \hat{L} \leq \varepsilon$  we see that

$$\rho^{\hat{\pi}} \geq \hat{L}\mathbf{1} \geq \hat{U}\mathbf{1} - \varepsilon\mathbf{1} \geq \rho^* - \varepsilon\mathbf{1}$$

as desired. ■

## C.2. Proof of Theorem 1

We start by recalling the definitions of some objects which appear in Algorithm 1 for convenience. These definitions will be in effect for the entirety of this subsection. We define the empirical transition kernel  $\hat{P}(s' | s, a) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{S_{s,a}^i = s'\}$ , for all  $s' \in \mathcal{S}$ , using the  $n$  samples drawn from all state-action pairs within Algorithm 1. Let  $\mathcal{H} = \{\gamma : \text{there exists an integer } k \text{ such that } \sqrt{n} \leq \frac{1}{1-\gamma} = 2^k \leq n\}$ .

Then for all  $\gamma \in \mathcal{H}$  we define the policy  $\tilde{\pi}_\gamma$  and value function  $\tilde{V}_\gamma$  as the outputs of  $\text{SOLVEDMDP}(\hat{P}, r, \gamma, \frac{1}{n})$ .

Now let  $\gamma^*$  be the smallest member of  $\mathcal{H}$  that corresponds to an effective horizon at least  $\sqrt{n(\|h^*\|_{\text{span}} + 1)}$ , that is,

$$\gamma^* = \inf \left\{ \gamma \in \mathcal{H} : \frac{1}{1-\gamma} \geq \sqrt{n(\|h^*\|_{\text{span}} + 1)} \right\}. \quad (53)$$

**Lemma 14** *If  $n \geq 4$  then  $\mathcal{H}$  is nonempty. Furthermore if also  $n \geq 16(\|h^*\|_{\text{span}} + 1)$ , then*

$$\sqrt{n(\|h^*\|_{\text{span}} + 1)} \leq \frac{1}{1 - \gamma^*} \leq 2\sqrt{n(\|h^*\|_{\text{span}} + 1)}. \quad (54)$$

**Proof** First we note that  $\mathcal{H}$  is nonempty if the interval  $[\sqrt{n}, n]$  contains some power of 2, which is ensured if  $n \geq 2\sqrt{n}$ , or equivalently if  $n \geq 4$ . It is also simple to check for  $n \in \{1, 2, 3\}$  that  $\mathcal{H}$  is also nonempty.

Next, the smallest power of 2 which is  $\geq \sqrt{n(\|h^*\|_{\text{span}} + 1)}$  is at most  $2\sqrt{n(\|h^*\|_{\text{span}} + 1)}$ , so to guarantee it is contained in  $\mathcal{H}$  we need it to be  $\leq$  the largest element of  $\mathcal{H}$ , which is at least  $n/2$ . Therefore if

$$2\sqrt{n(\|h^*\|_{\text{span}} + 1)} \leq n/2 \iff n \geq 16(\|h^*\|_{\text{span}} + 1)$$

then (54) will be satisfied. ■

**Lemma 15** *Define the function  $\alpha(\delta, \tilde{n}) = 24\sqrt{16 \log\left(\frac{24SA\tilde{n}^5}{\delta}\right) \log_2(\log_2(\tilde{n} + 4))}$ . Fix  $\delta > 0$ . Then with probability at least  $1 - \delta$ , we have for all  $\gamma \in \mathcal{H}$  that*

$$\|V_\gamma^{\pi_\gamma^*} - \hat{V}_\gamma^{\pi_\gamma^*}\|_\infty \leq \frac{\alpha(\delta, n)}{1 - \gamma} \sqrt{\frac{\|V_\gamma^{\pi_\gamma^*}\|_{\text{span}} + 1}{n}} \quad (55)$$

and also the subroutine on line 7 outputs a policy  $\tilde{\pi}_\gamma$  such that

$$\hat{V}_\gamma^{\tilde{\pi}_\gamma} \geq \hat{V}_\gamma^* - \frac{1}{n} \mathbf{1} \quad (56)$$

$$\|\tilde{V}_\gamma - \hat{V}_\gamma^*\|_\infty \leq \frac{1}{n} \quad (57)$$

$$\|\hat{V}_\gamma^{\tilde{\pi}_\gamma} - V_\gamma^{\tilde{\pi}_\gamma}\|_\infty \leq \frac{\alpha(\delta, n)}{1 - \gamma} \sqrt{\frac{\|\hat{V}_\gamma^{\tilde{\pi}_\gamma}\|_{\text{span}} + 1}{n}}. \quad (58)$$

**Proof** Following identical steps as to the proof of Lemma 6, but with  $n$  in place of  $n_i$ , up until equations (23) and (24), we obtain that with probability at least  $1 - \delta$  the two inequalities

$$\|V_\gamma^{\pi_\gamma^*} - \hat{V}_\gamma^{\pi_\gamma^*}\|_\infty \leq \frac{24 \log_2 \log_2(n + 4)}{1 - \gamma} \sqrt{16 \frac{\|V_\gamma^*\|_{\text{span}} + 1}{n} \log\left(\frac{24SA n^4}{\delta}\right)}$$

and

$$\|V_\gamma^{\tilde{\pi}_\gamma} - \hat{V}_\gamma^{\tilde{\pi}_\gamma}\|_\infty \leq \frac{24 \log_2 \log_2(n + 4)}{1 - \gamma} \sqrt{16 \frac{\|\hat{V}_\gamma^{\tilde{\pi}_\gamma}\|_{\text{span}} + 1}{n} \log\left(\frac{24SA n^4}{\delta}\right)}$$

both hold. We obtain the desired conclusion by noting that

$$\alpha(\delta, n) \geq 24 \log_2 \log_2(n + 4) \sqrt{16 \log\left(\frac{24SA n^4}{\delta}\right)}.$$
■

**Lemma 16** *Under the event described in Lemma 15, we have*

$$\rho^{\tilde{\pi}_\gamma} \geq \hat{L}(\gamma)\mathbf{1} = (1-\gamma) \min_s \tilde{V}_\gamma(s)\mathbf{1} - 2\frac{1-\gamma}{n}\mathbf{1} - \alpha(\delta, n) \sqrt{\frac{\|\tilde{V}_\gamma\|_{\text{span}} + \frac{3}{n} + 1}{n}}\mathbf{1} \quad (59)$$

for all  $\gamma \in \mathcal{H}$ .

**Proof** The proof is very similar to the first part of the proof of Lemma 9. Fix  $\gamma \in \mathcal{H}$ . Using inequality (15) from Lemma 5, then using the triangle inequality, then (56), then (58), then the triangle inequality again, then (57), we have

$$\begin{aligned} \rho^{\tilde{\pi}_\gamma} &\geq (1-\gamma) \min_s V_{\gamma}^{\tilde{\pi}_\gamma}(s)\mathbf{1} \\ &\geq (1-\gamma) \min_s \hat{V}_{\gamma}^{\tilde{\pi}_\gamma}(s)\mathbf{1} - (1-\gamma) \left\| \hat{V}_{\gamma}^{\tilde{\pi}_\gamma} - V_{\gamma}^{\tilde{\pi}_\gamma} \right\|_{\infty} \mathbf{1} \\ &\geq (1-\gamma) \min_s \hat{V}_{\gamma}^{\star}(s)\mathbf{1} - \frac{1-\gamma}{n}\mathbf{1} - (1-\gamma) \left\| \hat{V}_{\gamma}^{\tilde{\pi}_\gamma} - V_{\gamma}^{\tilde{\pi}_\gamma} \right\|_{\infty} \mathbf{1} \\ &\geq (1-\gamma) \min_s \hat{V}_{\gamma}^{\star}(s)\mathbf{1} - \frac{1-\gamma}{n}\mathbf{1} - \alpha(\delta, n) \sqrt{\frac{\|\hat{V}_{\gamma}^{\tilde{\pi}_\gamma}\|_{\text{span}} + 1}{n}}\mathbf{1} \\ &\geq (1-\gamma) \min_s \tilde{V}_{\gamma}(s)\mathbf{1} - (1-\gamma) \left\| \hat{V}_{\gamma}^{\star} - \tilde{V}_{\gamma} \right\|_{\infty} \mathbf{1} - \frac{1-\gamma}{n}\mathbf{1} - \alpha(\delta, n) \sqrt{\frac{\|\hat{V}_{\gamma}^{\tilde{\pi}_\gamma}\|_{\text{span}} + 1}{n}}\mathbf{1} \\ &\geq (1-\gamma) \min_s \tilde{V}_{\gamma}(s)\mathbf{1} - 2\frac{1-\gamma}{n}\mathbf{1} - \alpha(\delta, n) \sqrt{\frac{\|\hat{V}_{\gamma}^{\tilde{\pi}_\gamma}\|_{\text{span}} + 1}{n}}\mathbf{1}. \end{aligned} \quad (60)$$

Now, nearly identically to the bound (43), we can use the requirements (56) and (57) to bound

$$\begin{aligned} \left\| \hat{V}_{\gamma}^{\tilde{\pi}_\gamma} \right\|_{\text{span}} &\leq \left\| \hat{V}_{\gamma}^{\star} \right\|_{\text{span}} + \left\| \hat{V}_{\gamma}^{\star} - \hat{V}_{\gamma}^{\tilde{\pi}_\gamma} \right\|_{\text{span}} \\ &\leq \left\| \tilde{V}_{\gamma} \right\|_{\text{span}} + \left\| \tilde{V}_{\gamma} - \hat{V}_{\gamma}^{\star} \right\|_{\text{span}} + \left\| \hat{V}_{\gamma}^{\star} - \hat{V}_{\gamma}^{\tilde{\pi}_\gamma} \right\|_{\text{span}} \\ &\leq \left\| \tilde{V}_{\gamma} \right\|_{\text{span}} + 2 \left\| \tilde{V}_{\gamma} - \hat{V}_{\gamma}^{\star} \right\|_{\infty} + \frac{1}{n} \\ &\leq \left\| \tilde{V}_{\gamma} \right\|_{\text{span}} + \frac{3}{n}. \end{aligned} \quad (61)$$

Finally combining (61) with (60), we obtain that

$$\rho^{\tilde{\pi}_\gamma} \geq (1-\gamma) \min_s \tilde{V}_{\gamma}(s)\mathbf{1} - 2\frac{1-\gamma}{n}\mathbf{1} - \alpha(\delta, n) \sqrt{\frac{\|\tilde{V}_{\gamma}\|_{\text{span}} + \frac{3}{n} + 1}{n}}\mathbf{1} = \hat{L}(\gamma)\mathbf{1}$$

as desired. ■

**Lemma 17** Under the event in Lemma 15, for all  $\gamma \in \mathcal{H}$ , we have

$$\left\| \widehat{V}_\gamma^\star - V_\gamma^\star \right\|_\infty \leq \frac{2\alpha(\delta, n)^2}{(1-\gamma)^2 n} + \frac{4\alpha(\delta, n)}{1-\gamma} \sqrt{\frac{\|V_\gamma^\star\|_{\text{span}} + 1 + \frac{1}{n}}{n}} + \frac{4}{n}. \quad (62)$$

**Proof** Similarly to Lemma 8, this follows immediately from Lemma 7 as its conditions are satisfied under the event in Lemma 15 (by setting  $m = n$  and  $\beta = \alpha(\delta, n)$  for all  $\gamma \in \mathcal{H}$ ).  $\blacksquare$

**Lemma 18** Suppose  $n \geq 16(\|h^\star\|_{\text{span}} + 1)$ . Under the event described in Lemma 15, we have

$$\widehat{L}(\gamma^\star) \mathbf{1} \geq \rho^\star - 30\alpha(\delta, n)^2 \sqrt{\frac{\|h^\star\|_{\text{span}} + 1}{n}} \mathbf{1}.$$

**Proof** First we bound  $\left\| \widehat{V}_\gamma^\star - V_\gamma^\star \right\|_\infty$  in terms of  $\|h^\star\|_{\text{span}}$ . Using the bound (62) from Lemma 17 and substituting  $\gamma = \gamma^\star$ , as well as using Lemma 14 to bound  $\frac{1}{1-\gamma^\star}$ , we have

$$\begin{aligned} \left\| \widehat{V}_{\gamma^\star}^\star - V_{\gamma^\star}^\star \right\|_\infty &\leq \frac{2\alpha(\delta, n)^2}{(1-\gamma^\star)^2 n} + \frac{4\alpha(\delta, n)}{1-\gamma^\star} \sqrt{\frac{\|V_{\gamma^\star}^\star\|_{\text{span}} + 1 + \frac{1}{n}}{n}} + \frac{4}{n} \\ &\leq \frac{2\alpha(\delta, n)^2}{n} 4n \left( \|h^\star\|_{\text{span}} + 1 \right) + 8\alpha(\delta, n) \sqrt{n \left( \|h^\star\|_{\text{span}} + 1 \right)} \sqrt{\frac{\|V_{\gamma^\star}^\star\|_{\text{span}} + 1 + \frac{1}{n}}{n}} + \frac{4}{n} \\ &\leq 8\alpha(\delta, n)^2 \left( \|h^\star\|_{\text{span}} + 1 \right) + 8\alpha(\delta, n) \sqrt{n \left( \|h^\star\|_{\text{span}} + 1 \right)} \sqrt{\frac{2\|h^\star\|_{\text{span}} + 2}{n}} + \frac{1}{4} \\ &\leq 8\alpha(\delta, n)^2 \left( \|h^\star\|_{\text{span}} + 1 \right) + (8\sqrt{2} + 1/4)\alpha(\delta, n) \left( \|h^\star\|_{\text{span}} + 1 \right) \\ &\leq 20\alpha(\delta, n)^2 \left( \|h^\star\|_{\text{span}} + 1 \right) \end{aligned} \quad (63)$$

also using the facts that  $n \geq 16$ , that  $\alpha(\delta, n) \geq 1$ , and that  $\|V_{\gamma^\star}^\star\|_{\text{span}} \leq 2\|h^\star\|_{\text{span}}$  (Wei et al., 2020, Lemma 2).

We can use the triangle inequality, that  $\|\cdot\|_{\text{span}} \leq 2\|\cdot\|_\infty$ , triangle inequality again, that  $\|V_{\gamma^\star}^\star\|_{\text{span}} \leq 2\|h^\star\|_{\text{span}}$ , (57) and (63), and then that  $\alpha(\delta, n) \geq 1$  and  $n \geq 16$  to bound

$$\begin{aligned} \left\| \widetilde{V}_{\gamma^\star} \right\|_{\text{span}} &\leq \|V_{\gamma^\star}^\star\|_{\text{span}} + \left\| \widetilde{V}_{\gamma^\star} - V_{\gamma^\star}^\star \right\|_{\text{span}} \\ &\leq \|V_{\gamma^\star}^\star\|_{\text{span}} + 2 \left\| \widetilde{V}_{\gamma^\star} - V_{\gamma^\star}^\star \right\|_\infty \\ &\leq 2\|h^\star\|_{\text{span}} + 2 \left\| \widetilde{V}_{\gamma^\star} - \widehat{V}_{\gamma^\star}^\star \right\|_\infty + 2 \left\| \widehat{V}_{\gamma^\star}^\star - V_{\gamma^\star}^\star \right\|_\infty \\ &\leq 2\|h^\star\|_{\text{span}} + \frac{2}{n} + 40\alpha(\delta, n)^2 \left( \|h^\star\|_{\text{span}} + 1 \right) \\ &\leq 42\alpha(\delta, n)^2 \left( \|h^\star\|_{\text{span}} + 1 \right). \end{aligned} \quad (64)$$

Now using the inequalities (63) and (64) to lower-bound  $\widehat{L}(\gamma^*)$ , we obtain

$$\begin{aligned}
\widehat{L}(\gamma^*)\mathbf{1} &= (1 - \gamma^*) \min_s \widetilde{V}_{\gamma^*}(s)\mathbf{1} - 2\frac{1 - \gamma^*}{n}\mathbf{1} - \alpha(\delta, n) \sqrt{\frac{\|\widetilde{V}_{\gamma^*}\|_{\text{span}} + \frac{3}{n} + 1}{n}}\mathbf{1} \\
&\geq (1 - \gamma^*) \min_s \widehat{V}_{\gamma^*}(s)\mathbf{1} - 3\frac{1 - \gamma^*}{n}\mathbf{1} - \alpha(\delta, n) \sqrt{\frac{\|\widetilde{V}_{\gamma^*}\|_{\text{span}} + \frac{3}{n} + 1}{n}}\mathbf{1} \\
&\geq (1 - \gamma^*) \min_s V_{\gamma^*}^*(s)\mathbf{1} - (1 - \gamma^*) \|\widehat{V}_{\gamma^*} - V_{\gamma^*}^*\|_{\infty} \mathbf{1} - 3\frac{1 - \gamma^*}{n}\mathbf{1} \\
&\quad - \alpha(\delta, n) \sqrt{\frac{\|\widetilde{V}_{\gamma^*}\|_{\text{span}} + \frac{3}{n} + 1}{n}}\mathbf{1} \\
&\geq (1 - \gamma^*) \min_s V_{\gamma^*}^*(s)\mathbf{1} - (1 - \gamma^*) 21\alpha(\delta, n)^2 \left( \|h^*\|_{\text{span}} + 1 \right) \mathbf{1} - 3\frac{1 - \gamma^*}{n}\mathbf{1} \\
&\quad - \alpha(\delta, n)^2 \sqrt{\frac{42 \left( \|h^*\|_{\text{span}} + 1 \right) + \frac{3}{n} + 1}{n}}\mathbf{1} \\
&\geq \rho^* - 2(1 - \gamma^*) \|h^*\|_{\text{span}} \mathbf{1} - (1 - \gamma^*) 21\alpha(\delta, n)^2 \left( \|h^*\|_{\text{span}} + 1 \right) \mathbf{1} - 3\frac{1 - \gamma^*}{n}\mathbf{1} \\
&\quad - \alpha(\delta, n)^2 \sqrt{\frac{44 \left( \|h^*\|_{\text{span}} + 1 \right)}{n}}\mathbf{1} \\
&\geq \rho^* - (1 - \gamma^*) 22\alpha(\delta, n)^2 \left( \|h^*\|_{\text{span}} + 1 \right) \mathbf{1} - \alpha(\delta, n)^2 \sqrt{\frac{44 \left( \|h^*\|_{\text{span}} + 1 \right)}{n}}\mathbf{1} \\
&\geq \rho^* - \frac{1}{\sqrt{n \left( \|h^*\|_{\text{span}} + 1 \right)}} 22\alpha(\delta, n)^2 \left( \|h^*\|_{\text{span}} + 1 \right) \mathbf{1} - \alpha(\delta, n)^2 \sqrt{\frac{44 \left( \|h^*\|_{\text{span}} + 1 \right)}{n}}\mathbf{1} \\
&\geq \rho^* - 30\alpha(\delta, n)^2 \sqrt{\frac{\|h^*\|_{\text{span}} + 1}{n}}\mathbf{1}
\end{aligned}$$

where in the first inequality we used that  $\|\widehat{V}_{\gamma^*} - \widetilde{V}_{\gamma^*}\|_{\infty} \leq \frac{1}{n}$  from (57) and triangle inequality, in the second inequality we used triangle inequality, in the third we used (63) and (64), in the fourth we used that  $\min_s V_{\gamma^*}^*(s) \geq \frac{1}{1 - \gamma^*} \rho^* - \|h^*\|_{\text{span}}$  from (Wei et al., 2020, Lemma 2) and that  $n \geq 16$ , in the fifth we again used that  $n \geq 16$ , and in the sixth we used (54) to upper bound  $(1 - \gamma^*)$ . ■

Now we combine these intermediate results to prove Theorem 1.

**Proof of Theorem 1** Under the event in Lemma 15 which holds with probability at least  $1 - \delta$ , we have by Lemma 16 that  $\rho^{\widehat{\gamma}} \geq \widehat{L}(\widehat{\gamma})\mathbf{1}$ . Since we trivially have  $\rho^{\widehat{\gamma}} \geq 0\mathbf{1}$ , this implies  $\rho^{\widehat{\gamma}} \geq \max\{\widehat{L}(\widehat{\gamma}), 0\}\mathbf{1} = \widehat{\rho}$ . Now to lower bound  $\widehat{\rho}$  we consider two cases. First, if  $n \geq 16(\|h^*\|_{\text{span}} + 1)$ ,



then we have

$$\hat{\rho} = \max\{\hat{L}(\hat{\gamma}), 0\} \mathbf{1} \geq \hat{L}(\hat{\gamma}) \geq \hat{L}(\gamma^*) \mathbf{1} \geq \rho^* - 30\alpha(\delta, n)^2 \sqrt{\frac{\|h^*\|_{\text{span}} + 1}{n}} \mathbf{1}$$

where the second inequality step is by the definition of  $\hat{\gamma}$ , and third inequality is from Lemma 18. Next, if  $n < 16(\|h^*\|_{\text{span}} + 1)$ , then

$$\hat{\rho} = \max\{\hat{L}(\hat{\gamma}), 0\} \mathbf{1} \geq 0 \mathbf{1} \geq \rho^* - 30\alpha(\delta, n)^2 \sqrt{\frac{\|h^*\|_{\text{span}} + 1}{n}} \mathbf{1}$$

where the final inequality is because  $\rho^* \leq 1$  so the condition on  $n$  ensures the RHS is  $< 0$  (note  $\alpha(\delta, n) \geq 1$ ). Finally, we let  $C_3 = 30$ .  $\blacksquare$

### C.3. Proof of Lemma 3

**Proof of Lemma 3** First we check that  $T$  ensures that  $\gamma^T \leq \frac{(1-\gamma)^{2\varepsilon}}{3}$ . We have

$$T \geq \frac{\log(\frac{3}{(1-\gamma)^{2\varepsilon}})}{1-\gamma} \geq \frac{\log(\frac{3}{(1-\gamma)^{2\varepsilon}})}{\log(1/\gamma)}$$

using the fact that  $\frac{1}{1-\gamma} \geq \frac{1}{\log(1/\gamma)}$  for  $\gamma \in (0, 1)$ . Thus

$$\gamma^T \leq \gamma^{\frac{\log(\frac{3}{(1-\gamma)^{2\varepsilon}})}{\log(1/\gamma)}} = (e^{-1})^{\log(\frac{3}{(1-\gamma)^{2\varepsilon}})} = \frac{(1-\gamma)^{2\varepsilon}}{3}. \quad (65)$$

By (Fruit et al., 2018, Lemma 16)  $\text{Clip}_M$  is non-expansive with respect to  $\|\cdot\|_\infty$ , so since  $\mathcal{T}_\gamma$  is  $\gamma$ -contractive, we have that  $\mathcal{L} = \text{Clip}_M \circ \mathcal{T}_\gamma$  is also  $\gamma$ -contractive. It is a standard fact that this implies  $\mathcal{L}$  has a unique fixed point, which we name  $V_{\gamma, M}^*$ .

Thus

$$\begin{aligned} \|V^T - V_{\gamma, M}^*\|_\infty &= \|\mathcal{L}(V^{T-1}) - \mathcal{L}(V_{\gamma, M}^*)\|_\infty \\ &\leq \gamma \|V^{T-1} - V_{\gamma, M}^*\|_\infty \leq \dots \leq \gamma^T \|V^0 - V_{\gamma, M}^*\|_\infty \leq \gamma^T \frac{1}{1-\gamma}, \end{aligned}$$

which is  $\leq \frac{(1-\gamma)\varepsilon}{3} \leq \varepsilon$  using (65). Similarly, we have

$$\begin{aligned} \|V^T - \mathcal{L}(V^T)\|_\infty &= \|\mathcal{L}(V^{T-1}) - \mathcal{L}(V^T)\|_\infty \leq \gamma \|V^{T-1} - V^T\|_\infty = \gamma \|V^{T-1} - \mathcal{L}(V^{T-1})\|_\infty \\ &\leq \dots \leq \gamma^T \|V^0 - \mathcal{L}(V^0)\|_\infty \leq \gamma^T \frac{1}{1-\gamma} \leq \frac{(1-\gamma)\varepsilon}{3} \end{aligned}$$

again using (65). By definition of  $\tilde{r}$ , we immediately have that  $\tilde{r} \leq r$ . Also by construction of  $\tilde{r}$  we have that

$$\tilde{r}_{\hat{\pi}} + \gamma P_{\hat{\pi}} V^T = \min \left\{ r_{\hat{\pi}} + \gamma P_{\hat{\pi}} V^T, M \mathbf{1} + \min_{s'} V^T(s') \mathbf{1} \right\},$$

where the min is elementwise. Thus

$$\begin{aligned}
& \|\mathcal{L}(V^T) - \tilde{r}_{\hat{\pi}} - \gamma P_{\hat{\pi}} V^T\|_{\infty} \\
&= \left\| \min \left\{ r_{\hat{\pi}} + \gamma P_{\hat{\pi}} V^T, M\mathbf{1} + \min_{s'} (r_{\hat{\pi}} + \gamma P_{\hat{\pi}} V^T)(s')\mathbf{1} \right\} - \min \left\{ r_{\hat{\pi}} + \gamma P_{\hat{\pi}} V^T, M\mathbf{1} + \min_{s'} V^T(s')\mathbf{1} \right\} \right\|_{\infty} \\
&\leq \left\| M\mathbf{1} + \min_{s'} (r_{\hat{\pi}} + \gamma P_{\hat{\pi}} V^T)(s')\mathbf{1} - \left( M\mathbf{1} + \min_{s'} V^T(s')\mathbf{1} \right) \right\|_{\infty} \\
&= \left| \min_{s'} (r_{\hat{\pi}} + \gamma P_{\hat{\pi}} V^T)(s') - \min_{s'} V^T(s') \right| \\
&= \left| \min_{s'} \text{Clip}_M(r_{\hat{\pi}} + \gamma P_{\hat{\pi}} V^T)(s') - \min_{s'} V^T(s') \right| \\
&= \left| \min_{s'} \text{Clip}_M(\mathcal{T}_{\gamma}(V^T))(s') - \min_{s'} V^T(s') \right| \\
&= \left| \min_{s'} \mathcal{L}(V^T)(s') - \min_{s'} V^T(s') \right| \\
&\leq \|\mathcal{L}(V^T) - V^T\|_{\infty}, \tag{66}
\end{aligned}$$

where we use the elementary fact that  $|\min\{a, b\} - \min\{c, d\}| \leq \max\{|a - c|, |b - d|\}$  for  $a, b, c, d \in \mathbb{R}$ , as well as the fact that  $\text{Clip}_M$  does not change the minimum entry of its input.

Now we can calculate that

$$\begin{aligned}
\|V_{\gamma, \tilde{r}}^{\hat{\pi}} - V^T\|_{\infty} &\leq \|V_{\gamma, \tilde{r}}^{\hat{\pi}} - \mathcal{L}(V^T)\|_{\infty} + \|\mathcal{L}(V^T) - V^T\|_{\infty} \\
&\leq \|V_{\gamma, \tilde{r}}^{\hat{\pi}} - \tilde{r}_{\hat{\pi}} - \gamma P_{\hat{\pi}} V^T\|_{\infty} + \|\mathcal{L}(V^T) - \tilde{r}_{\hat{\pi}} - \gamma P_{\hat{\pi}} V^T\|_{\infty} + \|\mathcal{L}(V^T) - V^T\|_{\infty} \\
&= \|\tilde{r}_{\hat{\pi}} - \gamma P_{\hat{\pi}} V_{\gamma, \tilde{r}}^{\hat{\pi}} - \tilde{r}_{\hat{\pi}} - \gamma P_{\hat{\pi}} V^T\|_{\infty} + \|\mathcal{L}(V^T) - \tilde{r}_{\hat{\pi}} - \gamma P_{\hat{\pi}} V^T\|_{\infty} + \|\mathcal{L}(V^T) - V^T\|_{\infty} \\
&\stackrel{\text{(I)}}{\leq} \|\tilde{r}_{\hat{\pi}} - \gamma P_{\hat{\pi}} V_{\gamma, \tilde{r}}^{\hat{\pi}} - \tilde{r}_{\hat{\pi}} - \gamma P_{\hat{\pi}} V^T\|_{\infty} + 2\|\mathcal{L}(V^T) - V^T\|_{\infty} \\
&\stackrel{\text{(II)}}{\leq} \gamma \|V_{\gamma, \tilde{r}}^{\hat{\pi}} - V^T\|_{\infty} + 2\|\mathcal{L}(V^T) - V^T\|_{\infty},
\end{aligned}$$

where in (I) we used (66) and in (II) we used that  $\|P_{\hat{\pi}}\|_{\infty \rightarrow \infty} \leq 1$ . Thus by rearranging we have that

$$\|V_{\gamma, \tilde{r}}^{\hat{\pi}} - V^T\|_{\infty} \leq \frac{2\|\mathcal{L}(V^T) - V^T\|_{\infty}}{1 - \gamma} \leq \frac{2}{3}\varepsilon$$

using (65). Then we can bound

$$\|V_{\gamma, \tilde{r}}^{\hat{\pi}} - V_{\gamma, M}^{\star}\|_{\infty} \leq \|V_{\gamma, \tilde{r}}^{\hat{\pi}} - V^T\|_{\infty} + \|V^T - V_{\gamma, M}^{\star}\|_{\infty} \leq \frac{2}{3}\varepsilon + (1 - \gamma)\frac{\varepsilon}{3} \leq \varepsilon$$

as desired. Also,  $\|V_{\gamma, M}^{\star}\|_{\text{span}} \leq M$  since it is equal to  $\text{Clip}_M(\mathcal{T}_{\gamma}(V_{\gamma, M}^{\star}))$  and the output of  $\text{Clip}_M$  clearly has span bounded by  $M$ . Therefore

$$\|V_{\gamma, \tilde{r}}^{\hat{\pi}}\|_{\text{span}} \leq \|V_{\gamma, M}^{\star}\|_{\text{span}} + \|V_{\gamma, \tilde{r}}^{\hat{\pi}} - V_{\gamma, M}^{\star}\|_{\text{span}} \leq \|V_{\gamma, M}^{\star}\|_{\text{span}} + 2\|V_{\gamma, \tilde{r}}^{\hat{\pi}} - V_{\gamma, M}^{\star}\|_{\infty} \leq M + 2\varepsilon.$$

Finally, letting  $\mathcal{T}_{\gamma,r'}^{\pi'}$  be the Bellman consistency/evaluation operator for the policy  $\pi'$  in the DMDP  $(P, r', \gamma)$ , we have  $\mathcal{T}_{\gamma,r'}^{\pi'}(V_{\gamma,r'}^{\pi'}) = V_{\gamma,r'}^{\pi'}$ . Also, since  $\|V_{\gamma,r'}^{\pi'}\|_{\text{span}} \leq M$ , we have  $\text{Clip}_M(\mathcal{T}_{\gamma,r'}^{\pi'}(V_{\gamma,r'}^{\pi'})) = V_{\gamma,r'}^{\pi'}$ . Letting  $\mathcal{L}^{\pi'} = \text{Clip}_M \circ \mathcal{T}_{\gamma,r'}^{\pi'}$ , we have that  $\mathcal{L}^{\pi'}$  is a  $\gamma$  contraction with respect to  $\|\cdot\|_{\infty}$  (since as discussed above,  $\text{Clip}_M$  is  $\|\cdot\|_{\infty}$ -nonexpansive, and because  $\mathcal{T}_{\gamma,r'}^{\pi'}$  is well-known to be  $\gamma$ -contractive with respect to  $\|\cdot\|_{\infty}$ ). Thus  $\mathcal{L}^{\pi'}$  has a unique fixed point, which must be  $V_{\gamma,r'}^{\pi'}$  (since we have already verified  $V_{\gamma,r'}^{\pi'}$  is a fixed point), and furthermore Picard iteration will converge to this fixed point. The operator  $\text{Clip}_M$  is also monotonic (in the sense that for any  $x, y \in \mathbb{R}^S$ ,  $x \leq y \implies \text{Clip}_M(x) \leq \text{Clip}_M(y)$ ) as shown in (Fruit et al., 2018, Lemma 16). It is well-known that  $\mathcal{T}_{\gamma}$  and  $\mathcal{T}_{\gamma,r'}^{\pi'}$  are also monotonic, which immediately implies that the compositions  $\mathcal{L}$  and  $\mathcal{L}^{\pi'}$  are monotonic. It is immediate that for any  $x \in \mathbb{R}^S$  we have  $\mathcal{T}_{\gamma,r'}^{\pi'}(x) \leq \mathcal{T}_{\gamma}(x)$ , since for any  $x \in \mathbb{R}^S$  we have

$$\mathcal{T}_{\gamma,r'}^{\pi'}(x) = M^{\pi}(r' + \gamma Px) \leq M^{\pi}(r + \gamma Px) \leq M(r + \gamma Px) = \mathcal{T}_{\gamma}(x)$$

crucially using the fact that  $r' \leq r$  and also monotonicity of  $M^{\pi}$ . This thus implies that  $\mathcal{L}^{\pi'}(x) = \text{Clip}_M(\mathcal{T}_{\gamma,r'}^{\pi'}(x)) \leq \text{Clip}_M(\mathcal{T}_{\gamma}(x)) \leq \mathcal{L}(x)$  for any  $x \in \mathbb{R}^S$  (using monotonicity of  $\text{Clip}_M$ ). Therefore  $\mathcal{L}^{\pi'}(\mathbf{0}) \leq \mathcal{L}(\mathbf{0})$ , and if it holds for some integer  $i \geq 1$  that  $(\mathcal{L}^{\pi'})^{(i)}(\mathbf{0}) \leq \mathcal{L}^{(i)}(\mathbf{0})$  then we can use the above-discussed properties to obtain that

$$(\mathcal{L}^{\pi'})^{(i+1)}(\mathbf{0}) = \mathcal{L}^{\pi'}\left((\mathcal{L}^{\pi'})^{(i)}(\mathbf{0})\right) \leq \mathcal{L}\left((\mathcal{L}^{\pi'})^{(i)}(\mathbf{0})\right) \leq \mathcal{L}\left(\mathcal{L}^{(i)}(\mathbf{0})\right) = \mathcal{L}^{(i+1)}(\mathbf{0}),$$

so we have by induction that  $(\mathcal{L}^{\pi'})^{(i)}(\mathbf{0}) \leq \mathcal{L}^{(i)}(\mathbf{0})$  holds for all  $i$ , and thus

$$V_{\gamma,r'}^{\pi'} = \lim_{i \rightarrow \infty} (\mathcal{L}^{\pi'})^{(i)}(\mathbf{0}) \leq \lim_{i \rightarrow \infty} \mathcal{L}^{(i)}(\mathbf{0}) = V_{\gamma,M}^{\star}.$$

Thus we have

$$V_{\gamma}^{\widehat{\pi}} \geq V_{\gamma,\tilde{r}}^{\widehat{\pi}} \geq V_{\gamma,M}^{\star} - \varepsilon \mathbf{1} \geq V_{\gamma,r'}^{\pi'} - \varepsilon \mathbf{1}$$

as desired, where the first inequality is because  $r \geq \tilde{r}$ , the second inequality uses that  $\|V_{\gamma,\tilde{r}}^{\widehat{\pi}} - V_{\gamma,M}^{\star}\|_{\infty} \leq \varepsilon$ , and the final inequality uses the above argument.  $\blacksquare$

#### C.4. Proof of Theorem 4

First we develop bounds on the error between value functions within  $P$  and  $\widehat{P}$ .

**Lemma 19** *Fix a policy  $\pi$ ,  $\gamma \in (0, 1)$ , and  $\delta, n > 0$ . Then with probability at least  $1 - \delta$ , we have that*

$$\|V_{\gamma}^{\pi} - \widehat{V}_{\gamma}^{\pi}\|_{\infty} \leq \frac{24 \log_2 \log_2(\frac{1}{1-\gamma} + 4)}{1 - \gamma} \sqrt{\frac{\|V_{\gamma}^{\pi}\|_{\text{span}} + 1}{n}} 16 \log \left( \frac{12SA n}{(1 - \gamma)^2 \delta} \right).$$

**Proof** We can reuse the proof of (Zurek and Chen, 2024, Theorem 9). Specifically, (Zurek and Chen, 2024, Equation 31) is stated for an optimal policy in the DMDP  $(P, r, \gamma)$ , but completely identical

arguments actually hold for any fixed policy  $\pi$ . Therefore (Zurek and Chen, 2024, Equation 31) yields that with probability at least  $1 - \delta$  it holds that

$$\|V_\gamma^\pi - \widehat{V}_\gamma^\pi\|_\infty \leq \frac{24 \log_2 \log_2(\frac{1}{1-\gamma} + 4)}{1-\gamma} \sqrt{\frac{\|V_\gamma^\pi\|_{\text{span}} + 1}{n} 16 \log \left( \frac{12SA n}{(1-\gamma)^2 \delta} \right)}.$$

This completes the proof. ■

We would like to have a similar bound involving policies output by the span-constrained planning procedure applied to  $\widehat{P}$ , but since such policies are probabilistically dependent on  $\widehat{P}$ , much more effort is needed. In the absence of span-constrained planning, Agarwal et al. (2020) introduce the absorbing MDP construction to overcome such issues, and our approach is strongly inspired by theirs, although ours requires new modifications which are specific to the span-constrained planning procedure Algorithm 3.

First we must define the key quantities involved in our construction. Fix  $s \in \mathcal{S}$  and  $u \in [0, 1]$ . We define the MDP transition kernel  $\widehat{P}^{(s)}$  as

$$\widehat{P}^{(s)}(s' | s'', a) = \begin{cases} \widehat{P}(s' | s'', a) & s'' \neq s \\ 1 & s'' = s' = s \\ 0 & s'' = s \text{ and } s' \neq s \end{cases}.$$

In words,  $\widehat{P}^{(s)}$  is identical to  $\widehat{P}$  except state  $s$  is made to be absorbing. We also define the reward function  $r^{(s,u)}$  as

$$r^{(s,u)}(s', a) = \begin{cases} r(s', a) & s' \neq s \\ u & s' = s \end{cases}.$$

In words,  $r^{(s,u)}$  is identical to  $r$  except state  $s$  gives reward  $u$  for all actions. For any  $\gamma \in (0, 1)$  and any  $M > 0$ , let  $\widehat{\mathcal{T}}_\gamma^{(s,u)} : \mathbb{R}^{\mathcal{S}} \rightarrow \mathbb{R}^{\mathcal{S}}$  be the  $\gamma$ -discounted Bellman operator for the DMDP  $(\widehat{P}^{(s)}, r^{(s,u)}, \gamma)$ , that is,

$$\widehat{\mathcal{T}}_\gamma^{(s,u)}(x) = M(r^{(s,u)} + \gamma \widehat{P}^{(s)}x)$$

for all  $x \in \mathbb{R}^{\mathcal{S}}$ . Let  $\widehat{V}_{\gamma,M}^{(s,u)}$  be the unique fixed point of  $\text{Clip}_M \circ \widehat{\mathcal{T}}_\gamma^{(s,u)}$  (this fixed point exists and is unique because  $\text{Clip}_M \circ \widehat{\mathcal{T}}_\gamma^{(s,u)}$  is a  $\gamma$ -contraction due to Lemma 3).

Now we can summarize the key properties of this construction in the following lemma.

**Lemma 20** *Fixing  $\gamma \in (0, 1)$  and  $M > 0$ , for any  $u, u' \in [0, 1]$  we have*

$$\|\widehat{V}_{\gamma,M}^{(s,u)} - \widehat{V}_{\gamma,M}^{(s,u')}\|_\infty \leq \frac{|u - u'|}{1 - \gamma}.$$

Also, letting  $u^*(s) = \widehat{\mathcal{T}}_\gamma(\widehat{V}_{\gamma,M}^*)(s) - \gamma \widehat{V}_{\gamma,M}^*(s)$ , we have that  $u^*(s) \in [0, 1]$  and

$$\widehat{V}_{\gamma,M}^{(s,u^*(s))} = \widehat{V}_{\gamma,M}^*,$$

where  $\widehat{V}_{\gamma,M}^*$  is the unique fixed point of  $\text{Clip}_M \circ \widehat{\mathcal{T}}_\gamma$ .

Consequently, there exists a finite set  $U$  with  $|U| = \left\lceil \frac{1}{2(1-\gamma)\varepsilon} \right\rceil$  such that almost surely, for any  $s \in \mathcal{S}$  there exists  $u \in U$  such that  $\left\| \widehat{V}_{\gamma,M}^* - \widehat{V}_{\gamma,M}^{(s,u)} \right\|_\infty \leq \varepsilon$ .

Lastly, for any  $u \in [0, 1]$ ,  $s \in \mathcal{S}$ , and  $a \in \mathcal{A}$ ,  $\widehat{V}_{\gamma,M}^{(s,u)}$  is independent of the samples  $S_{s,a}^1, \dots, S_{s,a}^n$  used to construct  $\widehat{P}(\cdot \mid s, a)$ .

**Proof** For the first inequality, we can calculate that

$$\begin{aligned}
\left\| \widehat{V}_{\gamma,M}^{(s,u)} - \widehat{V}_{\gamma,M}^{(s,u')} \right\|_\infty &= \left\| \text{Clip}_M \left( \widehat{\mathcal{T}}_\gamma^{(s,u)} \left( \widehat{V}_{\gamma,M}^{(s,u)} \right) \right) - \text{Clip}_M \left( \widehat{\mathcal{T}}_\gamma^{(s,u')} \left( \widehat{V}_{\gamma,M}^{(s,u')} \right) \right) \right\|_\infty \\
&\leq \left\| \widehat{\mathcal{T}}_\gamma^{(s,u)} \left( \widehat{V}_{\gamma,M}^{(s,u)} \right) - \widehat{\mathcal{T}}_\gamma^{(s,u')} \left( \widehat{V}_{\gamma,M}^{(s,u')} \right) \right\|_\infty \\
&= \left\| M \left( r^{(s,u)} + \gamma \widehat{P}^{(s)} \widehat{V}_{\gamma,M}^{(s,u)} \right) - M \left( r^{(s,u')} + \gamma \widehat{P}^{(s)} \widehat{V}_{\gamma,M}^{(s,u')} \right) \right\|_\infty \\
&\leq \left\| r^{(s,u)} - r^{(s,u')} + \gamma \widehat{P}^{(s)} \left( \widehat{V}_{\gamma,M}^{(s,u)} - \widehat{V}_{\gamma,M}^{(s,u')} \right) \right\|_\infty \\
&\leq \left\| r^{(s,u)} - r^{(s,u')} \right\|_\infty + \gamma \left\| \widehat{V}_{\gamma,M}^{(s,u)} - \widehat{V}_{\gamma,M}^{(s,u')} \right\|_\infty \\
&= |u - u'| + \gamma \left\| \widehat{V}_{\gamma,M}^{(s,u)} - \widehat{V}_{\gamma,M}^{(s,u')} \right\|_\infty.
\end{aligned}$$

Rearranging, we obtain that

$$\left\| \widehat{V}_{\gamma,M}^{(s,u)} - \widehat{V}_{\gamma,M}^{(s,u')} \right\|_\infty \leq \frac{|u - u'|}{1 - \gamma}$$

as desired.

Now we show the second statement in the lemma. First we show that  $\widehat{V}_{\gamma,M}^{(s,u^*(s))} = \widehat{V}_{\gamma,M}^*$ . To show this, it suffices to show that

$$\widehat{\mathcal{T}}_\gamma \left( \widehat{V}_{\gamma,M}^* \right) = \widehat{\mathcal{T}}_\gamma^{(s,u^*(s))} \left( \widehat{V}_{\gamma,M}^* \right), \quad (67)$$

since this implies that

$$\widehat{V}_{\gamma,M}^* = \text{Clip}_M \left( \widehat{\mathcal{T}}_\gamma \left( \widehat{V}_{\gamma,M}^* \right) \right) = \text{Clip}_M \left( \widehat{\mathcal{T}}_\gamma^{(s,u^*(s))} \left( \widehat{V}_{\gamma,M}^* \right) \right)$$

(using the fact that  $\widehat{V}_{\gamma,M}^*$  is a fixed point of  $\text{Clip}_M \circ \widehat{\mathcal{T}}_\gamma$  in the first equality), meaning that  $\widehat{V}_{\gamma,M}^*$  is a fixed point of  $\text{Clip}_M \circ \widehat{\mathcal{T}}_\gamma^{(s,u^*(s))}$ , and since the unique fixed point of  $\text{Clip}_M \circ \widehat{\mathcal{T}}_\gamma^{(s,u^*(s))}$  is  $\widehat{V}_{\gamma,M}^{(s,u^*(s))}$ , this would imply that  $\widehat{V}_{\gamma,M}^{(s,u^*(s))} = \widehat{V}_{\gamma,M}^*$ . Now focusing on (67), we note that by construction of  $\widehat{\mathcal{T}}_\gamma^{(s,u^*(s))}$  we already have  $\widehat{\mathcal{T}}_\gamma \left( \widehat{V}_{\gamma,M}^* \right) (s') = \widehat{\mathcal{T}}_\gamma^{(s,u^*(s))} \left( \widehat{V}_{\gamma,M}^* \right) (s')$  for all  $s' \neq s$ , so it remains to show the equality for the  $s$ th coordinates. We have the desired equality

$$\widehat{\mathcal{T}}_\gamma^{(s,u^*(s))} \left( \widehat{V}_{\gamma,M}^* \right) (s) = u^*(s) + \gamma \widehat{V}_{\gamma,M}^*(s) = \widehat{\mathcal{T}}_\gamma \left( \widehat{V}_{\gamma,M}^* \right) (s)$$

using the definition of  $\widehat{\mathcal{T}}_\gamma^{(s,u^*(s))}$  and then the definition of  $u^*(s)$ , so (67) holds and thus  $\widehat{V}_{\gamma,M}^{(s,u^*(s))} = \widehat{V}_{\gamma,M}^*$ .

Next we check that  $u^*(s) \in [0, 1]$ . We consider two cases, either  $\widehat{\mathcal{T}}_\gamma(\widehat{V}_{\gamma,M}^*)(s) = \widehat{V}_{\gamma,M}^*(s)$  (the clipping operator does not affect entry  $s$ ) or  $\widehat{\mathcal{T}}_\gamma(\widehat{V}_{\gamma,M}^*)(s) > \widehat{V}_{\gamma,M}^*(s)$ . In the first case we immediately have

$$u^*(s) = \widehat{\mathcal{T}}_\gamma(\widehat{V}_{\gamma,M}^*)(s) - \gamma \widehat{V}_{\gamma,M}^*(s) = (1 - \gamma) \widehat{V}_{\gamma,M}^*(s)$$

which is  $\in [0, 1]$  since  $\widehat{V}_{\gamma,M}^*(s) \in [0, \frac{1}{1-\gamma}]$ . (This fact can be seen by noting that  $\widehat{\mathcal{T}}_\gamma(x) \geq (\text{Clip}_M \circ \widehat{\mathcal{T}}_\gamma)(x)$  for any  $x$  so by a monotonicity argument we have that  $\widehat{V}_{\gamma,M}^* \leq \widehat{V}_\gamma^* \leq \frac{1}{1-\gamma} \mathbf{1}$ .) In the second case, since the clipping affects entry  $s$ , this means that  $\widehat{V}_{\gamma,M}^*(s)$  must be the largest (not necessarily the uniquely largest) entry of  $\widehat{V}_{\gamma,M}^*$ . Therefore we have by the definition of  $\widehat{\mathcal{T}}_\gamma(\widehat{V}_{\gamma,M}^*)$  that for some  $a^* \in \mathcal{A}$ ,

$$\widehat{\mathcal{T}}_\gamma(\widehat{V}_{\gamma,M}^*)(s) = r(s, a^*) + \gamma \widehat{P}_{s,a^*} \widehat{V}_{\gamma,M}^* \leq r(s, a^*) + \gamma \widehat{V}_{\gamma,M}^*(s)$$

since  $\widehat{P}_{s,a^*} \widehat{V}_{\gamma,M}^*$  is an average of entries of  $\widehat{V}_{\gamma,M}^*$  which must be less than the largest entry  $\widehat{V}_{\gamma,M}^*(s)$ . After rearranging we have that

$$u^*(s) = \widehat{\mathcal{T}}_\gamma(\widehat{V}_{\gamma,M}^*)(s) - \gamma \widehat{V}_{\gamma,M}^*(s) \leq r(s, a^*) \leq 1,$$

and to lower bound  $u^*(s)$  we can simply use that in this case we have assumed  $\widehat{\mathcal{T}}_\gamma(\widehat{V}_{\gamma,M}^*)(s) > \widehat{V}_{\gamma,M}^*(s)$ , so

$$u^*(s) = \widehat{\mathcal{T}}_\gamma(\widehat{V}_{\gamma,M}^*)(s) - \gamma \widehat{V}_{\gamma,M}^*(s) > (1 - \gamma) \widehat{V}_{\gamma,M}^*(s) \geq 0.$$

Now we show the penultimate statement. Letting  $U$  be a set of  $\left\lceil \frac{1}{2(1-\gamma)\varepsilon} \right\rceil$  equally spaced points in  $[0, 1]$ , for any  $z \in [0, 1]$  there exists  $u \in U$  such that  $|z - u| \leq (1 - \gamma)\varepsilon$ . Therefore for any  $s$ , letting  $U(s)$  be the closest element of  $U$  to  $u^*(s)$ , we have that

$$\left\| \widehat{V}_{\gamma,M}^* - \widehat{V}_{\gamma,M}^{(s,U(s))} \right\|_\infty = \left\| \widehat{V}_{\gamma,M}^{(s,u^*(s))} - \widehat{V}_{\gamma,M}^{(s,U(s))} \right\|_\infty \leq \frac{|u^*(s) - U(s)|}{1 - \gamma} \leq \varepsilon$$

as desired.

Finally, for any  $s \in \mathcal{S}, a \in \mathcal{A}, u \in [0, 1]$ , the independence of  $\widehat{V}_{\gamma,M}^{(s,u)}$  from the samples  $S_{s,a}^1, \dots, S_{s,a}^n$  is immediate from the construction, since  $\widehat{V}_{\gamma,M}^{(s,u)}$  uses the transition kernel  $\widehat{P}^{(s)}$  which is independent of  $S_{s,a}^1, \dots, S_{s,a}^n$ .  $\blacksquare$

Now we can use this construction to prove the desired error bound. First we state the desired bound:

**Lemma 21** Fix  $\gamma \in (0, 1)$  and  $\delta, M, n > 0$ . Let  $\widehat{\pi}, \widetilde{V}, \widetilde{r}$  be the output of the Span-Constrained Planning Algorithm 3 with inputs  $(\widehat{P}, r, \gamma), M, \frac{1}{n}$ . Then with probability at least  $1 - \delta$ , we have that

$$\left\| V_{\gamma, \widetilde{r}}^{\widehat{\pi}} - \widehat{V}_{\gamma, \widetilde{r}}^{\widehat{\pi}} \right\|_\infty \leq \frac{24 \log_2 \log_2 \left( \frac{1}{1-\gamma} + 4 \right)}{1 - \gamma} \sqrt{\frac{\left\| \widehat{V}_{\gamma, \widetilde{r}}^{\widehat{\pi}} \right\|_{\text{span}} + 1}{n}} 16 \log \left( \frac{12 S A n}{(1 - \gamma)^2 \delta} \right).$$

At a high level, to prove Lemma 21, we take advantage of the formal similarity between our absorbing MDP construction, designed for span-constrained MDPs, and the original construction of Agarwal et al. (2020). Specifically, we are able to reuse certain proof steps from Zurek and Chen (2024) which utilize the absorbing MDP construction of Agarwal et al. (2020), since although their statements involve different objects, their proofs utilize certain properties of said objects which all hold in an identical manner for the objects we consider. First, by examining the specific properties used within the proof of (Zurek and Chen, 2024, Lemma 20), the following result is actually shown:

**Lemma 22** *Suppose there exists a random variable  $\hat{V}^*$ , a set  $U$ , as well as random variables  $(\hat{V}_{s,u}^*)_{s \in \mathcal{S}, u \in U}$  such that*

1. *For any  $s \in \mathcal{S}$ ,  $u \in U$ , and  $a \in \mathcal{A}$ ,  $\hat{V}_{s,u}^*$  is independent of the samples  $S_{s,a}^1, \dots, S_{s,a}^n$  used to construct  $\hat{P}(\cdot \mid s, a)$ .*
2. *Almost surely, for any  $s \in \mathcal{S}$ , there exists  $u(s) \in U$  such that  $\left\| \hat{V}^* - \hat{V}_{s,u(s)}^* \right\|_\infty \leq \frac{1}{n}$ .*
3.  $|U| \leq \left\lceil \frac{n}{2(1-\gamma)} \right\rceil$ .

Also suppose  $n \geq 4$ . With probability at least  $1 - \delta$ , the following holds: for all policies  $\pi$  and all  $X \in \mathbb{R}^{\mathcal{S}}$  which satisfy  $\left\| X - \hat{V}^* \right\|_\infty \leq \frac{1}{n}$  and that  $\|X\|_{\text{span}} \leq \frac{1}{1-\gamma}$ , letting  $\bar{V} = X - (\min_s X(s)) \mathbf{1}$ , for all  $k = 0, \dots, \left\lceil \log_2 \log_2 \left( \|X\|_{\text{span}} + 4 \right) \right\rceil$ , we have

$$\left| \left( \hat{P}_\pi - P_\pi \right) (\bar{V})^{\circ 2^k} \right| \leq \sqrt{\frac{\beta \mathbb{V}_{\hat{P}_\pi} \left[ (\bar{V})^{\circ 2^k} \right]}{n}} + \frac{\beta \cdot 2^k}{n} (\|\bar{V}\|_\infty + 1)^{2^k} \mathbf{1} \quad (68)$$

where  $\beta = 16 \log \left( 12 \frac{SAn}{(1-\gamma)^{2\delta}} \right)$ .

Here  $x^{\circ p}$  denotes the elementwise  $p$ th power of any vector  $x$ . With Lemma 22, as well as another key result from Zurek and Chen (2024) which utilizes a conclusion in the form of Lemma 22 to prove an error bound between value functions, we can now prove Lemma 21.

**Proof of Lemma 21** By Lemma 20, it is immediate that the assumptions of Lemma 22 are satisfied for  $\hat{V}^* \leftarrow \hat{V}_{\gamma, M}^*$  and  $\hat{V}_{s,u}^* \leftarrow \hat{V}_{\gamma, M}^{(s,u)}$ . We would like to obtain the conclusion (68) for  $X \leftarrow \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}}$ , which we now argue satisfies the assumptions. First, by Lemma 3, we have that (almost surely)  $\left\| \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}} - \hat{V}_{\gamma, M}^* \right\|_\infty \leq \frac{1}{n}$ . Second, we have  $0 \leq \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}} \leq \frac{1}{1-\gamma} \mathbf{1}$ , which implies  $\left\| \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}} \right\|_{\text{span}} \leq \frac{1}{1-\gamma}$ . Therefore, assuming that  $n \geq 4$ , by Lemma 22 we have with probability at least  $1 - \delta$  that, letting  $\bar{V} = \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}} - \left( \min_s \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}}(s) \right) \mathbf{1}$  and  $\ell = \left\lceil \log_2 \log_2 \left( \left\| \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}} \right\|_{\text{span}} + 4 \right) \right\rceil$ , for all  $k = 0, \dots, \ell$ , we have

$$\left| \left( \hat{P}_\pi - P_\pi \right) (\bar{V})^{\circ 2^k} \right| \leq \sqrt{\frac{\beta \mathbb{V}_{\hat{P}_\pi} \left[ (\bar{V})^{\circ 2^k} \right]}{n}} + \frac{\beta \cdot 2^k}{n} (\|\bar{V}\|_\infty + 1)^{2^k} \mathbf{1} \quad (69)$$

for  $\beta = 16 \log \left( 12 \frac{SAn}{(1-\gamma)^{2\delta}} \right)$ .

Now we can immediately apply (69) with (Zurek and Chen, 2024, Lemma 16) to obtain that

$$\left\| V_{\gamma, \tilde{r}}^{\hat{\pi}} - \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}} \right\|_{\infty} \leq \frac{4(\ell+1)\beta}{(1-\gamma)n} \left\| \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}} \right\|_{\text{span}} + \frac{2(\ell+1)}{1-\gamma} \sqrt{\frac{2\beta(\left\| \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}} \right\|_{\text{span}} + 1)}{n}}. \quad (70)$$

Now all that remains is to simplify (70) which can be done in an identical manner as to (Zurek and Chen, 2024, Proof of Theorem 9) to obtain that

$$\left\| V_{\gamma, \tilde{r}}^{\hat{\pi}} - \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}} \right\|_{\infty} \leq \frac{24 \log_2 \log_2 \left( \frac{1}{1-\gamma} + 4 \right)}{1-\gamma} \sqrt{\frac{\left\| \hat{V}_{\gamma, \tilde{r}}^{\hat{\pi}} \right\|_{\text{span}} + 1}{n} 16 \log \left( \frac{12SAn}{(1-\gamma)^2 \delta} \right)}.$$

■

Now we can apply these error bounds to the analysis of Algorithm 4. We recall the definitions of some objects which appear in Algorithm 4 for convenience, which will be in effect for the remainder of this subsection. We define the empirical transition kernel  $\hat{P}(s' | s, a) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{S_{s,a}^i = s'\}$ , for all  $s' \in \mathcal{S}$ , using the  $n$  samples drawn from all state-action pairs within Algorithm 4. For all integers  $i \in \{2, \dots, \lceil \log_2 n \rceil\}$  we define  $M_i = 2^i$ , and we define  $\gamma_i$  so that  $\frac{1}{1-\gamma_i} = \max \left\{ \frac{\sqrt{nM_i}}{\alpha(\delta, n)2\sqrt{2}}, 1 \right\}$ . For each  $i \in \{2, \dots, \lceil \log_2 n \rceil\}$  we also define  $\tilde{\pi}_i$ ,  $\tilde{V}_i$ , and  $\tilde{r}_i$  as the outputs from the Span-Constrained Planning Algorithm 3 with input DMDP  $(\hat{P}, r, \gamma_i)$ , input span constraint bound  $M_i$ , and input target error  $\frac{1}{n}$ . For the remainder of this proof we fix a policy  $\pi$  such that  $\rho^\pi$  is a constant vector. ( $\pi$  will later be chosen to optimize a certain deterministic quantity.)

**Lemma 23** Define the function  $\alpha(\delta, \tilde{n}) = 24\sqrt{16 \log \left( \frac{24SAn\tilde{n}^5}{\delta} \right) \log_2(\log_2(\tilde{n} + 4))}$ . Fix  $\delta > 0$ . Then with probability at least  $1 - \delta$ , we have for all  $i = 2, \dots, \lceil \log_2 n \rceil$  that

$$\left\| V_{\gamma_i}^{\pi} - \hat{V}_{\gamma_i}^{\pi} \right\|_{\infty} \leq \frac{\alpha(\delta, n)}{1-\gamma_i} \sqrt{\frac{\left\| V_{\gamma_i}^{\pi} \right\|_{\text{span}} + 1}{n}} \quad (71)$$

and also the Span-Constrained Planning Algorithm 3 used within subroutine on line 8 outputs a policy  $\tilde{\pi}_i$ , approximate value function  $\tilde{V}_i$ , and reward function  $\tilde{r}_i$  such that

$$\left\| \hat{V}_{\gamma_i, \tilde{r}_i}^{\tilde{\pi}_i} - \hat{V}_{\gamma_i, M_i}^{\star} \right\|_{\infty} \leq \frac{1}{n} \quad (72)$$

$$\left\| \tilde{V}_i - \hat{V}_{\gamma_i, M_i}^{\star} \right\|_{\infty} \leq \frac{1}{n} \quad (73)$$

$$\left\| \hat{V}_{\gamma_i, \tilde{r}_i}^{\tilde{\pi}_i} - V_{\gamma_i, \tilde{r}_i}^{\tilde{\pi}_i} \right\|_{\infty} \leq \frac{\alpha(\delta, n)}{1-\gamma_i} \sqrt{\frac{\left\| \hat{V}_{\gamma_i, \tilde{r}_i}^{\tilde{\pi}_i} \right\|_{\text{span}} + 1}{n}}. \quad (74)$$

**Proof** First, we note that the properties (72) and (73) immediately follow from Lemma 3. By a union bound and Lemmas 19 and 21, we have that with probability at least  $1 - 2(\lceil \log_2 n \rceil - 1)\delta'$ , for all  $i = 2, \dots, \lceil \log_2 n \rceil$ , it holds that

$$\left\| V_{\gamma_i}^{\pi} - \hat{V}_{\gamma_i}^{\pi} \right\|_{\infty} \leq \frac{24 \log_2 \log_2 \left( \frac{1}{1-\gamma_i} + 4 \right)}{1-\gamma_i} \sqrt{\frac{\left\| V_{\gamma_i}^{\pi} \right\|_{\text{span}} + 1}{n} 16 \log \left( \frac{12SAn}{(1-\gamma_i)^2 \delta'} \right)}.$$



and

$$\left\| V_{\gamma_i, \tilde{r}_i}^{\tilde{\pi}_i} - \widehat{V}_{\gamma_i, \tilde{r}_i}^{\tilde{\pi}_i} \right\|_{\infty} \leq \frac{24 \log_2 \log_2 \left( \frac{1}{1-\gamma_i} + 4 \right)}{1-\gamma_i} \sqrt{\frac{\left\| \widehat{V}_{\gamma_i, \tilde{r}_i}^{\tilde{\pi}_i} \right\|_{\text{span}} + 1}{n} 16 \log \left( \frac{12 S A n}{(1-\gamma_i)^2 \delta'} \right)}$$

Since

$$\frac{1}{1-\gamma_i} \leq \frac{1}{1-\gamma_{\lceil \log_2 n \rceil}} \leq \max \left\{ 1, \frac{\sqrt{n M_{\lceil \log_2 n \rceil}}}{\alpha(\delta, n) 2\sqrt{2}} \right\} \leq \max \left\{ 1, \frac{\sqrt{n 2n}}{\alpha(\delta, n) 2\sqrt{2}} \right\} \leq n$$

since  $\alpha(\delta, n) \geq 1$  and  $n \geq 1$ , and also noting that  $\lceil \log_2 n \rceil - 1 \leq \log_2 n \leq n$ , by taking  $\delta' = \frac{\delta}{2(\lceil \log_2 n \rceil - 1)}$  we obtain that with probability at least  $1 - \delta$  both

$$\begin{aligned} \left\| V_{\gamma_i}^{\pi} - \widehat{V}_{\gamma_i}^{\pi} \right\|_{\infty} &\leq \frac{24 \log_2 \log_2 (n+4)}{1-\gamma_i} \sqrt{\frac{\left\| V_{\gamma_i}^{\pi} \right\|_{\text{span}} + 1}{n} 16 \log \left( \frac{24 S A n^3 (\lceil \log_2 n \rceil - 1)}{\delta} \right)} \\ &\leq \frac{\alpha(\delta, n)}{1-\gamma_i} \sqrt{\frac{\left\| V_{\gamma_i}^{\pi} \right\|_{\text{span}} + 1}{n}} \end{aligned}$$

and similarly

$$\left\| \widehat{V}_{\gamma_i, \tilde{r}_i}^{\tilde{\pi}_i} - V_{\gamma_i, \tilde{r}_i}^{\tilde{\pi}_i} \right\|_{\infty} \leq \frac{\alpha(\delta, n)}{1-\gamma_i} \sqrt{\frac{\left\| \widehat{V}_{\gamma_i, \tilde{r}_i}^{\tilde{\pi}_i} \right\|_{\text{span}} + 1}{n}}$$

for all  $i = 2, \dots, \lceil \log_2 n \rceil$ . ■

**Lemma 24** *Under the event described in Lemma 23, we have*

$$\rho^{\tilde{\pi}_i} \geq \widehat{L}(i) \mathbf{1} = (1-\gamma_i) \min_s \tilde{V}_i(s) \mathbf{1} - 2 \frac{1-\gamma_i}{n} \mathbf{1} - \alpha(\delta, n) \sqrt{\frac{M_i + \frac{2}{n} + 1}{n}} \mathbf{1} \quad (75)$$

for all  $i = 2, \dots, \lceil \log_2 n \rceil$ .

**Proof** The proof is very similar to the first part of the proof of Lemma 16. Fix  $i \in \{2, \dots, \lceil \log_2 n \rceil\}$ .

First, by using (72) and the fact that  $\left\| \widehat{V}_{\gamma_i, M_i}^{\star} \right\|_{\text{span}} \leq M_i$ , we have

$$\left\| \widehat{V}_{\gamma_i, \tilde{r}_i}^{\tilde{\pi}_i} \right\|_{\text{span}} \leq \left\| \widehat{V}_{\gamma_i, M_i}^{\star} \right\|_{\text{span}} + \left\| \widehat{V}_{\gamma_i, \tilde{r}_i}^{\tilde{\pi}_i} - \widehat{V}_{\gamma_i, M_i}^{\star} \right\|_{\text{span}} \leq \left\| \widehat{V}_{\gamma_i, M_i}^{\star} \right\|_{\text{span}} + 2 \left\| \widehat{V}_{\gamma_i, \tilde{r}_i}^{\tilde{\pi}_i} - \widehat{V}_{\gamma_i, M_i}^{\star} \right\|_{\infty} \leq M_i + \frac{2}{n}. \quad (76)$$

We note that by triangle inequality and (72) and (73) we have  $\|\tilde{V}_i - \hat{V}_{\gamma_i, \tilde{r}_i}^{\pi_i}\|_\infty \leq \frac{2}{n}$ . Using this bound on  $\|\tilde{V}_i - \hat{V}_{\gamma_i, \tilde{r}_i}^{\pi_i}\|_\infty$ , (74), and (76), we have

$$\begin{aligned}
\rho^{\pi_i} &\geq (1 - \gamma_i) \min_s V_{\gamma_i}^{\pi_i}(s) \mathbf{1} \\
&\geq (1 - \gamma_i) \min_s V_{\gamma_i, \tilde{r}_i}^{\pi_i}(s) \mathbf{1} \\
&\geq (1 - \gamma_i) \min_s \hat{V}_{\gamma_i, \tilde{r}_i}^{\pi_i}(s) \mathbf{1} - (1 - \gamma_i) \|V_{\gamma_i, \tilde{r}_i}^{\pi_i} - \hat{V}_{\gamma_i, \tilde{r}_i}^{\pi_i}\|_\infty \mathbf{1} \\
&\geq (1 - \gamma_i) \min_s \tilde{V}_i(s) \mathbf{1} - (1 - \gamma_i) \|\tilde{V}_i - \hat{V}_{\gamma_i, \tilde{r}_i}^{\pi_i}\|_\infty \mathbf{1} - (1 - \gamma_i) \|V_{\gamma_i, \tilde{r}_i}^{\pi_i} - \hat{V}_{\gamma_i, \tilde{r}_i}^{\pi_i}\|_\infty \mathbf{1} \\
&\geq (1 - \gamma_i) \min_s \tilde{V}_i(s) \mathbf{1} - 2 \frac{1 - \gamma_i}{n} \mathbf{1} - \alpha(\delta, n) \sqrt{\frac{\|\hat{V}_{\gamma_i, \tilde{r}_i}^{\pi_i}\|_{\text{span}} + 1}{n}} \mathbf{1} \\
&\geq (1 - \gamma_i) \min_s \tilde{V}_i(s) \mathbf{1} - 2 \frac{1 - \gamma_i}{n} \mathbf{1} - \alpha(\delta, n) \sqrt{\frac{M_i + \frac{2}{n} + 1}{n}} \mathbf{1} \\
&= \hat{L}(i) \mathbf{1}.
\end{aligned}$$

■

Before continuing we need to establish some relationships between DMDP and AMDP quantities for the policy  $\pi$ . This lemma is similar to (Wei et al., 2020, Lemma 2) but for a general policy.

**Lemma 25** *Suppose  $\rho^\pi$  is constant. Then*

1.  $\left\| \frac{1}{1-\gamma} \rho^\pi - V_\gamma^\pi \right\|_\infty \leq \|h^\pi\|_{\text{span}}$
2.  $\|V_\gamma^\pi\|_{\text{span}} \leq 2 \|h^\pi\|_{\text{span}}$ .

**Proof** For the first statement, we have that

$$V_\gamma^\pi = (I - \gamma P_\pi)^{-1} r_\pi = (I - \gamma P_\pi)^{-1} (\rho^\pi + h^\pi - P_\pi h^\pi) = (I - \gamma P_\pi)^{-1} \rho^\pi + (I - \gamma P_\pi)^{-1} (I - P_\pi) h^\pi.$$

Since  $P_\pi \rho^\pi = \rho^\pi$ , by the Neumann series the first term is equal to  $\frac{1}{1-\gamma} \rho^\pi$ . By a standard calculation (e.g. (Zurek and Chen, 2025, Lemma 20)) the second term satisfies  $\|(I - \gamma P_\pi)^{-1} (I - P_\pi) h^\pi\|_\infty \leq \|h^\pi\|_{\text{span}}$ . Therefore we have that  $\left\| \frac{1}{1-\gamma} \rho^\pi - V_\gamma^\pi \right\|_\infty \leq \|(I - \gamma P_\pi)^{-1} (I - P_\pi) h^\pi\|_\infty \leq \|h^\pi\|_{\text{span}}$ .

For the second statement, since  $\rho^\pi$  is constant, we have that

$$\|V_\gamma^\pi\|_{\text{span}} = \left\| V_\gamma^\pi - \frac{1}{1-\gamma} \rho^\pi \right\|_{\text{span}} \leq 2 \left\| V_\gamma^\pi - \frac{1}{1-\gamma} \rho^\pi \right\|_\infty \leq 2 \|h^\pi\|_{\text{span}}.$$

■

**Lemma 26** *Suppose that for some integer  $i \geq 2$  we have  $\|h^\pi\|_{\text{span}} + 1 \leq M_i/4$ , and also that  $n \geq 2\alpha(\delta, n)^2$ . Then under the event in Lemma 23, we have*

$$\hat{L}(i) \mathbf{1} \geq \rho^\pi - \alpha(\delta, n)(2 + \sqrt{2}) \sqrt{\frac{M_i}{n}} \mathbf{1}.$$

**Proof** The fact that  $n \geq 2\alpha(\delta, n)^2$  implies that we always have

$$\frac{1}{1-\gamma_i} = \max \left\{ \frac{\sqrt{nM_i}}{\alpha(\delta, n)2\sqrt{2}}, 1 \right\} = \frac{\sqrt{nM_i}}{\alpha(\delta, n)2\sqrt{2}}$$

for all  $i \geq 2$ , since  $M_i = 2^i \geq 4$ . Now, by triangle inequality,  $\|\cdot\|_{\text{span}} \leq 2\|\cdot\|_{\infty}$ , (71), Lemma 25, and the above expression for  $\frac{1}{1-\gamma_i}$ , we calculate that

$$\begin{aligned} \|\widehat{V}_{\gamma_i}^{\pi}\|_{\text{span}} &\leq \|V_{\gamma_i}^{\pi}\|_{\text{span}} + \|\widehat{V}_{\gamma_i}^{\pi} - V_{\gamma_i}^{\pi}\|_{\text{span}} \\ &\leq \|V_{\gamma_i}^{\pi}\|_{\text{span}} + 2\|\widehat{V}_{\gamma_i}^{\pi} - V_{\gamma_i}^{\pi}\|_{\infty} \\ &\leq \|V_{\gamma_i}^{\pi}\|_{\text{span}} + 2\frac{\alpha(\delta, n)}{1-\gamma_i} \sqrt{\frac{\|V_{\gamma_i}^{\pi}\|_{\text{span}} + 1}{n}} \\ &\leq 2\|h^{\pi}\|_{\text{span}} + 2\frac{\alpha(\delta, n)}{1-\gamma_i} \sqrt{\frac{2\|h^{\pi}\|_{\text{span}} + 1}{n}} \\ &= 2\|h^{\pi}\|_{\text{span}} + \frac{\sqrt{nM_i}}{\sqrt{2}} \sqrt{\frac{2\|h^{\pi}\|_{\text{span}} + 1}{n}} \\ &\leq \frac{M_i}{2} + \frac{\sqrt{M_i}}{\sqrt{2}} \sqrt{M_i/2} \\ &= M_i. \end{aligned}$$

Consequently by Lemma 3 we have that  $\widehat{V}_{\gamma_i, M_i}^{\star} \geq \widehat{V}_{\gamma_i}^{\pi}$ , and thus

$$\widetilde{V}_i \geq \widehat{V}_{\gamma_i, M_i}^{\star} - \frac{1}{n} \mathbf{1} \geq \widehat{V}_{\gamma_i}^{\pi} - \frac{1}{n} \mathbf{1}.$$

Now we lower-bound  $\widehat{V}_{\gamma_i}^{\pi}$ . We have that

$$\begin{aligned} \widehat{V}_{\gamma_i}^{\pi} &\geq V_{\gamma_i}^{\pi} - \|\widehat{V}_{\gamma_i}^{\pi} - V_{\gamma_i}^{\pi}\|_{\infty} \mathbf{1} \\ &\geq \frac{1}{1-\gamma_i} \rho^{\pi} - \left\| \frac{1}{1-\gamma_i} \rho^{\pi} - V_{\gamma_i}^{\pi} \right\|_{\infty} \mathbf{1} - \|\widehat{V}_{\gamma_i}^{\pi} - V_{\gamma_i}^{\pi}\|_{\infty} \mathbf{1} \\ &\geq \frac{1}{1-\gamma_i} \rho^{\pi} - \|h^{\pi}\|_{\text{span}} \mathbf{1} - \frac{\alpha(\delta, n)}{1-\gamma_i} \sqrt{\frac{\|V_{\gamma_i}^{\pi}\|_{\text{span}} + 1}{n}} \mathbf{1} \end{aligned} \tag{77}$$

using Lemma 25 and (71) again. Thus

$$\begin{aligned} \widetilde{V}_i &\geq \widehat{V}_{\gamma_i}^{\pi} - \frac{1}{n} \mathbf{1} \\ &\geq \frac{1}{1-\gamma_i} \rho^{\pi} - \|h^{\pi}\|_{\text{span}} \mathbf{1} - \frac{\alpha(\delta, n)}{1-\gamma_i} \sqrt{\frac{\|V_{\gamma_i}^{\pi}\|_{\text{span}} + 1}{n}} \mathbf{1} - \frac{1}{n} \mathbf{1} \\ &\geq \frac{1}{1-\gamma_i} \rho^{\pi} - \frac{M_i}{4} \mathbf{1} - \frac{\alpha(\delta, n)}{1-\gamma_i} \sqrt{\frac{2\|h^{\pi}\|_{\text{span}} + 1}{n}} \\ &\geq \frac{1}{1-\gamma_i} \rho^{\pi} - \frac{M_i}{4} \mathbf{1} - \frac{\alpha(\delta, n)}{1-\gamma_i} \sqrt{\frac{M_i/2}{n}} \end{aligned}$$

where in the first inequality we combine (73) with Lemma 3, which states that  $\widehat{V}_{\gamma_i, M_i}^* \geq \widehat{V}_{\gamma_i}^\pi$  since (as we calculated above) we have  $\|\widehat{V}_{\gamma_i}^\pi\|_{\text{span}} \leq M_i$ . In the second inequality we use (77), in the third we use Lemma 25, and in the final inequality we use the assumption  $\|h^\pi\|_{\text{span}} + 1 \leq M_i/4$ . Therefore, we have

$$\begin{aligned}
\widehat{L}(i) &= (1 - \gamma_i) \min_s \widetilde{V}_i(s) - 2 \frac{1 - \gamma_i}{n} - \alpha(\delta, n) \sqrt{\frac{M_i + \frac{2}{n} + 1}{n}} \\
&\geq \rho^\pi(s_0) - (1 - \gamma_i) \frac{M_i}{4} \mathbf{1} - \alpha(\delta, n) \sqrt{\frac{M_i/2}{n}} - 2 \frac{1 - \gamma_i}{n} - \alpha(\delta, n) \sqrt{\frac{M_i + \frac{2}{n} + 1}{n}} \\
&= \rho^\pi(s_0) - \frac{\alpha(\delta, n) 2\sqrt{2}}{\sqrt{nM_i}} \frac{M_i}{4} \mathbf{1} - \alpha(\delta, n) \sqrt{\frac{M_i/2}{n}} - 2 \frac{1 - \gamma_i}{n} - \alpha(\delta, n) \sqrt{\frac{M_i + \frac{2}{n} + 1}{n}} \\
&\geq \rho^\pi(s_0) - \alpha(\delta, n) \sqrt{2} \sqrt{\frac{M_i}{n}} - \frac{1}{\sqrt{2}} \frac{1}{n} \frac{1}{2\sqrt{n}} - \alpha(\delta, n) \sqrt{\frac{M_i + \frac{2}{n} + 1}{n}} \\
&\geq \rho^\pi(s_0) - \alpha(\delta, n) (2 + \sqrt{2}) \sqrt{\frac{M_i}{n}},
\end{aligned}$$

where  $s_0$  is an arbitrary state (since  $\rho^\pi$  is constant), and we used that  $\alpha(\delta, n) \geq 1$  and  $M_i \geq 4$ . ■

**Proof of Theorem 4** We assume that the event in Lemma 23 holds, which occurs with probability at least  $1 - \delta$ . Lemma 24 immediately implies that

$$\rho^{\widehat{\pi}} = \rho^{\widehat{\pi}_i} \geq \widehat{L}(i) \mathbf{1}.$$

We also trivially have that  $\rho^{\widehat{\pi}} \geq \mathbf{0}$ , which implies that

$$\rho^{\widehat{\pi}} \geq \max\{\widehat{L}(i), 0\} \mathbf{1} = \widehat{\rho}.$$

Now we lower-bound  $\widehat{\rho}$ . First note that if  $n < 2\alpha(\delta, n)^2$ , then

$$10\alpha(\delta, n) \sqrt{\frac{\|h^\pi\|_{\text{span}} + 1}{n}} \geq 10/\sqrt{2} > 1,$$

so the desired conclusion holds trivially since we must have  $\widehat{\rho} \geq \mathbf{0}$ , and  $\rho^\pi \leq \mathbf{1}$ . Thus we can now consider the situation  $n \geq 2\alpha(\delta, n)^2$ . We note that the smallest  $i = 2$  causes  $M_2/8 = \frac{1}{2} < 1$ . Therefore by the construction of the  $M_i$ 's, either there exists  $i$  such that

$$M_i/8 \leq \|h^\pi\|_{\text{span}} + 1 \leq M_i/4 \tag{78}$$

or we have that  $\|h^\pi\|_{\text{span}} + 1 > M_i/4$  for all  $i$ . Since the largest  $i = \lceil \log_2 n \rceil$  causes  $M_i/4 > n/4$ , in this second case we thus have that

$$10\alpha(\delta, n) \sqrt{\frac{\|h^\pi\|_{\text{span}} + 1}{n}} \geq 10\alpha(\delta, n) \frac{1}{2} > 1,$$

so again the desired conclusion holds trivially since we must have  $\hat{\rho} \geq \mathbf{0}$ . In the first case that (78) holds for some  $i$ , since also  $n \geq 2\alpha(\delta, n)^2$ , by Lemma 26 we have that

$$\begin{aligned} \hat{\rho} &\geq \hat{L}(\hat{i})\mathbf{1} \geq \hat{L}(i)\mathbf{1} \geq \rho^\pi - \alpha(\delta, n)(2 + \sqrt{2})\sqrt{\frac{M_i}{n}}\mathbf{1} \\ &\geq \rho^\pi - (2 + \sqrt{2})\sqrt{8}\alpha(\delta, n)\sqrt{\frac{\|h^\pi\|_{\text{span}} + 1}{n}} \\ &\geq \rho^\pi - 10\alpha(\delta, n)\sqrt{\frac{\|h^\pi\|_{\text{span}} + 1}{n}}. \end{aligned}$$

Since we have proven this for arbitrary fixed  $\pi$ , we can apply this result to a policy  $\pi \in \sup_{\pi: \rho^\pi \text{ constant}} \rho^\pi(s_0) - 10\alpha(\delta, n)\sqrt{\frac{\|h^\pi\|_{\text{span}} + 1}{n}}$  (for an arbitrary state  $s_0$ ) to obtain the desired conclusion. Finally, we can set  $C_4 = 10$ .  $\blacksquare$

## Appendix D. Examples

In this section we provide the two examples mentioned in Subsection 3.3 of situations where the guarantee of Theorem 4 could be much better than the minimax rate. In both examples each solid line represents a single action, and an expression “ $R = \dots$ ” denotes the reward for said action. If the line splits into multiple dashed arrows then this indicates that the next-state transition from following this action is stochastic, and the numbers next to each dashed line indicate the transition probabilities. Otherwise if the line does not split then it indicates a deterministic next-state transition.

First we provide the simpler of the two examples, where  $\|h^\star\|_{\text{span}}$  is arbitrarily larger than  $\inf_{\pi: \rho^\pi = \rho^\star} \|h^\pi\|_{\text{span}}$ .

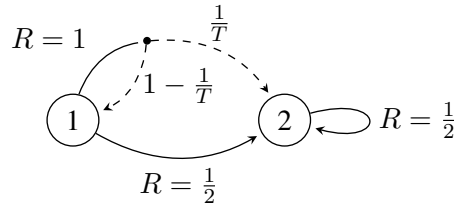


Figure 1: An MDP where  $\|h^\star\|_{\text{span}}$  can be arbitrarily larger than  $\inf_{\pi: \rho^\pi = \rho^\star} \|h^\pi\|_{\text{span}}$ .

**Theorem 27** *Consider the MDP displayed in Figure 1. For any  $T \geq 1$ , we have  $\|h^\star\|_{\text{span}} = \frac{T}{2}$  and  $\inf_{\pi: \rho^\pi = \rho^\star} \|h^\pi\|_{\text{span}} = 0$ .*

Next, we provide an instance where  $\|h^\star\|_{\text{span}}$  can be arbitrarily large but a policy  $\pi$  with an arbitrarily low level of suboptimality  $\rho^\star - \rho^\pi$  can have  $\|h^\pi\|_{\text{span}} = O(1)$ .

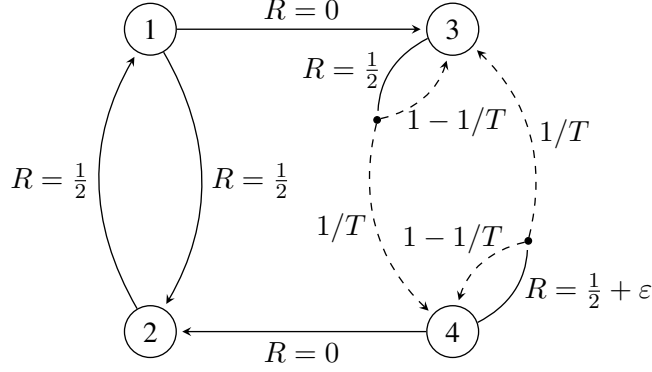


Figure 2: An MDP where  $\|h^*\|_{\text{span}}$  can be arbitrarily larger than  $\|h^\pi\|_{\text{span}}$  for some near-optimal policy  $\pi$  satisfying  $\rho^\pi = \rho^* - \frac{\epsilon}{2}\mathbf{1}$ .

**Theorem 28** Consider the MDP displayed in Figure 2. For any  $T \geq 1$  and  $\epsilon > 0$ , we have that  $\|h^*\|_{\text{span}} = \frac{\epsilon T}{2} + \epsilon + \frac{1}{2}$ , but there exists some policy  $\pi$  with constant gain such that  $\rho^\pi = \rho^* - \frac{\epsilon}{2}\mathbf{1}$  and  $\|h^\pi\|_{\text{span}} = \frac{1}{2}$ .

#### D.1. Proofs of Theorems 27 and 28

**Proof of Theorem 27** It is easy to see that state 1 is transient under all policies and state 2 is absorbing, so all policies are gain-optimal and have gain  $\frac{1}{2}\mathbf{1}$ . Only state 1 has multiple possible actions, so it suffices to consider the two policies  $\pi_{\text{up}}$ , which takes the “up” action which has nonzero probability of returning to state 1, and  $\pi_{\text{down}}$ , which leads to an immediate transition to state 2. It is trivial to see that  $h^{\pi_{\text{down}}} = \mathbf{0}$ , so we have that  $\|h^{\pi_{\text{down}}}\|_{\text{span}} = 0$ . To compute  $h^{\pi_{\text{up}}}$ , we must have  $h^{\pi_{\text{up}}}(2) = 0$ , so we can then calculate that

$$\begin{aligned} h^{\pi_{\text{up}}}(1) + \frac{1}{2} &= 1 + \left(1 - \frac{1}{T}\right)h^{\pi_{\text{up}}}(1) + \frac{1}{T}h^{\pi_{\text{up}}}(2) \\ \implies h^{\pi_{\text{up}}}(1) &= \frac{T}{2} + h^{\pi_{\text{up}}}(2) = \frac{T}{2} \end{aligned}$$

and thus  $\|h^{\pi_{\text{up}}}\|_{\text{span}} = T/2$ . Since these are the only two stationary deterministic policies, one of them must be Blackwell-optimal, and since they have equal gain and elementwise  $h^{\pi_{\text{up}}} \geq h^{\pi_{\text{down}}}$  (with a strict inequality in state 1), we have that  $\pi_{\text{up}} = \pi^*$  and  $\|h^*\|_{\text{span}} = \|h^{\pi_{\text{up}}}\|_{\text{span}} = T/2$ . Since  $\pi_{\text{down}}$  is gain-optimal it is immediate from  $h^{\pi_{\text{down}}} = \mathbf{0}$  that  $\inf_{\pi: \rho^\pi = \rho^*} \|h^\pi\|_{\text{span}} = 0$ . ■

**Proof of Theorem 28** There are two states, 1 and 4, where multiple actions are possible. We name the actions in state 1 the “down” action (which leads to 2) and the “right” action (which leads to 3), and we name the actions in state 4 the “up” action (which has positive probability of leading to 3) and the “left” action (which leads to 2). A deterministic stationary policy can be specified by its two choices between the actions available at states 1 and 4. We thus use  $\pi_{DL}$  to indicate the policy which takes the down action in state 1 and the left action in state 4, and likewise for other choices of  $\{D, R\} \times \{U, L\}$ .

It is easy to see that since  $\varepsilon > 0$  the unique gain-optimal policy is  $\pi_{RU}$  which has  $\rho^{\pi_{RU}} = \rho^* = \frac{1+\varepsilon}{2}\mathbf{1}$ . Thus this policy must also be Blackwell-optimal. We now compute  $h^* = h^{\pi_{RU}}$ . We have

$$\begin{aligned} h^*(3) + \frac{1+\varepsilon}{2} &= \frac{1}{2} + (1 - \frac{1}{T})h^*(3) + \frac{1}{T}h^*(4) \\ \implies h^*(3) &= h^*(4) - \frac{\varepsilon T}{2} \end{aligned}$$

and since the stationary distribution of  $\pi_{RU}$  has equal mass on states 3 and 4 we must have  $h^*(3) + h^*(4) = 0$ , which implies that  $h^*(3) = -\frac{\varepsilon T}{4}$  and  $h^*(4) = \frac{\varepsilon T}{4}$ . It is then easy to check that  $h^*(1) = h^*(3) - \frac{1+\varepsilon}{2} = \frac{-\varepsilon T - 2\varepsilon - 2}{4}$  and  $h^*(2) = h^*(1) - \frac{\varepsilon}{2} = \frac{-\varepsilon T - 4\varepsilon - 2}{4}$ , so we have that  $\|h^*\|_{\text{span}} = \frac{\varepsilon T}{4} - \frac{-\varepsilon T - 4\varepsilon - 2}{4} = \frac{\varepsilon T}{2} + \varepsilon + \frac{1}{2}$ .

Next we consider the policy  $\pi_{DL}$ . It is easy to see that  $\rho^{\pi_{DL}} = \frac{1}{2}\mathbf{1}$ , which is constant and which satisfies  $\rho^* - \rho^{\pi_{DL}} = \frac{\varepsilon}{2}\mathbf{1}$ . Now we compute  $\|h^{\pi_{DL}}\|_{\text{span}}$ . It is immediate to see that  $h^{\pi_{DL}}(1) = h^{\pi_{DL}}(2) = 0$ . Then we can calculate that  $h^{\pi_{DL}}(4) = h^{\pi_{DL}}(2) - \frac{1}{2} = -\frac{1}{2}$ , and finally that

$$\begin{aligned} h^{\pi_{DL}}(3) + \frac{1}{2} &= \frac{1}{2} + (1 - \frac{1}{T})h^{\pi_{DL}}(3) + \frac{1}{T}h^{\pi_{DL}}(4) \\ \implies h^{\pi_{DL}}(3) &= h^{\pi_{DL}}(4) = -\frac{1}{2}. \end{aligned}$$

Therefore we have that  $\|h^{\pi_{DL}}\|_{\text{span}} = 0 - -\frac{1}{2} = \frac{1}{2}$ . ■