

Breaking the Sample Size Barrier in Model-Based Reinforcement Learning with a Generative Model

Gen Li*
Tsinghua

Yuting Wei†
CMU

Yuejie Chi‡
CMU

Yuantao Gu*
Tsinghua

Yuxin Chen§
Princeton

May 2020; Revised: September 2020

Abstract

We investigate the sample efficiency of reinforcement learning in a γ -discounted infinite-horizon Markov decision process (MDP) with state space \mathcal{S} and action space \mathcal{A} , assuming access to a generative model. Despite a number of prior work tackling this problem, a complete picture of the trade-offs between sample complexity and statistical accuracy is yet to be determined. In particular, prior results suffer from a sample size barrier, in the sense that their claimed statistical guarantees hold only when the sample size exceeds at least $\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^2}$ (up to some log factor). The current paper overcomes this barrier by certifying the minimax optimality of model-based reinforcement learning as soon as the sample size exceeds the order of $\frac{|\mathcal{S}||\mathcal{A}|}{1-\gamma}$ (modulo some log factor). More specifically, a *perturbed* model-based planning algorithm provably finds an ε -optimal policy with an order of $\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2} \log \frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)\varepsilon}$ samples for any $\varepsilon \in (0, \frac{1}{1-\gamma}]$. Along the way, we derive improved (instance-dependent) guarantees for model-based policy evaluation. To the best of our knowledge, this work provides the first minimax-optimal guarantee in a generative model that accommodates the entire range of sample sizes (beyond which finding a meaningful policy is information theoretically impossible).

Keywords: model-based reinforcement learning, minimaxity, planning, policy evaluation, instance-dependent guarantees, generative model

Contents

| | | |
|----------|---|----------|
| 1 | Introduction | 2 |
| 2 | Problem formulation | 4 |
| 2.1 | Models and background | 4 |
| 2.2 | Notation | 5 |
| 3 | Main results | 5 |
| 3.1 | Theoretical guarantees for model-based reinforcement learning | 5 |
| 3.2 | Comparisons with prior work and implications | 7 |
| 4 | Other related work | 8 |
| 5 | Analysis | 8 |
| 5.1 | Matrix notation and Bellman equations | 8 |
| 5.2 | Analysis: model-based policy evaluation | 9 |
| 5.3 | Analysis: model-based planning | 10 |

*Department of Electronic Engineering, Tsinghua University, Beijing 100084, China.

†Department of Statistics and Data Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

‡Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

§Department of Electrical Engineering, Princeton University, Princeton, NJ 08544, USA.

| | | |
|----------|--|-----------|
| 5.3.1 | Value function estimation for a policy obeying Bernstein-type conditions | 11 |
| 5.3.2 | Decoupling statistical dependency via (s, a) -absorbing MDPs | 12 |
| 5.3.3 | A tie-breaking argument | 14 |
| 5.3.4 | Proof of Theorem 1 | 14 |
| 6 | Discussion | 15 |
| A | Preliminary facts | 16 |
| B | Proofs of auxiliary lemmas | 17 |
| B.1 | Proofs of Lemma 1 and Lemma 2 | 17 |
| B.1.1 | Proof of Lemma 8 | 20 |
| B.2 | Proof of Lemma 3 | 21 |
| B.3 | Proof of Lemma 4 | 22 |
| B.4 | Proof of Lemma 5 | 22 |
| B.5 | Proof of Lemma 6 | 23 |

1 Introduction

Reinforcement learning (RL) (Sutton and Barto, 2018; Szepesvári, 2010), which is frequently modeled as learning and decision making in a Markov decision process (MDP), is garnering growing interest in recent years due to its remarkable success in practice. A core objective of RL is to search for a policy — based on a collection of noisy data samples — that approximately maximizes expected rewards in an MDP, without direct access to a precise description of the underlying model.¹ In contemporary applications, it is increasingly more common to encounter environments with prohibitively large state and action space, thus exacerbating the challenge of collecting enough samples to learn the model. To enable faithful policy learning in the sample-starved regime (i.e. the regime where the model complexity overwhelms the sample size), it is crucial to obtain a quantitative picture of the trade-off between sample complexity and statistical accuracy, and to design efficient algorithms that provably achieve the optimal trade-off.

Broadly speaking, there are two common algorithmic approaches: a model-based approach and a model-free one. In the model-based approach, one first learns to describe the unknown model using the data samples in hand, and then leverages the fitted model to perform planning — a task that can be accomplished by resorting to Bellman’s principle of optimality (Bellman, 1952; Puterman, 2014). An advantage of model-based algorithms is their flexibility: the learned model can be adapted to perform new ad-hoc tasks without revisiting the data samples. In comparison, the model-free approach attempts to compute the optimal policy (and the optimal value function) without learning the model explicitly, which lends itself to scenarios when a realistic model is difficult to construct or changes on the fly. Characterizing the sample complexities for both approaches has been the focal point of a large body of recent work, e.g. Agarwal et al. (2019); Azar et al. (2013); Jin et al. (2018); Kearns and Singh (1999); Sidford et al. (2018a,b); Tu and Recht (2018); Wainwright (2019a,b); Wang (2019).

In this paper, we pursue a comprehensive understanding of model-based RL, assuming access to a generative model — that is, a simulator that produces samples based on the transition kernel of the MDP for each state-action pair (Kakade, 2003; Kearns and Singh, 1999). To allow for more precise discussions, we focus our attention on an infinite-horizon discounted MDP with state space \mathcal{S} , action space \mathcal{A} and discount factor $0 < \gamma < 1$. We obtain N samples per state-action pair by querying the generative model. For an *arbitrary* target accuracy level $\varepsilon > 0$, a desired model-based planning algorithm should return an ε -optimal policy with a minimal number of calls to the generative model. Particular emphasis is placed on the sub-linear sampling scenario, in which the total sample size is smaller than the total number $|\mathcal{S}|^2|\mathcal{A}|$ of model parameters (so that it might be infeasible to estimate the model accurately).

Motivation: a sample size barrier. Several prior work was dedicated to investigating model-based RL with a generative model, which uncovered the minimax optimality of this approach for an already wide

¹Here and throughout, the “model” refers to the transition kernel and the rewards of the MDP taken collectively.

| Algorithm | Sample size range | Sample complexity | ε -range |
|--|---|--|---------------------------------------|
| Phased Q-learning Kearns and Singh (1999) | $[\frac{ S \mathcal{A} }{(1-\gamma)^5}, \infty)$ | $\frac{ S \mathcal{A} }{(1-\gamma)^7\varepsilon^2}$ | $(0, \frac{1}{1-\gamma}]$ |
| Empirical QVI Azar et al. (2013) | $[\frac{ S ^2 \mathcal{A} }{(1-\gamma)^2}, \infty)$ | $\frac{ S \mathcal{A} }{(1-\gamma)^3\varepsilon^2}$ | $(0, \frac{1}{\sqrt{(1-\gamma) S }}]$ |
| Sublinear randomized value iteration Sidford et al. (2018b) | $[\frac{ S \mathcal{A} }{(1-\gamma)^2}, \infty)$ | $\frac{ S \mathcal{A} }{(1-\gamma)^4\varepsilon^2}$ | $(0, \frac{1}{1-\gamma}]$ |
| Variance-reduced QVI Sidford et al. (2018a) | $[\frac{ S \mathcal{A} }{(1-\gamma)^3}, \infty)$ | $\frac{ S \mathcal{A} }{(1-\gamma)^3\varepsilon^2}$ | $(0, 1]$ |
| Randomized primal-dual method Wang (2019) | $[\frac{ S \mathcal{A} }{(1-\gamma)^2}, \infty)$ | $\frac{ S \mathcal{A} }{(1-\gamma)^4\varepsilon^2}$ | $(0, \frac{1}{1-\gamma}]$ |
| Empirical MDP + planning Agarwal et al. (2019) | $[\frac{ S \mathcal{A} }{(1-\gamma)^2}, \infty)$ | $\frac{ S \mathcal{A} }{(1-\gamma)^3\varepsilon^2}$ | $(0, \frac{1}{\sqrt{1-\gamma}}]$ |
| <i>Perturbed</i> empirical MDP + planning This paper | $[\frac{ S \mathcal{A} }{1-\gamma}, \infty)$ | $\frac{ S \mathcal{A} }{(1-\gamma)^3\varepsilon^2}$ | $(0, \frac{1}{1-\gamma}]$ |

Table 1: Comparisons with prior results (up to log factors) regarding finding an ε -optimal policy with a generative model. The sample size range and the ε -range stand for the range of sample size and optimality gap (e.g. ε -accuracy) for the claimed sample complexity to hold. Note that the results in ([Kearns and Singh, 1999](#); [Wang, 2019](#)) only hold for a restricted family of MDPs satisfying certain ergodicity assumptions. In addition, [Azar et al. \(2013\)](#) (resp. [Wainwright \(2019b\)](#)) showed that empirical QVI (resp. variance-reduced Q-learning) finds an ε -optimal Q-function estimate with sample complexity $\frac{|S||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}$ ($\varepsilon \in (0, 1]$) in a sample size range $[\frac{|S||\mathcal{A}|}{(1-\gamma)^3}, \infty)$, which did not translate directly to an ε -optimal policy.

regime ([Agarwal et al., 2019](#); [Azar et al., 2013](#)). However, the results therein often suffered from a sample complexity barrier that prevents us from obtaining a complete trade-off curve between sample complexity and statistical accuracy. For instance, the state-of-the-art result [Agarwal et al. \(2019\)](#) required the total sample size to at least exceed $\frac{|S||\mathcal{A}|}{(1-\gamma)^2}$ (up to some log factor), thus restricting the validity of the theory for broader contexts. In truth, this is not merely an issue for model-based planning; the same barrier already showed up when analyzing the simpler task of model-based policy evaluation ([Agarwal et al., 2019](#); [Pananjady and Wainwright, 2019](#)). Furthermore, a similar or even higher barrier emerged in prior theory for model-free methods; for instance, [Sidford et al. \(2018b\)](#); [Wainwright \(2019b\)](#) required the sample size to exceed $\frac{|S||\mathcal{A}|}{(1-\gamma)^3}$ modulo some log factor. In stark contrast, however, no lower bounds developed thus far preclude us from attaining reasonable statistical accuracy when going below the aforementioned sample complexity barriers.

Our contributions. Following the model-based approach, we propose to perform planning based on an empirical MDP learned from samples with mild *reward perturbation*. The perturbed model-based algorithm we propose provably finds an ε -optimal policy with an order of $\frac{|S||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}$ samples (up to log factor), which matches the minimax lower bound ([Azar et al., 2013](#)). Our result accommodates the full range of accuracy level ε (namely, $\varepsilon \in (0, \frac{1}{1-\gamma}]$), thus unveiling the minimaxity of our algorithm as soon as the sample size exceeds $\frac{|S||\mathcal{A}|}{1-\gamma}$ (modulo some log factor); encouragingly, this covers the *full* range of sample sizes that enable one to find a policy strictly better than a random guess. See Table 1 for detailed comparisons with prior literature. Along the way, we also derive (instance-dependent) statistical guarantees for policy evaluation, which strengthens state-of-the-art results by broadening the sample size range.

Our theory is established upon a novel combination of several key ideas: (1) a high-order expansion of the estimation error for value functions, coupled with fine-grained analysis for each term in the expansion; (2) the construction of auxiliary leave-one-out type (state-action-absorbing) MDPs — motivated by [Agarwal et al. \(2019\)](#) — that help decouple the complicated statistical dependency between the empirically optimal policy (as opposed to value functions) and data samples; (3) a tie-breaking argument guaranteeing that the empirically optimal policy stands out from all other policies under reward perturbation.

2 Problem formulation

2.1 Models and background

Basics of Markov decision processes. Consider an infinite-horizon discounted MDP represented by a quintuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, \gamma)$, where $\mathcal{S} := \{1, 2, \dots, |\mathcal{S}|\}$ denotes a finite set of states, $\mathcal{A} := \{1, 2, \dots, |\mathcal{A}|\}$ is a finite set of actions, $\gamma \in (0, 1)$ stands for the discount factor, and $r : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ represents the reward function, namely, $r(s, a)$ is the immediate reward received upon executing action a while in state s (here and throughout, we consider the normalized setting where the rewards lie within $[0, 1]$). In addition, $P : \mathcal{S} \times \mathcal{A} \mapsto \Delta(\mathcal{S})$ represents the probability transition kernel of the MDP, where $P(s'|s, a)$ denotes the probability of transiting from state s to state s' when action a is executed, and $\Delta(\mathcal{S})$ denotes the probability simplex over \mathcal{S} .

A deterministic policy is a mapping $\pi : \mathcal{S} \mapsto \mathcal{A}$ that maps a state to an action. The value function $V^\pi : \mathcal{S} \mapsto \mathbb{R}$ of a policy π is defined by

$$\forall s \in \mathcal{S} : \quad V^\pi(s) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s^t, a^t) \mid s^0 = s \right], \quad (1)$$

which is the expected discounted total reward starting from the initial state $s^0 = s$, with the actions taken according to the policy π (namely, $a^t = \pi(s^t)$ for all $t \geq 0$) and the trajectory generated based on the transition kernel (namely, $s^{t+1} \sim P(\cdot | s^t, a^t)$). It is easily seen that $0 \leq V^\pi(s) \leq \frac{1}{1-\gamma}$. The corresponding action-value function (or Q-function) $Q^\pi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ of a policy π is defined by

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad Q^\pi(s, a) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s^t, a^t) \mid s^0 = s, a^0 = a \right], \quad (2)$$

where the actions are taken according to the policy π after the initial action (i.e. $a^t = \pi(s^t)$ for all $t \geq 1$). It is well-known that there exists an optimal policy, denoted by π^* , that simultaneously maximizes $V^\pi(s)$ (resp. $Q^\pi(s, a)$) for all states $s \in \mathcal{S}$ (resp. state-action pairs $(s, a) \in (\mathcal{S} \times \mathcal{A})$) (Sutton and Barto, 2018). The corresponding value function $V^* := V^{\pi^*}$ (resp. action-value function $Q^* := Q^{\pi^*}$) is called the optimal value function (resp. optimal action-value function).

A generative model and an empirical MDP. The current paper focuses on a stylized generative model (also called a simulator) as studied in Kakade (2003); Kearns et al. (2002). Assuming access to this generative model, we collect N independent samples $s_{s,a}^i \sim P(\cdot | s, a)$, $i = 1, \dots, N$ for any state-action pair (s, a) , which allows us to construct an empirical transition kernel \hat{P} given by

$$\forall s' \in \mathcal{S}, \quad \hat{P}(s' | s, a) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{s_{s,a}^i = s'\}, \quad (3)$$

where $\mathbb{1}\{\cdot\}$ is the indicator function. In words, $\hat{P}(s' | s, a)$ counts the empirical frequency of transitions from (s, a) to state s' . The total sample size should be understood as $N^{\text{total}} := N|\mathcal{S}||\mathcal{A}|$. This leads to an empirical MDP $\widehat{\mathcal{M}} = (\mathcal{S}, \mathcal{A}, \hat{P}, r, \gamma)$ constructed from the data samples. We can define the value function and the action-value function of a policy π for $\widehat{\mathcal{M}}$ analogously, which we shall denote by \widehat{V}^π and \widehat{Q}^π , respectively. The optimal policy of $\widehat{\mathcal{M}}$ is denoted by $\widehat{\pi}^*$, with the optimal value function and Q-function denoted by $\widehat{V}^* := \widehat{V}^{\widehat{\pi}^*}$ and $\widehat{Q}^* := \widehat{Q}^{\widehat{\pi}^*}$, respectively.

Policy evaluation and planning. Given a few data samples in hand, the task of policy evaluation aims to compute or approximate the value function V^π under a given policy π . To be precise, for any target level $\varepsilon > 0$, the goal is to find an ε -accurate estimate V_{est}^π such that

$$\text{(policy evaluation)} \quad \forall s \in \mathcal{S} : \quad |V_{\text{est}}^\pi(s) - V^\pi(s)| \leq \varepsilon. \quad (4)$$

In contrast, the task of planning seeks to identify a policy that (approximately) maximizes the expected discounted reward given the data samples. Specifically, for any target level $\varepsilon > 0$, the aim is to compute an ε -accurate policy π_{est} obeying

$$\text{(planning)} \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad V^{\pi_{\text{est}}}(s) \geq V^*(s) - \varepsilon, \quad Q^{\pi_{\text{est}}}(s, a) \geq Q^*(s, a) - \varepsilon. \quad (5)$$

Naturally, one would hope to accomplish these tasks with as few samples as possible. Recall that for the normalized reward setting with $0 \leq r \leq 1$, the value function and Q-function fall within the range $[0, \frac{1}{1-\gamma}]$; this means that the range of the target level ε should be set to $\varepsilon \in [0, \frac{1}{1-\gamma}]$. The model-based approach typically starts by constructing an empirical MDP $\widehat{\mathcal{M}}$ based on all collected samples, and then “plugs in” this empirical model directly into the Bellman recursion to perform policy evaluation or planning, with prominent examples including Q-value iteration (QVI) and policy iteration (PI) (Bertsekas, 2017).

2.2 Notation

Let $\mathcal{X} := (|\mathcal{S}|, |\mathcal{A}|, \frac{1}{1-\gamma}, \frac{1}{\varepsilon})$. The notation $f(\mathcal{X}) = O(g(\mathcal{X}))$ means there exists a universal constant $C_1 > 0$ such that $f \leq C_1 g$, whereas the notation $f(\mathcal{X}) = \Omega(g(\mathcal{X}))$ means $g(\mathcal{X}) = O(f(\mathcal{X}))$. In addition, the notation $\widetilde{O}(\cdot)$ (resp. $\widetilde{\Omega}(\cdot)$) is defined in the same way as $O(\cdot)$ (resp. $\Omega(\cdot)$) except that it ignores logarithmic factors.

For any vector $\mathbf{a} = [a_i]_{1 \leq i \leq n} \in \mathbb{R}^n$, we overload the notation $\sqrt{\cdot}$ and $|\cdot|$ in an entry-wise manner such that $\sqrt{\mathbf{a}} := [\sqrt{a_i}]_{1 \leq i \leq n}$ and $|\mathbf{a}| := [|a_i|]_{1 \leq i \leq n}$. For any vectors $\mathbf{a} = [a_i]_{1 \leq i \leq n}$ and $\mathbf{b} = [b_i]_{1 \leq i \leq n}$, the notation $\mathbf{a} \geq \mathbf{b}$ (resp. $\mathbf{a} \leq \mathbf{b}$) means $a_i \geq b_i$ (resp. $a_i \leq b_i$) for all $1 \leq i \leq n$, and we let $\mathbf{a} \circ \mathbf{b} := [a_i b_i]_{1 \leq i \leq n}$ represent the Hadamard product. Additionally, we denote by $\mathbf{1}$ the all-one vector, and \mathbf{I} the identity matrix. For any matrix \mathbf{A} , we define the norm $\|\mathbf{A}\|_1 := \max_i \sum_j |A_{i,j}|$.

3 Main results

3.1 Theoretical guarantees for model-based reinforcement learning

As summarized in Table 1, the theory of all prior work required the sample size per state-action pair to at least exceed $N \geq \Omega(\frac{1}{(1-\gamma)^2})$. In order to break this sample size barrier, we propose to invoke the model-based planning approach applied to an empirical MDP with *perturbed rewards*.

Specifically, for each state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$, let us randomly perturb the immediate reward to obtain

$$r_p(s, a) = r(s, a) + \zeta(s, a), \quad \zeta(s, a) \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(0, \xi), \quad (6)$$

where $\text{Unif}(0, \xi)$ denotes the uniform distribution between 0 and some parameter $\xi > 0$ (to be specified momentarily).² For any policy π , we denote by \widehat{V}_p^π the corresponding value function of the perturbed empirical MDP $\widehat{\mathcal{M}}_p = (\mathcal{S}, \mathcal{A}, \widehat{P}, r_p, \gamma)$ with the probability transition kernel \widehat{P} and the perturbed reward function r_p . Let $\widehat{\pi}_p^*$ represent the optimal policy of $\widehat{\mathcal{M}}_p$, i.e.

$$\widehat{\pi}_p^* := \arg \max_{\pi} \widehat{V}_p^\pi. \quad (7)$$

Encouragingly, this policy results in a value function $V^{\widehat{\pi}_p^*}$ (resp. Q-function $Q^{\widehat{\pi}_p^*}$) that well approximates the true optimal value function V^* (resp. optimal Q-function Q^*), as asserted by the following theorem.

Theorem 1 (Perturbed model-based planning). *There exist some universal constants $c_0, c_1 > 0$ such that: for any $\delta > 0$ and any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the policy $\widehat{\pi}_p^*$ defined in (7) obeys*

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \quad V^{\widehat{\pi}_p^*}(s) \geq V^*(s) - \varepsilon \quad \text{and} \quad Q^{\widehat{\pi}_p^*}(s, a) \geq Q^*(s, a) - \gamma\varepsilon \quad (8)$$

with probability at least $1 - \delta$, provided that the perturbation size is $\xi = \frac{c_1(1-\gamma)\varepsilon}{|\mathcal{S}|^5|\mathcal{A}|^5}$ and that the sample size per state-action pair exceeds

$$N \geq \frac{c_0 \log\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)\varepsilon\delta}\right)}{(1-\gamma)^3\varepsilon^2}. \quad (9)$$

²Note that perturbation is only invoked when running the planning algorithms and does not require collecting new samples.

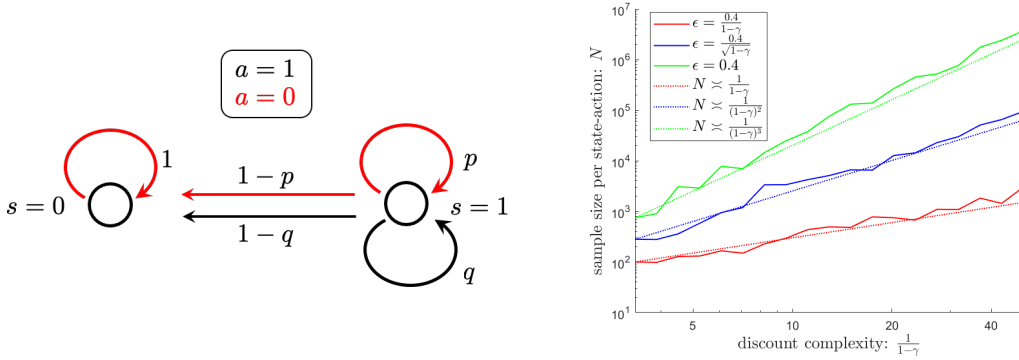


Figure 1: Numerical sample complexity per state-action pair N vs. discount complexity $\frac{1}{1-\gamma}$.

In addition, both the empirical QVI and PI algorithms w.r.t. $\widehat{\mathcal{M}}_p$ (cf. (Azar et al., 2013, Algorithms 1-2)) are able to recover $\widehat{\pi}_p^*$ perfectly within $O(\frac{1}{1-\gamma} \log(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)\varepsilon\delta}))$ iterations.

Remark 1. Theorem 1 holds unchanged if ξ is taken to be $\frac{c_1(1-\gamma)\varepsilon}{|\mathcal{S}||\mathcal{A}|^\alpha}$ for any $\alpha \geq 1$. The current paper picks the specific choice $\alpha = 5$ merely to convey that a very small level of perturbation suffices for our purpose.

Remark 2. Perturbation brings a side benefit: one can recover the optimal policy $\widehat{\pi}_p^*$ of the perturbed empirical MDP $\widehat{\mathcal{M}}_p$ exactly in a small number of iterations without incurring further ptimization errors. To give a flavor of the overall computational complexity, let us take QVI for example Azar et al. (2013). Recall that each iteration of QVI takes time proportional to the time taken to read \widehat{P} (which is a matrix with at most $N|\mathcal{S}||\mathcal{A}|$ nonzeros), hence the resulting computational complexity can be as low as $O(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^4\varepsilon^2} \log^2(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)\varepsilon\delta}))$.

The above theorem demonstrates that: the perturbed model-based planning approach finds an ε -optimal policy as soon as the total sample complexity exceeds the order of $\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}$ (modulo some log factor). It is worth emphasizing that, compared to prior literature, our result imposes no restriction on the range of ε and, in particular, we allow the accuracy level ε to go all the way up to $\frac{1}{1-\gamma}$. Our result is particularly useful in the regime with small-to-moderate sample sizes, since its validity is guaranteed as long as

$$N \geq \widetilde{\Omega}\left(\frac{1}{1-\gamma}\right). \quad (10)$$

Tackling the sample-limited regime (in particular, the scenario when $N \in [\frac{1}{1-\gamma}, \frac{1}{(1-\gamma)^2}]$) requires us to develop new analysis frameworks beyond prior theory, which we shall discuss in detail momentarily.

We remark that the work Azar et al. (2013) established a minimax lower bound of the same order as (9) (up to some log factor) in the regime $\varepsilon = O(1)$. A closer inspection of their analysis, however, reveals that their argument and bound hold true as long as $\varepsilon = O(\frac{1}{1-\gamma})$. This in turn corroborates the *minimax optimality* of our perturbed model-based approach for the full ε -range (which is previously unavailable), and demonstrates the information-theoretic infeasibility to learn a policy strictly better than a random guess if $N \leq \widetilde{O}(\frac{1}{1-\gamma})$. Put another way, the condition (10) contains the full range of “meaningful” sample sizes.

To further demonstrate the effectiveness of our model-based approach, we conduct a series of numerical experiments (motivated by Azar et al. (2013)), focusing on the effect of the discount complexity $\frac{1}{1-\gamma}$ upon the sample complexity. More specifically, consider the following MDP $(\mathcal{S}, \mathcal{A}, P, r, \gamma)$, where $\mathcal{S} = \{0, 1\}$, $\mathcal{A}_0 = \{0\}$, $\mathcal{A}_1 = \{0, 1\}$, $P(0|0, 0) = 1$, $P(1|1, 0) = p$, $P(1|1, 1) = q$ and $r(0, 0) = 0$, $r(1, 0) = r(1, 1) = 1$ (see Fig. 1 for an illustration), and we take the quantities p and q to be $p = \gamma + \frac{2\gamma(1-\gamma)^2\varepsilon}{(1+\gamma)^2}$ and $q = \gamma - \frac{2\gamma(1-\gamma)^2\varepsilon}{(1+\gamma)^2}$ as suggested by Azar et al. (2013). As depicted in Fig. 1, the numerical sample complexity per state-action pair N scales on the order of $\frac{1}{(1-\gamma)^3\varepsilon^2}$ for varying choices of ε , which is consistent with our theory.

Finally, we single out an intermediate result in the analysis of Theorem 1 concerning model-based policy evaluation, which might be of interest on its own. Specifically, for any fixed policy π independent of the data, this task concerns value function estimation via the plug-in estimate \widehat{V}^π (i.e. the value function of

the empirical \mathcal{M} under this policy). However simple as this might seem, existing theoretical underpinnings of this approach remain suboptimal, unless the sample size is already sufficiently large. Our result is the following, which does not require enforcing reward perturbation.

Theorem 2 (Model-based policy evaluation). *Fix any policy π . There exists some universal constant $c_0 > 0$ such that: for any $0 < \delta < 1$ and any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, one has*

$$\forall s \in \mathcal{S} : \quad |\widehat{V}^\pi(s) - V^\pi(s)| \leq \varepsilon \quad (11)$$

with probability at least $1 - \delta$, provided that the sample size per state-action pair exceeds

$$N \geq c_0 \frac{\log\left(\frac{|\mathcal{S}| \log \frac{1}{1-\gamma}}{\delta}\right)}{(1-\gamma)^3 \varepsilon^2}. \quad (12)$$

In words, this theorem reveals that \widehat{V}^π begins to outperform a random guess as soon as $N \geq \tilde{\Omega}\left(\frac{1}{1-\gamma}\right)$. The sample complexity bound (12) enjoys *full coverage* of the ε -range $(0, \frac{1}{1-\gamma}]$, and matches the minimax lower bound derived in (Pananjady and Wainwright, 2019, Theorem 2(b)) up to only a $\log \log \frac{1}{1-\gamma}$ factor. In addition, a recent line of work investigated instance-dependent guarantees for policy evaluation (Khamaru et al. (2020); Pananjady and Wainwright (2019)). While this is not our focus, our analysis does uncover an instance-dependent bound with a broadened sample size range. See Lemma 1 and the discussion thereafter.

3.2 Comparisons with prior work and implications

In order to discuss the novelty of our results in context, we take a moment to compare them with prior theory. See Table 1 for a more complete list of comparisons.

Prior bounds for planning and policy learning. None of the prior results with a generative model (including both model-based or model-free approaches) was capable of efficiently finding the desired policy while accommodating the full sample size range (10). For instance, the state-of-the-art analysis for the model-based approach Agarwal et al. (2019) required the sample size to at least exceed

$$N \geq \tilde{\Omega}\left(\frac{1}{(1-\gamma)^2}\right), \quad (13)$$

whereas the theory for the variance-reduced model-free approach Sidford et al. (2018a); Wainwright (2019b) imposed the sample size requirement

$$N \geq \tilde{\Omega}\left(\frac{1}{(1-\gamma)^3}\right). \quad (14)$$

In fact, it was previously unknown what is achievable in the sample size range $N \in [\frac{1}{1-\gamma}, \frac{1}{(1-\gamma)^2}]$. In contrast, our results confirm the minimax-optimal statistical performance of the model-based approach with full coverage of the ε -range and the sample size range.

Remark 3. We briefly point out why the sample size barrier (13) appeared in the analysis of Agarwal et al. (2019). Take (Agarwal et al., 2019, Section 4.3) for example: the contraction factor $\gamma \sqrt{\frac{8 \log(|\mathcal{S}||\mathcal{A}|/(1-\gamma)\delta)}{N}} \frac{1}{1-\gamma}$ therein needs to be smaller than 1, thereby requiring $N \geq \tilde{\Omega}((1-\gamma)^{-2})$.

Prior bounds for policy evaluation. Regarding value function estimation for any fixed policy π , the prior results Agarwal et al. (2019); Azar et al. (2013); Pananjady and Wainwright (2019) for the plug-in approach all operated under the assumption that $N \geq \tilde{\Omega}\left(\frac{1}{(1-\gamma)^2}\right)$, which is more stringent than our result by a factor of at least $\frac{1}{1-\gamma}$. In addition, our sample complexity matches the state-of-the-art guarantees in the regime where $\varepsilon \leq \frac{1}{\sqrt{1-\gamma}}$ (Agarwal et al., 2019; Pananjady and Wainwright, 2019), while extending them to the range $\varepsilon \in [\frac{1}{\sqrt{1-\gamma}}, \frac{1}{1-\gamma}]$ uncovered in these previous papers.

4 Other related work

Classical analyses of reinforcement learning algorithms have largely focused on asymptotic performance (e.g. Jaakkola et al. (1994); Szepesvári (1998); Tsitsiklis and Van Roy (1997); Tsitsiklis (1994)). Leveraging the toolkit of concentration inequalities, a number of recent papers have shifted attention towards understanding the performance in the non-asymptotic and finite-time settings. A highly incomplete list includes Azar et al. (2017); Beck and Srikant (2012); Bhandari et al. (2018); Bradtke and Barto (1996); Cai et al. (2019); Chen et al. (2020); Dalal et al. (2018); Even-Dar and Mansour (2003); Fan et al. (2019); Gupta et al. (2019); Jin et al. (2018); Kaledin et al. (2020); Kearns and Singh (1999); Khamaru et al. (2020); Lakshminarayanan and Szepesvari (2018); Li et al. (2020); Mou et al. (2020); Qu and Wierman (2020); Shah and Xie (2018); Sidford et al. (2018a); Srikant and Ying (2019); Strehl et al. (2006); Wainwright (2019b); Xu and Gu (2020); Xu et al. (2019), a large fraction of which is concerned with model-free algorithms.

The generative model (or simulator) adopted in this paper was first proposed in Kearns and Singh (1999), which has been invoked in Agarwal et al. (2019); Azar et al. (2012, 2013); Kakade (2003); Kearns et al. (2002); Kearns and Singh (1999); Khamaru et al. (2020); Lattimore and Hutter (2012); Pananjady and Wainwright (2019); Sidford et al. (2018a,b); Wainwright (2019b); Wang (2019); Yang and Wang (2019), to name just a few. In particular, Azar et al. (2013) developed the minimax lower bound on the sample complexity $N = \Omega\left(\frac{|\mathcal{S}||\mathcal{A}|\log(|\mathcal{S}||\mathcal{A}|)}{(1-\gamma)^3\varepsilon^2}\right)$ necessary for finding an ε -optimal policy, and showed that, for any $\varepsilon \in (0, 1)$, a model-based approach (e.g. applying QVI or PI to the empirical MDP) can estimate the optimal Q-function to within an ε -accuracy given near-minimal samples. Note, however, that directly translating this result to the policy guarantees leads to an additional factor of $\frac{1}{1-\gamma}$ in estimation accuracy and of $\frac{1}{(1-\gamma)^2}$ in sample complexity. In light of this, Azar et al. (2013) further showed that a near-optimal sample complexity is possible for policy learning if the sample size is at least on the order of $\frac{|\mathcal{S}|^2|\mathcal{A}|}{(1-\gamma)^2}$ which, however, is no longer sub-linear in the model complexity. A recent breakthrough Agarwal et al. (2019) substantially improved the model-based guarantee with the aid of auxiliary state-absorbing MDPs, extending the range of sample complexity to $\left[\frac{|\mathcal{S}||\mathcal{A}|\log(|\mathcal{S}||\mathcal{A}|)}{(1-\gamma)^2}, \infty\right)$. Our analysis is motivated in part by Agarwal et al. (2019), but also relies on several other novel techniques to complete the picture.

Finally, we remark that the construction of state-absorbing MDPs or state-action-absorbing MDPs falls under the category of “leave-one-out” type analysis, which is particularly effective in decoupling complicated statistical dependency in various statistical estimation problems Agarwal et al. (2019); Chen et al. (2019); El Karoui (2015); Ma et al. (2020); Pananjady and Wainwright (2019). The application of such an analysis framework to MDPs should be attributed to Agarwal et al. (2019). Other applications to Markov chains include Chen et al. (2019); Pananjady and Wainwright (2019).

5 Analysis

This section presents the key ideas for proving our main results, following an introduction of some convenient matrix notation.

5.1 Matrix notation and Bellman equations

It is convenient to present our proof based on some matrix notation for MDPs. Denoting by $e_1, \dots, e_{|\mathcal{S}|} \in \mathbb{R}^{|\mathcal{S}|}$ the standard basis vectors, we can define:

- $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$: a vector representing the reward function r (so that $r_{(s,a)} = r(s, a)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$).
- $\mathbf{V}^\pi \in \mathbb{R}^{|\mathcal{S}|}$: a vector representing the value function V^π (so that $V_s^\pi = V^\pi(s)$ for all $s \in \mathcal{S}$).
- $\mathbf{Q}^\pi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$: a vector representing the Q-function Q^π (so that $Q_{(s,a)}^\pi = Q^\pi(s, a)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$).
- $\mathbf{V}^* \in \mathbb{R}^{|\mathcal{S}|}$ and $\mathbf{Q}^* \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$: representing the optimal value function V^* and optimal Q-function Q^* .
- $\mathbf{P} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$: a matrix representing the probability transition kernel P , where the (s, a) -th row of \mathbf{P} is a probability vector representing $P(\cdot|s, a)$. Denote $\mathbf{P}_{s,a}$ as the (s, a) -th row of the transition matrix \mathbf{P} .

- $\mathbf{\Pi}^\pi \in \{0, 1\}^{|\mathcal{S}| \times |\mathcal{S}| \times |\mathcal{A}|}$: a projection matrix associated with a given policy π taking the following form

$$\mathbf{\Pi}^\pi = \begin{pmatrix} \mathbf{e}_{\pi(1)}^\top & & & \\ & \mathbf{e}_{\pi(2)}^\top & & \\ & & \ddots & \\ & & & \mathbf{e}_{\pi(|\mathcal{S}|)}^\top \end{pmatrix}. \quad (15)$$

- $\mathbf{P}^\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}| \times |\mathcal{S}| \times |\mathcal{A}|}$ and $\mathbf{P}_\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$: two *square* probability transition matrices induced by the policy π over the state-action pairs and the states respectively, defined by

$$\mathbf{P}^\pi := \mathbf{P}\mathbf{\Pi}^\pi \quad \text{and} \quad \mathbf{P}_\pi := \mathbf{\Pi}^\pi \mathbf{P}. \quad (16)$$

- $\mathbf{r}_\pi \in \mathbb{R}^{|\mathcal{S}|}$: a reward vector restricted to the actions chosen by the policy π , namely, $r_\pi(s) = r(s, \pi(s))$ for all $s \in \mathcal{S}$ (or simply, $\mathbf{r}_\pi = \mathbf{\Pi}^\pi \mathbf{r}$).

Armed with the above matrix notation, we can write, for any policy π , the *Bellman consistency equation* as

$$\mathbf{Q}^\pi = \mathbf{r} + \gamma \mathbf{P}\mathbf{V}^\pi = \mathbf{r} + \gamma \mathbf{P}^\pi \mathbf{Q}^\pi, \quad (17)$$

which implies that

$$\mathbf{Q}^\pi = (\mathbf{I} - \gamma \mathbf{P}^\pi)^{-1} \mathbf{r}; \quad (18)$$

$$\mathbf{V}^\pi = \mathbf{r}_\pi + \gamma \mathbf{P}_\pi \mathbf{V}^\pi \quad \text{and} \quad \mathbf{V}^\pi = (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{r}_\pi. \quad (19)$$

For a vector $\mathbf{V} = [V_i]_{1 \leq i \leq |\mathcal{S}|} \in \mathbb{R}^{|\mathcal{S}|}$, we define the vector $\text{Var}_{\mathbf{P}}(\mathbf{V}) \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$ whose entries are given by

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \quad [\text{Var}_{\mathbf{P}}(\mathbf{V})]_{(s,a)} := \sum_{s' \in \mathcal{S}} P(s'|s, a) V_{s'}^2 - \left(\sum_{s' \in \mathcal{S}} P(s'|s, a) V_{s'} \right)^2,$$

i.e. the variance of \mathbf{V} w.r.t. $P(\cdot|s, a)$. This can be expressed using our matrix notation as follows

$$\text{Var}_{\mathbf{P}}(\mathbf{V}) = \mathbf{P}(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}\mathbf{V}) \circ (\mathbf{P}\mathbf{V}).$$

Similarly, for any given policy π we define

$$\text{Var}_{\mathbf{P}_\pi}(\mathbf{V}) = \mathbf{\Pi}^\pi \text{Var}_{\mathbf{P}}(\mathbf{V}) = \mathbf{P}_\pi(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}_\pi \mathbf{V}) \circ (\mathbf{P}_\pi \mathbf{V}) \in \mathbb{R}^{|\mathcal{S}|}.$$

We shall also define $\widehat{\mathbf{V}}^\pi$, $\widehat{\mathbf{Q}}^\pi$, $\widehat{\mathbf{V}}^*$, $\widehat{\mathbf{Q}}^*$, $\widehat{\mathbf{P}}$, $\widehat{\mathbf{P}}^\pi$, $\widehat{\mathbf{P}}_\pi$, $\text{Var}_{\widehat{\mathbf{P}}}(\mathbf{V})$, $\text{Var}_{\widehat{\mathbf{P}}_\pi}(\mathbf{V})$ w.r.t. the empirical MDP $\widehat{\mathcal{M}}$ in an analogous fashion.

5.2 Analysis: model-based policy evaluation

We start with the simpler task of policy evaluation, which also plays a crucial role in the analysis of planning. To establish our guarantees in Theorem 2, we aim to prove the following result. Here, we recall that the true value function under a policy π and the model-based empirical estimate are given respectively by

$$\mathbf{V}^\pi = (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{r}_\pi \quad \text{and} \quad \widehat{\mathbf{V}}^\pi = (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} \mathbf{r}_\pi. \quad (20)$$

Lemma 1. *Fix any policy π . Consider any $0 < \delta < 1$, and suppose $N \geq \frac{32e^2}{1-\gamma} \log\left(\frac{4|\mathcal{S}| \log(\frac{e}{1-\gamma})}{\delta}\right)$. Then with probability at least $1 - \delta$, the vectors defined in (20) obey*

$$\begin{aligned} \|\widehat{\mathbf{V}}^\pi - \mathbf{V}^\pi\|_\infty &\leq 4\gamma \sqrt{\frac{2 \log\left(\frac{4|\mathcal{S}| \log(\frac{e}{1-\gamma})}{\delta}\right)}{N}} \left\| (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^\pi]} \right\|_\infty + \frac{2\gamma \log\left(\frac{4|\mathcal{S}| \log(\frac{e}{1-\gamma})}{\delta}\right)}{(1-\gamma)N} \|\mathbf{V}^\pi\|_\infty \\ &\leq 6 \sqrt{\frac{2 \log\left(\frac{4|\mathcal{S}| \log(\frac{e}{1-\gamma})}{\delta}\right)}{N(1-\gamma)^3}}. \end{aligned} \quad (21)$$

Proof. The key proof idea is to resort to a high-order successive expansion of $\widehat{\mathbf{V}}^\pi - \mathbf{V}^\pi$, followed by fine-grained analysis of each term up to a certain logarithmic order. See Appendix B.1. \square

Clearly, Theorem 2 is a straightforward consequence of Lemma 1. Further, we strengthen the result by providing an additional instance-dependent bound (see the first line of (21) that depends on the true instance $\mathbf{P}_\pi, \mathbf{V}^\pi$), which is often tighter than the worst-case bound stated in the second line of (21). Our contribution can be better understood when compared with Pananjady and Wainwright (2019). Assuming that there is no noise in the rewards, our instance-dependent guarantee matches Pananjady and Wainwright (2019, Theorem 1(a)) up to some $\log \log \frac{1}{1-\gamma}$ factor, while being capable of covering the full sample size range $N \geq \widetilde{\Omega}(\frac{1}{1-\gamma})$. In contrast, Pananjady and Wainwright (2019, Theorem 1) is only valid when $N \geq \widetilde{\Omega}(\frac{1}{(1-\gamma)^2})$.

Proof ideas. We now briefly and informally describe the key proof ideas. As a starting point, the elementary identities (20) allow us to obtain

$$\begin{aligned} \widehat{\mathbf{V}}^\pi - \mathbf{V}^\pi &= (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} \mathbf{r}_\pi - \mathbf{V}^\pi \\ &= (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} (\mathbf{I} - \gamma \mathbf{P}_\pi) \mathbf{V}^\pi - (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi) \mathbf{V}^\pi \\ &= \gamma (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^\pi. \end{aligned} \quad (22)$$

Due to the complicated dependency between $(\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1}$ and $(\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^\pi$, a natural strategy is to control these two terms separately and then to combine bounds; see Agarwal et al. (2019, Lemma 5) for an introduction. This simple approach, however, leads to sub-optimal statistical guarantees.

In order to refine the statistical analysis, we propose to further expand (22) in a similar way to deduce

$$\begin{aligned} (22) &= \gamma (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^\pi + \gamma \left\{ (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} - (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \right\} (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^\pi \\ &= \gamma (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^\pi + \gamma^2 (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^\pi, \end{aligned} \quad (23)$$

where the last line holds due to the same reason as (22) (basically it can be seen by replacing \mathbf{r}_π with $(\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^\pi$ in (22)). This can be viewed as a “second-order” expansion, with (22) being a “first-order” counterpart. The advantage is that: the first term in (23) becomes easier to cope with than its counterpart (22), owing to the independence between $(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1}$ and $(\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^\pi$. However, the second term in (23) remains difficult to control optimally. To remedy this issue, we shall continue to expand it to higher order (up to some logarithmic order), which eventually allows for optimal control of the estimation error.

Another crucial issue is that: in order to obtain fine-grained analyses on each term in the expansion (except for the first-order term), a common approach is to combine the Bernstein inequality with a classical entrywise bound on a quantity taking the form $(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})}$ (which dates back to Azar et al. (2013)). Such a classical bound in prior literature, however, is not sufficiently tight for our purpose, which calls for refinement; see Lemma 8. Details are deferred to Appendix B.1.

5.3 Analysis: model-based planning

This subsection moves on to establishing our theory for model-based planning (cf. Theorem 1) and outlines the key ideas. In what follows, we shall start by analyzing the unperturbed version, which will elucidate the role of reward perturbation in our analysis.

We first make note of the following elementary decomposition:

$$\begin{aligned} \mathbf{V}^* - \mathbf{V}^{\widehat{\pi}^*} &= (\widehat{\mathbf{V}}^{\widehat{\pi}^*} - \mathbf{V}^{\widehat{\pi}^*}) + (\widehat{\mathbf{V}}^{\pi^*} - \widehat{\mathbf{V}}^{\widehat{\pi}^*}) + (\mathbf{V}^* - \widehat{\mathbf{V}}^{\pi^*}) \\ &\leq (\widehat{\mathbf{V}}^{\widehat{\pi}^*} - \mathbf{V}^{\widehat{\pi}^*}) + (\mathbf{V}^{\pi^*} - \widehat{\mathbf{V}}^{\pi^*}), \end{aligned} \quad (24)$$

where the inequality follows from the optimality of $\widehat{\pi}^*$ w.r.t. $\widehat{\mathbf{V}}$ (so that $\widehat{\mathbf{V}}^{\pi^*} \leq \widehat{\mathbf{V}}^{\widehat{\pi}^*}$) and the definition $\mathbf{V}^* = \mathbf{V}^{\pi^*}$. This leaves us with two terms to control.

Step 1: bounding $\|\mathbf{V}^{\pi^*} - \widehat{\mathbf{V}}^{\pi^*}\|_\infty$. Given that π^* is independent of the data, we can carry out this step using Lemma 1. Specifically, taking $\pi = \pi^*$ in Lemma 1 yields that, with probability at least $1 - \delta$,

$$\|\widehat{\mathbf{V}}^{\pi^*} - \mathbf{V}^{\pi^*}\|_\infty \leq 6\sqrt{\frac{2\log\left(\frac{4|\mathcal{S}|\log\frac{e}{1-\gamma}}{\delta}\right)}{N(1-\gamma)^3}}. \quad (25)$$

Step 2: bounding $\|\widehat{\mathbf{V}}^{\widehat{\pi}^*} - \mathbf{V}^{\widehat{\pi}^*}\|_\infty$. Extending the result in Step 1 to $\|\widehat{\mathbf{V}}^{\widehat{\pi}^*} - \mathbf{V}^{\widehat{\pi}^*}\|_\infty$ is considerably more challenging, primarily due to the complicated statistical dependency between $(\mathbf{V}^{\widehat{\pi}^*}, \widehat{\mathbf{V}}^{\widehat{\pi}^*})$ and the data matrix $\widehat{\mathbf{P}}$. The recent work Agarwal et al. (2019) developed a clever “leave-one-out” type argument by constructing some auxiliary state-absorbing MDPs to decouple the statistical dependency when $\varepsilon < 1/\sqrt{1-\gamma}$. However, their argument falls short of accommodating the full range of ε . To address this challenge, our analysis consists of the following two steps, both of which require new ideas beyond Agarwal et al. (2019).

- *Decoupling statistical dependency between $\widehat{\pi}^*$ and $\widehat{\mathbf{P}}$.* Instead of attempting to decouple the statistical dependency between $\widehat{\mathbf{V}}^{\widehat{\pi}^*}$ and $\widehat{\mathbf{P}}$ as in Agarwal et al. (2019), we focus on decoupling the statistical dependency between the policy $\widehat{\pi}^*$ and $\widehat{\mathbf{P}}$. If this can be achieved, then the proof strategy adopted in Step 1 for a fixed policy becomes applicable (see Section 5.3.1). A key ingredient of this step lies in the construction of a collection of auxiliary state-action-absorbing MDPs (motivated by Agarwal et al. (2019)), which allows us to get hold of $\|\mathbf{V}^{\widehat{\pi}^*} - \widehat{\mathbf{V}}^{\widehat{\pi}^*}\|_\infty$. See Section 5.3.2 for details, with a formal bound delivered in Lemma 5.
- *Tie-breaking via reward perturbation.* A shortcoming of the above-mentioned approach, however, is that it relies crucially on the separability of $\widehat{\pi}^*$ from other policies; in other words, the proof might fail if $\widehat{\pi}^*$ is non-unique or not sufficiently distinguishable from others. Consequently, it remains to ensure that the optimal policy $\widehat{\pi}^*$ stands out from all the rest for all MDPs of interest. As it turns out, this can be guaranteed with high probability by slightly perturbing the reward function so as to break the ties. See Section 5.3.3 for details.

In the sequel, we shall flesh out these key ideas.

5.3.1 Value function estimation for a policy obeying Bernstein-type conditions

Before discussing how to decouple statistical dependency, we record a useful result that plays an important role in the analysis. Specifically, Lemma 1 can be generalized beyond the family of fixed policies (namely, those independent of $\widehat{\mathbf{P}}$), as long as a certain Bernstein-type condition — to be formalized in (27) — is satisfied. To make it precise, we need to introduce a set of auxiliary vectors as follows

$$\begin{aligned} \mathbf{r}^{(0)} &:= \mathbf{r}_\pi, & \mathbf{V}^{(0)} &:= (\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\mathbf{r}^{(0)}, \\ \mathbf{r}^{(l)} &:= \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(l-1)}]}, & \mathbf{V}^{(l)} &:= (\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\mathbf{r}^{(l)}, \quad l \geq 1. \end{aligned} \quad (26)$$

Our generalization of Lemma 1 is as follows, which does *not* require statistical independence between the policy π and the data $\widehat{\mathbf{P}}$. Here, we remind the reader of the notation $|\mathbf{z}| := [z_1, \dots, z_n]^\top$ and $\sqrt{\mathbf{z}} := [\sqrt{z_1}, \dots, \sqrt{z_n}]^\top$ for any vector $\mathbf{z} \in \mathbb{R}^n$.

Lemma 2. *Suppose that there exists some quantity $\beta_1 > 0$ such that $\{\mathbf{V}^{(l)}\}$ (cf. (26)) obeys*

$$\left| (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi)\mathbf{V}^{(l)} \right| \leq \sqrt{\frac{\beta_1}{N}} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(l)}]} + \frac{\beta_1 \|\mathbf{V}^{(l)}\|_\infty}{N} \mathbf{1}, \quad \text{for all } 0 \leq l \leq \log\left(\frac{e}{1-\gamma}\right). \quad (27)$$

Suppose that $N > \frac{16e^2}{1-\gamma}\beta_1$. Then the vectors $\mathbf{V}^\pi = (\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\mathbf{r}_\pi$ and $\widehat{\mathbf{V}}^\pi = (\mathbf{I} - \gamma\widehat{\mathbf{P}}_\pi)^{-1}\mathbf{r}_\pi$ satisfy

$$\|\widehat{\mathbf{V}}^\pi - \mathbf{V}^\pi\|_\infty \leq \frac{6}{1-\gamma} \sqrt{\frac{\beta_1}{N(1-\gamma)}}. \quad (28)$$

While the Bernstein-type condition (27) clearly holds for some reasonably small β_1 if π is independent of $\widehat{\mathbf{P}}$, it might remain valid if π exhibits fairly “weak” statistical dependency on the data samples. This is a key step that paves the way for our subsequent analysis of $\widehat{\pi}^*$.

5.3.2 Decoupling statistical dependency via (s, a) -absorbing MDPs

We are now positioned to demonstrate how to control $\|\widehat{\mathbf{V}}^{\widehat{\pi}^*} - \mathbf{V}^{\widehat{\pi}^*}\|_\infty$ w.r.t. the optimal policy $\widehat{\pi}^*$ to $\widehat{\mathbf{V}}$. A crucial technical challenge lies in how to decouple the complicated statistical dependency between the optimal policy $\widehat{\pi}^*$ and the $\widehat{\mathbf{V}}^*$ (which heavily relies on the data samples). Towards this, we resort to a leave-one-row-out argument built upon a collection of auxiliary MDPs, largely motivated by the novel construction in (Agarwal et al., 2019, Section 4.2). In comparison to Agarwal et al. (2019) that introduces state-absorbing MDPs (so that a state s is absorbing regardless of the subsequent actions chosen), our construction is a set of state-action-absorbing MDPs, in which a state s is absorbing only when a designated action a is always executed at the state s .

Construction of (s, a) -absorbing MDPs. For each state-action pair (s, a) and each scalar u with $|u| \leq 1/(1 - \gamma)$, we construct an auxiliary MDP $\mathcal{M}_{s,a,u}$ — it is identical to the original \mathcal{M} except that it is absorbing in state s if we always choose action a in state s . More specifically, the probability transition kernel associated with $\mathcal{M}_{s,a,u}$ (denoted by $P_{\mathcal{M}_{s,a,u}}$) can be specified by

$$\begin{aligned} P_{\mathcal{M}_{s,a,u}}(s | s, a) &= 1, \\ P_{\mathcal{M}_{s,a,u}}(s' | s, a) &= 0, & \text{for all } s' \neq s, \\ P_{\mathcal{M}_{s,a,u}}(\cdot | s', a') &= P_{\mathcal{M}}(\cdot | s', a'), & \text{for all } (s', a') \neq (s, a), \end{aligned} \quad (29)$$

where $P_{\mathcal{M}}$ is the probability transition kernel w.r.t. the original \mathcal{M} . Meanwhile, the instant reward received at (s, a) in $\mathcal{M}_{s,a,u}$ is set to be u , while the rewards at all other state-action pairs stay unchanged. We can define $\widehat{\mathcal{M}}_{s,a,u}$ analogously (so that its probability transition matrix is identical to $\widehat{\mathbf{P}}$ except that the (s, a) -th row becomes absorbing). The main advantage of this construction is that: for any fixed u , the MDP $\widehat{\mathcal{M}}_{s,a,u}$ is statistically independent of $\widehat{\mathbf{P}}_{s,a}$ (the row of $\widehat{\mathbf{P}}$ corresponding to the state-action pair (s, a) , determined by the samples collected for the (s, a) pair).

To streamline notation, we let $\mathbf{Q}_{s,a,u}^\pi$ represent the Q-function of $\mathcal{M}_{s,a,u}$ under a policy π , denote by $\pi_{s,a,u}^*$ the optimal policy associated with $\mathcal{M}_{s,a,u}$, and let $\mathbf{Q}_{s,a,u}^*$ be the Q-function under this optimal policy $\pi_{s,a,u}^*$. The notation $\mathbf{V}_{s,a,u}^\pi$ and $\mathbf{V}_{s,a,u}^*$ regarding value functions, as well as their counterparts (i.e. $\widehat{\mathbf{Q}}_{s,a,u}^\pi$, $\widehat{\mathbf{Q}}_{s,a,u}^*$, $\widehat{\mathbf{V}}_{s,a,u}^\pi$, $\widehat{\mathbf{V}}_{s,a,u}^*$, $\widehat{\pi}_{s,a,u}^*$) in the empirical MDP $\widehat{\mathcal{M}}$, can be defined in an analogous fashion.

Remark 4. The careful reader will remark that the instant reward u is constrained to reside within $[-\frac{1}{1-\gamma}, \frac{1}{1-\gamma}]$ rather than the usual range $[0, 1]$. Fortunately, none of the subsequent steps that involve u requires u to lie within $[0, 1]$.

Intimate connections between the auxiliary MDPs and the original MDP. In the following, we introduce a result that connects the Q-function and the value function of the absorbing MDP with those of the original MDP. The idea is motivated by Agarwal et al. (2019, Lemma 7).

Lemma 3. *Setting $u^* := r(s, a) + \gamma(\mathbf{P}\mathbf{V}^*)_{s,a} - \gamma\mathbf{V}^*(s)$, one has*

$$\mathbf{Q}_{s,a,u^*}^* = \mathbf{Q}^* \quad \text{and} \quad \mathbf{V}_{s,a,u^*}^* = \mathbf{V}^*. \quad (30)$$

Proof. See Appendix B.2. □

Remark 5. Lemma 3 does not rely on the particular form of \mathbf{P} , and can be directly generalized to the empirical model $\widehat{\mathbf{P}}$ and the auxiliary MDPs built upon $\widehat{\mathbf{P}}$.

In words, by properly setting the instant reward $u = u^*$ (which can be easily shown to reside within $[-\frac{1}{1-\gamma}, \frac{1}{1-\gamma}]$), one guarantees that the (s, a) -absorbing MDP and the original MDP have the same Q-function and value function under the respective optimal policies.

Representing $\widehat{\pi}^*$ via a small set of policies independent of $\widehat{\mathbf{P}}_{s,a}$. With Lemma 3 in place, it is tempting to use $\widehat{\mathcal{M}}_{s,a,\widehat{u}^*}$ with $\widehat{u}^* := r(s, a) + \gamma(\widehat{\mathbf{P}}\widehat{\mathbf{V}}^*)_{s,a} - \gamma\widehat{\mathbf{V}}^*(s)$ to replace the original $\widehat{\mathcal{M}}$. The rationale is simple: given that the probability transition matrix of $\widehat{\mathcal{M}}_{s,a,\widehat{u}^*}$ does not rely upon $\widehat{\mathbf{P}}_{s,a}$, the statistical

dependency between $\widehat{\mathcal{M}}_{s,a,\widehat{u}^*}$ and $\widehat{\mathbf{P}}_{s,a}$ is now fully embedded into a single parameter \widehat{u}^* . This motivates us to decouple the statistical dependency effectively by constructing an epsilon-net (see, e.g., [Vershynin \(2018\)](#)) w.r.t. this single parameter. The aim is to locate a point u_0 over a small fixed set such that (i) it is close to \widehat{u}^* , and (ii) its associated optimal policy is identical to the original optimal policy $\widehat{\pi}^*$.

It turns out that this aim can be accomplished as long as the original Q-function $\widehat{\mathbf{Q}}^*$ satisfies a sort of separation condition (which indicates that there is no tie when it comes to the optimal policy). To make it precise, given any $0 < \omega < 1$, our separation condition is characterized through the following event

$$\mathcal{B}_\omega := \left\{ \widehat{\mathbf{Q}}^*(s, \widehat{\pi}^*(s)) - \max_{a: a \neq \widehat{\pi}^*(s)} \widehat{\mathbf{Q}}^*(s, a) \geq \omega \text{ for all } s \in \mathcal{S} \right\}. \quad (31)$$

Clearly, on the event \mathcal{B}_ω , the optimal policy $\widehat{\pi}^*$ is unique, since for each s the action $\widehat{\pi}^*(s)$ results in a strictly higher Q-value compared to any other action. With this separation condition in mind, our result is stated below. Here and throughout, we define an epsilon-net of the interval $[-\frac{1}{1-\gamma}, \frac{1}{1-\gamma}]$ as follows

$$\mathcal{N}_\epsilon := \left\{ -n_\epsilon\epsilon, \dots, -\epsilon, 0, \epsilon, \dots, n_\epsilon\epsilon \right\}, \quad \text{for the largest integer } n_\epsilon \text{ obeying } n_\epsilon\epsilon < \frac{1}{1-\gamma}, \quad (32)$$

which has cardinality at most $\frac{2}{(1-\gamma)\epsilon}$.

Lemma 4. *Consider any $\omega > 0$, and suppose the event \mathcal{B}_ω (cf. (31)) holds. Then for any pair $(s, a) \in \mathcal{S} \times \mathcal{A}$, there exists a point $u_0 \in \mathcal{N}_{(1-\gamma)\omega/4}$, such that*

$$\widehat{\pi}^* = \widehat{\pi}_{s,a,u_0}^*. \quad (33)$$

Proof. See Appendix B.3. □

Deriving an optimal error bound under the separation condition. Armed with the above bounds, we are ready to derive the desired error bound by combining Lemma 2 and Lemma 4.

Lemma 5. *Given $0 < \omega < 1$ and $\delta > 0$, suppose that \mathcal{B}_ω (defined in (31)) occurs with probability at least $1 - \delta$. Then with probability at least $1 - 3\delta$,*

$$\|\widehat{\mathbf{V}}^{\widehat{\pi}^*} - \mathbf{V}^{\widehat{\pi}^*}\|_\infty \leq 6\sqrt{\frac{2 \log\left(\frac{32|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\omega\delta}\right)}{N(1-\gamma)^3}} \quad \text{and} \quad \mathbf{V}^* - \mathbf{V}^{\widehat{\pi}^*} \leq 12\sqrt{\frac{2 \log\left(\frac{32|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\omega\delta}\right)}{N(1-\gamma)^3}} \mathbf{1}, \quad (34)$$

provided that $N \geq \frac{c_0 \log\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)\delta\omega}\right)}{1-\gamma}$ for some sufficiently large constant $c_0 > 0$.

Proof. See Appendix B.4. □

5.3.3 A tie-breaking argument

Unfortunately, the separation condition specified in \mathcal{B}_ω (cf. (31)) does not always hold. In order to accommodate all possible MDPs of interest without imposing such a special separation condition, we put forward a perturbation argument allowing one to generate a new MDP that (i) satisfies the separation condition, and that (ii) is sufficiently close to the original MDP.

Specifically, let us represent the proposed reward perturbation (6) in a vector form as follows

$$\mathbf{r}_p := \mathbf{r} + \boldsymbol{\zeta}, \quad (35)$$

where $\boldsymbol{\zeta} = [\zeta(s, a)]_{(s,a) \in \mathcal{S} \times \mathcal{A}}$ is an $|\mathcal{S}||\mathcal{A}|$ -dimensional vector composed of independent entries with each $\zeta(s, a) \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(0, \xi)$. We aim to show that: by randomly perturbing the reward function, we can “break the tie” in the Q-function and ensure sufficient separation of Q-values associated with different actions.

To formalize our result, we find it convenient to introduce additional notation. Denote by π_p^* the optimal policy of the MDP $\mathcal{M}_p = (\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r}_p, \gamma)$, and Q_p^* its optimal state-action value function. We can define \widehat{Q}_p^* and $\widehat{\pi}_p^*$ analogously for the MDP $\widehat{\mathcal{M}}_p = (\mathcal{S}, \mathcal{A}, \widehat{\mathbf{P}}, \mathbf{r}_p, \gamma)$. Our result is phrased as follows.

Lemma 6. Consider the perturbed reward vector defined in expression (35). With probability at least $1 - \delta$,

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} \text{ with } a \neq \pi_p^*(s) : \quad Q_p^*(s, \pi_p^*(s)) - Q_p^*(s, a) > \frac{\xi \delta (1 - \gamma)}{3|\mathcal{S}||\mathcal{A}|^2}. \quad (36)$$

This result holds unchanged if (Q_p^*, π_p^*) is replaced by $(\widehat{Q}_p^*, \widehat{\pi}_p^*)$.

Proof. See Appendix B.5. □

Lemma 6 reveals that at least a polynomially small degree of separation ($\omega = \frac{\xi \delta (1 - \gamma)}{3|\mathcal{S}||\mathcal{A}|^2}$) arises upon random perturbation (with size ξ) of the reward function. As we shall see momentarily, this level of separation suffices for our purpose.

5.3.4 Proof of Theorem 1

Let us consider the randomly perturbed reward function as in (35). For any policy π , we denote by \mathbf{V}_p^π (resp. $\widehat{\mathbf{V}}_p^\pi$) the corresponding value function vector in the MDP with probability transition matrix \mathbf{P} (resp. $\widehat{\mathbf{P}}$) and reward vector \mathbf{r}_p . Note that π_p^* (resp. $\widehat{\pi}_p^*$) denotes the optimal policy that maximizes \mathbf{V}_p^π (resp. $\widehat{\mathbf{V}}_p^\pi$).

In view of Lemma 6, with probability at least $1 - \delta$ one has the separation

$$\left| \widehat{Q}_p^*(s, \widehat{\pi}_p^*(s)) - \widehat{Q}_p^*(s, a) \right| > \frac{\xi \delta (1 - \gamma)}{3|\mathcal{S}||\mathcal{A}|^2} \quad (37)$$

uniformly over all s and $a \neq \widehat{\pi}_p^*(s)$. With this separation in place, taking $\omega := \frac{\xi \delta (1 - \gamma)}{3|\mathcal{S}||\mathcal{A}|^2}$ in Lemma 5 yields

$$\left\| \mathbf{V}_p^{\pi_p^*} - \mathbf{V}_p^{\widehat{\pi}_p^*} \right\|_\infty \leq 12 \sqrt{\frac{2 \log \left(\frac{96|\mathcal{S}|^2|\mathcal{A}|^3}{(1-\gamma)^4 \xi \delta^2} \right)}{N(1-\gamma)^3}}. \quad (38)$$

In addition, the value functions under any policy π obeys

$$\mathbf{V}^\pi - \mathbf{V}_p^\pi = \mathbf{\Pi}^\pi \left((\mathbf{I} - \mathbf{P}^\pi)^{-1} \mathbf{r} - (\mathbf{I} - \mathbf{P}^\pi)^{-1} \mathbf{r}_p \right),$$

which taken collectively with the facts $\|\mathbf{r} - \mathbf{r}_p\|_\infty \leq \xi$ and $\|(\mathbf{I} - \gamma \mathbf{P}^\pi)^{-1}\|_1 \leq \frac{1}{1-\gamma}$ gives

$$\left\| \mathbf{V}^\pi - \mathbf{V}_p^\pi \right\|_\infty \leq \|(\mathbf{I} - \gamma \mathbf{P}^\pi)^{-1}\|_1 \|\mathbf{r} - \mathbf{r}_p\|_\infty \leq \frac{1}{1-\gamma} \xi.$$

Specializing the above relation to π^* and $\widehat{\pi}_p^*$ gives

$$\left\| \mathbf{V}^{\pi^*} - \mathbf{V}_p^{\pi^*} \right\|_\infty \leq \frac{1}{1-\gamma} \xi \quad \text{and} \quad \left\| \mathbf{V}^{\widehat{\pi}_p^*} - \mathbf{V}_p^{\widehat{\pi}_p^*} \right\|_\infty \leq \frac{1}{1-\gamma} \xi. \quad (39)$$

Now let us consider the following decomposition

$$\begin{aligned} \mathbf{V}^{\widehat{\pi}_p^*} - \mathbf{V}^* &= (\mathbf{V}^{\widehat{\pi}_p^*} - \mathbf{V}_p^{\widehat{\pi}_p^*}) + (\mathbf{V}_p^{\widehat{\pi}_p^*} - \mathbf{V}_p^{\pi_p^*}) + (\mathbf{V}_p^{\pi_p^*} - \mathbf{V}_p^{\pi^*}) + (\mathbf{V}_p^{\pi^*} - \mathbf{V}^*) \\ &\geq (\mathbf{V}^{\widehat{\pi}_p^*} - \mathbf{V}_p^{\widehat{\pi}_p^*}) + (\mathbf{V}_p^{\widehat{\pi}_p^*} - \mathbf{V}_p^{\pi_p^*}) + (\mathbf{V}_p^{\pi^*} - \mathbf{V}^*), \end{aligned}$$

where the last step follows from the optimality of π_p^* w.r.t. \mathbf{V}_p . Taking this collectively with the inequalities (38) and (39), one shows that with probability greater than $1 - 3\delta$,

$$\mathbf{V}^{\widehat{\pi}_p^*} - \mathbf{V}^* \geq - \left(\frac{2}{1-\gamma} \xi + 12 \sqrt{\frac{2 \log \left(\frac{96|\mathcal{S}|^2|\mathcal{A}|^3}{\xi (1-\gamma)^4 \delta^2} \right)}{N(1-\gamma)^3}} \right) \mathbf{1}$$

By taking $\xi = \frac{(1-\gamma)\varepsilon}{3|\mathcal{S}|^5|\mathcal{A}|^5}$ and $N \geq \frac{c_0 \log(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^{\delta\varepsilon}})}{(1-\gamma)^{3\varepsilon^2}}$ for some constant $c_0 > 0$ large enough, we can ensure that $\mathbf{0} \geq \mathbf{V}^{\hat{\pi}_p^*} - \mathbf{V}^* \geq -\varepsilon \mathbf{1}$ as claimed. Regarding the Q-functions, the Bellman equation gives

$$\mathbf{Q}^{\hat{\pi}_p^*} - \mathbf{Q}^* = \mathbf{r} + \gamma \mathbf{P} \mathbf{V}^{\hat{\pi}_p^*} - (\mathbf{r} + \gamma \mathbf{P} \mathbf{V}^*) = \gamma \mathbf{P} (\mathbf{V}^{\hat{\pi}_p^*} - \mathbf{V}^*).$$

Consequently, one has

$$\mathbf{Q}^{\hat{\pi}_p^*} - \mathbf{Q}^* \geq -(\gamma \|\mathbf{P}\|_1 \|\mathbf{V}^{\hat{\pi}_p^*} - \mathbf{V}^*\|_\infty) \mathbf{1} \geq -\gamma \varepsilon \mathbf{1}.$$

Finally, we demonstrate that both the empirical QVI and PI w.r.t. $\widehat{\mathcal{M}}_p$ are guaranteed to find $\hat{\pi}_p^*$ in a few iterations. Suppose for the moment that we can obtain a policy π_k obeying

$$\|\widehat{\mathbf{Q}}_p^{\pi_k} - \widehat{\mathbf{Q}}_p^*\|_\infty < \frac{\xi \delta (1-\gamma)}{8|\mathcal{S}||\mathcal{A}|^2}. \quad (40)$$

Then for any $s \in \mathcal{S}$ and any action $a \neq \hat{\pi}_p^*(s)$ one has

$$\begin{aligned} & \widehat{\mathbf{Q}}_p^{\pi_k}(s, \hat{\pi}_p^*(s)) - \widehat{\mathbf{Q}}_p^{\pi_k}(s, a) \\ &= \widehat{\mathbf{Q}}_p^*(s, \hat{\pi}_p^*(s)) - \widehat{\mathbf{Q}}_p^*(s, a) + \left(\widehat{\mathbf{Q}}_p^{\pi_k}(s, \hat{\pi}_p^*(s)) - \widehat{\mathbf{Q}}_p^*(s, \hat{\pi}_p^*(s)) \right) - \left(\widehat{\mathbf{Q}}_p^{\pi_k}(s, a) - \widehat{\mathbf{Q}}_p^*(s, a) \right) \\ &\geq \widehat{\mathbf{Q}}_p^*(s, \hat{\pi}_p^*(s)) - \widehat{\mathbf{Q}}_p^*(s, a) - 2\|\widehat{\mathbf{Q}}_p^{\pi_k} - \widehat{\mathbf{Q}}_p^*\|_\infty \\ &> \frac{\xi \delta (1-\gamma)}{4|\mathcal{S}||\mathcal{A}|^2} - 2 \cdot \frac{\xi \delta (1-\gamma)}{8|\mathcal{S}||\mathcal{A}|^2} = 0, \end{aligned}$$

where the last line results from (37) and (40). In other words, we can perfectly recover the policy $\hat{\pi}_p^*$ from the estimate $\widehat{\mathbf{Q}}_p^{\pi_k}$, provided that (40) is satisfied. In addition, it has been shown that (Azar et al., 2013, Lemma 2) the greedy policy induced by k -th iteration of both algorithms — denoted by π_k — satisfies $\|\widehat{\mathbf{Q}}_p^{\pi_k} - \widehat{\mathbf{Q}}_p^*\|_\infty \leq \frac{2\gamma^{k+1}}{(1-\gamma)^2}$. Taking $\xi = \frac{c_1(1-\gamma)\varepsilon}{|\mathcal{S}|^5|\mathcal{A}|^5}$ and $k = \frac{c_2}{1-\gamma} \log(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)\varepsilon\delta})$ for some constant $c_2 > 0$ large enough, one guarantees that π_k satisfies (40), which in turn ensures perfect recovery of $\hat{\pi}_p^*$.

6 Discussion

This paper demonstrates that a *perturbed* model-based planning algorithm achieves the minimax sample complexity in a generative model, as soon as the sample size exceeds the order of $\frac{|\mathcal{S}||\mathcal{A}|}{1-\gamma}$ (modulo some log factor). Compared to prior literature, our result considerably broadens the sample size range, allowing us to pin down a complete trade-off curve between sample complexities and statistical accuracy.

The present work opens up several directions for future investigation, which we discuss in passing below.

- *Is perturbation necessary?* The planning algorithm analyzed here is applied to a perturbed variant of the original empirical MDP. This, however, gives rise to a natural question regarding the necessity of perturbation: can we achieve optimal planning performance directly using the original empirical MDP without perturbation? While we conjecture that the answer is affirmative, settling this conjecture requires developing new analysis techniques beyond the scope of this paper.
- *Improved analysis for model-free algorithms.* As mentioned previously, sample complexity barriers exist in prior theory for model-free approaches (e.g. Sidford et al. (2018a); Wainwright (2019b)). Our analysis might shed light on how to break such barriers for model-free approaches.
- *Episodic finite-horizon MDPs.* This work concentrates on infinite-horizon discounted MDPs, and we expect that the insights analysis can be extended to study the finite-horizon setting as well, where the samples are drawn in an episodic fashion (Dann and Brunskill (2015); Jiang and Agarwal (2018); Wang et al. (2020)). We leave this extension to future work.
- *Beyond the tabular setting.* This current paper focuses on the tabular setting with finite state and action spaces. While we improve the sample size range, the sample complexities might still be prohibitive when $|\mathcal{S}|$ and $|\mathcal{A}|$ are enormous. Therefore, it is desirable to further investigate settings where function approximation is employed to improve the efficiency (e.g. Duan and Wang (2020); Jin et al. (2019)).

Acknowledgements

Y. Wei is supported in part by the grants NSF CCF-2007911 and DMS-2015447. Y. Chi is supported in part by the grants ONR N00014-18-1-2142 and N00014-19-1-2404, ARO W911NF-18-1-0303, and NSF CCF-1806154 and CCF-2007911. Y. Chen is supported in part by the grants AFOSR YIP award FA9550-19-1-0030, ONR N00014-19-1-2120, ARO YIP award W911NF-20-1-0097, ARO W911NF-18-1-0303, NSF CCF-1907661, DMS-2014279 and IIS-1900140, and the Princeton SEAS Innovation Award. We thank Qiwen Cui for pointing out an issue in Appendix B.5 in an early version of this paper, and thank Shicong Cen, Chen Cheng and Cong Ma for numerous discussions about reinforcement learning.

A Preliminary facts

We begin by recording a few elementary facts about \mathbf{P}^π and \mathbf{P}_π (see definitions in (16)). These are standard results and we omit the proofs for brevity.

Lemma 7. *For any policy π , any probability transition matrix $\mathbf{P} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$ and any $0 < \gamma < 1$, one has*

- (a) $(\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1} = \sum_{i=0}^{\infty} (\gamma\mathbf{P}_\pi)^i$;
- (b) All entries of the matrix $(\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}$ are non-negative;
- (c) $\|(\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\|_1 \leq 1/(1 - \gamma)$;
- (d) $(1 - \gamma)(\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\mathbf{1} = \mathbf{1}$;
- (e) For any non-negative vectors $\mathbf{0} \leq \mathbf{r}_1 \leq \mathbf{r}_2$ of compatible dimension, one has $\mathbf{0} \leq (\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\mathbf{r}_1 \leq (\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\mathbf{r}_2$.

The above results continue to hold if \mathbf{P}_π is replaced by \mathbf{P}^π .

B Proofs of auxiliary lemmas

B.1 Proofs of Lemma 1 and Lemma 2

Auxiliary notation and preliminaries. Before proceeding, we define several $|\mathcal{S}|$ -dimensional auxiliary vectors $\mathbf{r}^{(i)}$, $\mathbf{V}^{(i)}$, $\widehat{\mathbf{V}}^{(i)}$ ($1 \leq i \leq m$) recursively as follows

$$\begin{aligned}
 \mathbf{r}^{(0)} &:= \mathbf{r}_\pi, & \mathbf{V}^{(0)} &:= (\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\mathbf{r}^{(0)}, & \widehat{\mathbf{V}}^{(0)} &:= (\mathbf{I} - \gamma\widehat{\mathbf{P}}_\pi)^{-1}\mathbf{r}^{(0)}, \\
 \mathbf{r}^{(1)} &:= \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(0)}]}, & \mathbf{V}^{(1)} &:= (\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\mathbf{r}^{(1)}, & \widehat{\mathbf{V}}^{(2)} &:= (\mathbf{I} - \gamma\widehat{\mathbf{P}}_\pi)^{-1}\mathbf{r}^{(1)}, \\
 &\vdots & &\vdots & &\vdots \\
 \mathbf{r}^{(m)} &:= \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(m-1)}]}, & \mathbf{V}^{(m)} &:= (\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\mathbf{r}^{(m)}, & \widehat{\mathbf{V}}^{(m)} &:= (\mathbf{I} - \gamma\widehat{\mathbf{P}}_\pi)^{-1}\mathbf{r}^{(m)},
 \end{aligned} \tag{41}$$

where m will be specified momentarily.

A crucial quantity that appears repeatedly in analyzing the above terms is $\|(\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})}\|_\infty$, whose importance was already made apparent in the work Azar et al. (2013). A widely used upper bound on this quantity, originally due to (Azar et al., 2013, Lemma 8), is given by

$$\left\| (\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})} \right\|_\infty \leq \frac{2 \log 2}{\gamma(1 - \gamma)^{1.5}} \|\mathbf{r}\|_\infty. \tag{42}$$

This bound turns out to be loose for our purpose, and we develop an improved bound as follows, whose proof is deferred to Appendix B.1.1.

Lemma 8. *Consider any policy π and any probability transition matrix $\mathbf{P} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$. Let \mathbf{V} be a vector obeying $\mathbf{V} = (\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\mathbf{r}_\pi$ for some $|\mathcal{S}|$ -dimensional vector $\mathbf{r}_\pi \geq \mathbf{0}$. For any $0 < \gamma < 1$, one has*

$$\left\| (\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})} \right\|_\infty \leq \frac{4}{\gamma\sqrt{1 - \gamma}} \|\mathbf{V}\|_\infty. \tag{43}$$

Remark 6. In comparison to the bound (42) derived in (Azar et al., 2013, Lemma 8), Lemma 8 offers an improved upper bound stated directly in terms of the properties of \mathbf{V} rather than those of \mathbf{r} .

As it turns out, Lemma 8 allows us to obtain an entrywise bound for $\mathbf{V}^{(l)}$ ($1 \leq l \leq m$). To begin with, the first term $\mathbf{V}^{(1)}$ satisfies

$$\|\mathbf{V}^{(1)}\|_\infty = \left\| (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(0)}]} \right\|_\infty \quad (44)$$

since $\mathbf{r}^{(1)} = \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(0)}]}$. Next, for any $l > 1$ one has

$$\begin{aligned} \|\mathbf{V}^{(l)}\|_\infty &= \|(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{r}^{(l)}\|_\infty = \left\| (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(l-1)}]} \right\|_\infty \\ &\leq \frac{4}{\gamma \sqrt{1-\gamma}} \|\mathbf{V}^{(l-1)}\|_\infty, \end{aligned}$$

where the second identity results from the definition of $\mathbf{r}^{(l)}$, and the last inequality comes from Lemma 8. As a consequence, applying this inequality recursively gives

$$\|\mathbf{V}^{(l)}\|_\infty \leq \left(\frac{4}{\gamma \sqrt{1-\gamma}} \right)^{l-1} \|\mathbf{V}^{(1)}\|_\infty. \quad (45)$$

Main proof. Equipped with the above facts, we are now in a position to prove the lemmas, for which we start with the more general one — Lemma 2. Consider any $0 \leq l \leq m$. We first observe that

$$\begin{aligned} \widehat{\mathbf{V}}^{(l)} - \mathbf{V}^{(l)} &= (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} \mathbf{r}^{(l)} - \mathbf{V}^{(l)} \\ &= (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} (\mathbf{I} - \gamma \mathbf{P}_\pi) \mathbf{V}^{(l)} - (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi) \mathbf{V}^{(l)} \\ &= \gamma (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^{(l)}, \end{aligned} \quad (46)$$

where the second line follows since, by definition, $(\mathbf{I} - \gamma \mathbf{P}_\pi) \mathbf{V}^{(l)} = \mathbf{r}^{(l)}$. Suppose that there exists some quantity $\beta_1 > 0$ such that the following condition

$$\left| (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^{(l)} \right| \leq \sqrt{\frac{\beta_1}{N}} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(l)}]} + \frac{\|\mathbf{V}^{(l)}\|_\infty \beta_1}{N} \mathbf{1} \quad (47)$$

holds uniformly for all $0 \leq l \leq m$. Then this combined with (46) gives

$$\begin{aligned} \|\widehat{\mathbf{V}}^{(l)} - \mathbf{V}^{(l)}\|_\infty &= \gamma \left\| (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^{(l)} \right\|_\infty \\ &\stackrel{(i)}{\leq} \gamma \left\| (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} \left| (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^{(l)} \right| \right\|_\infty \\ &\stackrel{(ii)}{\leq} \gamma \sqrt{\frac{\beta_1}{N}} \left\| (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(l)}]} \right\|_\infty + \frac{\gamma \|\mathbf{V}^{(l)}\|_\infty \beta_1}{N} \left\| (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} \mathbf{1} \right\|_\infty. \end{aligned}$$

Here, (i) follows since $(\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1}$ is a non-negative matrix, (ii) comes from (47) and the triangle inequality. Recalling the definition of $\mathbf{r}^{(l)}$ and $\widehat{\mathbf{V}}^{(l)}$ and invoking Lemma 7(d), we can further bound the above as

$$\begin{aligned} \|\widehat{\mathbf{V}}^{(l)} - \mathbf{V}^{(l)}\|_\infty &\leq \gamma \sqrt{\frac{\beta_1}{N}} \|\widehat{\mathbf{V}}^{(l+1)}\|_\infty + \frac{\gamma \|\mathbf{V}^{(l)}\|_\infty \beta_1}{(1-\gamma)N} \\ &\leq \gamma \sqrt{\frac{\beta_1}{N}} \|\widehat{\mathbf{V}}^{(l+1)} - \mathbf{V}^{(l+1)}\|_\infty + \gamma \sqrt{\frac{\beta_1}{N}} \|\mathbf{V}^{(l+1)}\|_\infty + \frac{\gamma \beta_1}{(1-\gamma)N} \|\mathbf{V}^{(l)}\|_\infty, \end{aligned} \quad (48)$$

where the last inequality is a consequence of the triangle inequality.

The above inequality (48) provides a recursive relation that in turn allows for an effective upper bound. Specifically, combining the inequalities (45) and (48) leads to

$$\|\widehat{\mathbf{V}}^{(0)} - \mathbf{V}^{(0)}\|_\infty \leq \gamma \sqrt{\frac{\beta_1}{N}} \|\widehat{\mathbf{V}}^{(1)} - \mathbf{V}^{(1)}\|_\infty + \gamma \sqrt{\frac{\beta_1}{N}} \|\mathbf{V}^{(1)}\|_\infty + \frac{\gamma \beta_1}{(1-\gamma)N} \|\mathbf{V}^{(0)}\|_\infty$$

$$=: b_1 \|\widehat{\mathbf{V}}^{(1)} - \mathbf{V}^{(1)}\|_\infty + b_1 \|\mathbf{V}^{(1)}\|_\infty + \frac{\gamma\beta_1}{(1-\gamma)N} \|\mathbf{V}^{(0)}\|_\infty,$$

and for $l \geq 1$,

$$\begin{aligned} \|\widehat{\mathbf{V}}^{(l)} - \mathbf{V}^{(l)}\|_\infty &\leq \gamma \sqrt{\frac{\beta_1}{N}} \|\widehat{\mathbf{V}}^{(l+1)} - \mathbf{V}^{(l+1)}\|_\infty + \left(4\sqrt{\frac{\beta_1}{(1-\gamma)N}} + \frac{\gamma\beta_1}{(1-\gamma)N}\right) \left(\frac{4}{\gamma\sqrt{1-\gamma}}\right)^{l-1} \|\mathbf{V}^{(1)}\|_\infty \\ &=: b_1 \|\widehat{\mathbf{V}}^{(l+1)} - \mathbf{V}^{(l+1)}\|_\infty + b_2 b_3^{l-1} \|\mathbf{V}^{(1)}\|_\infty. \end{aligned}$$

Here for notational simplicity, we introduce

$$b_1 := \gamma \sqrt{\frac{\beta_1}{N}}, \quad b_2 := 4\sqrt{\frac{\beta_1}{(1-\gamma)N}} + \frac{\gamma\beta_1}{(1-\gamma)N}, \quad b_3 := \frac{4}{\gamma\sqrt{1-\gamma}}.$$

Invoking the above recursive relation, we can arrange terms to reach

$$\|\widehat{\mathbf{V}}^{(0)} - \mathbf{V}^{(0)}\|_\infty \leq \underbrace{b_1^m \|\widehat{\mathbf{V}}^{(m)} - \mathbf{V}^{(m)}\|_\infty}_{=:\alpha_1} + \underbrace{\left(b_1 b_2 \sum_{l=0}^{m-2} (b_1 b_3)^l + b_1\right)}_{=:\alpha_2} \|\mathbf{V}^{(1)}\|_\infty + \frac{\gamma\beta_1}{(1-\gamma)N} \|\mathbf{V}^{(0)}\|_\infty. \quad (49)$$

Controlling the quantity α_2 . Now it suffices to control the two terms on the right-hand side of the inequality (49) separately, towards which we shall start with the quantity α_2 . Assuming that $N \geq 64\beta_1/(1-\gamma)$, one can easily verify $b_1 b_3 \leq 1/2$. The summation of the geometric sequence thus gives

$$\begin{aligned} \alpha_2 &:= b_1 b_2 \sum_{l=0}^{m-2} (b_1 b_3)^l + b_1 \leq \frac{b_1 b_2}{1 - b_1 b_3} + b_1 \leq 2b_1 b_2 + b_1 \\ &= \gamma \sqrt{\frac{\beta_1}{N}} \left\{ 1 + 8\sqrt{\frac{\beta_1}{(1-\gamma)N}} + \frac{2\gamma\beta_1}{(1-\gamma)N} \right\} \leq 3\gamma \sqrt{\frac{\beta_1}{N}}, \end{aligned} \quad (50)$$

where the last step holds with the assumption $N \geq \frac{64\beta_1}{1-\gamma}$.

Controlling the quantity α_1 . Next, we proceed to the quantity α_1 , which requires the control of $\|\widehat{\mathbf{V}}^{(m)} - \mathbf{V}^{(m)}\|_\infty$. In view of the identity (46), we obtain

$$\begin{aligned} \|\widehat{\mathbf{V}}^{(m)} - \mathbf{V}^{(m)}\|_\infty &= \gamma \left\| (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} (\mathbf{P}_\pi - \widehat{\mathbf{P}}_\pi) \mathbf{V}^{(m)} \right\|_\infty \\ &\leq \gamma \left\| (\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1} \left(\sqrt{\frac{\beta_1}{N}} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(m)}]} + \frac{\|\mathbf{V}^{(m)}\|_\infty \beta_1}{N} \mathbf{1} \right) \right\|_\infty, \end{aligned}$$

where the last inequality follows from the Bernstein-type condition (47) and the fact that $(\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1}$ has non-negative entries. By virtue of the simple relation $\sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(m)}]} \leq \|\mathbf{V}^{(m)}\|_\infty$ and the fact that $\|(\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1}\|_1 \leq \frac{1}{1-\gamma}$ (cf. Lemma 7(c)), it is further guaranteed that

$$\begin{aligned} \|\widehat{\mathbf{V}}^{(m)} - \mathbf{V}^{(m)}\|_\infty &\leq \gamma \|(\mathbf{I} - \gamma \widehat{\mathbf{P}}_\pi)^{-1}\|_1 \cdot \left\| \sqrt{\frac{\beta_1}{N}} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(m)}]} + \frac{\|\mathbf{V}^{(m)}\|_\infty \beta_1}{N} \mathbf{1} \right\|_\infty \\ &\leq \frac{\gamma}{1-\gamma} \left(\sqrt{\frac{\beta_1}{N}} + \frac{\beta_1}{N} \right) \|\mathbf{V}^{(m)}\|_\infty, \end{aligned}$$

which combined with the bound (45) yields

$$\|\widehat{\mathbf{V}}^{(m)} - \mathbf{V}^{(m)}\|_\infty \leq \frac{\gamma}{1-\gamma} \left(\sqrt{\frac{\beta_1}{N}} + \frac{\beta_1}{N} \right) \left(\frac{4}{\gamma\sqrt{1-\gamma}} \right)^{m-1} \|\mathbf{V}^{(1)}\|_\infty.$$

Putting the above bounds together yields

$$\begin{aligned}
\alpha_1 &:= b_1^m \|\widehat{\mathbf{V}}^{(m)} - \mathbf{V}^{(m)}\|_\infty \leq \left(\gamma \sqrt{\frac{\beta_1}{N}} \right)^m \frac{\gamma}{1-\gamma} \left(\sqrt{\frac{\beta_1}{N}} + \frac{\beta_1}{N} \right) \left(\frac{4}{\gamma \sqrt{1-\gamma}} \right)^{m-1} \|\mathbf{V}^{(1)}\|_\infty \\
&= \left(\sqrt{\frac{16\beta_1}{N(1-\gamma)}} \right)^{m-1} \left(\sqrt{\frac{\beta_1}{N}} + 1 \right) \frac{\gamma^2 \beta_1}{(1-\gamma)N} \|\mathbf{V}^{(1)}\|_\infty \\
&\leq \left(\frac{1}{e} \right)^{m-1} \frac{1.1\gamma^2 \beta_1}{(1-\gamma)N} \|\mathbf{V}^{(1)}\|_\infty, \tag{51}
\end{aligned}$$

where the last inequality holds provided that $N > \frac{16e^2}{1-\gamma} \beta_1$.

Putting all this together. Combining the inequalities (49), (50) and (51) gives

$$\|\widehat{\mathbf{V}}^{(0)} - \mathbf{V}^{(0)}\|_\infty \leq \frac{\gamma\beta_1}{(1-\gamma)N} \|\mathbf{V}^{(0)}\|_\infty + \left\{ 3\gamma \sqrt{\frac{\beta_1}{N}} + \left(\frac{1}{e} \right)^{m-1} \frac{1.1\gamma^2 \beta_1}{(1-\gamma)N} \right\} \|\mathbf{V}^{(1)}\|_\infty.$$

To finish up, set $m = \log(\frac{e}{1-\gamma})$ and assume that $N > \frac{16e^2}{1-\gamma} \beta_1$. Recognizing that $\mathbf{V}^{(0)} = \mathbf{V}^\pi$ and $\widehat{\mathbf{V}}^{(0)} = \widehat{\mathbf{V}}^\pi$, we arrive at

$$\begin{aligned}
\|\widehat{\mathbf{V}}^\pi - \mathbf{V}^\pi\|_\infty &= \|\widehat{\mathbf{V}}^{(0)} - \mathbf{V}^{(0)}\|_\infty \leq \frac{\gamma\beta_1}{(1-\gamma)N} \|\mathbf{V}^{(0)}\|_\infty + \left\{ 3\gamma \sqrt{\frac{\beta_1}{N}} + \left(\frac{1}{e} \right)^{m-1} \frac{1.1\gamma^2 \beta_1}{(1-\gamma)N} \right\} \|\mathbf{V}^{(1)}\|_\infty \\
&\leq \frac{\gamma\beta_1}{(1-\gamma)N} \|\mathbf{V}^\pi\|_\infty + 4\gamma \sqrt{\frac{\beta_1}{N}} \left\| (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^\pi]} \right\|_\infty, \tag{52}
\end{aligned}$$

provided that $N \geq \frac{16e^2 \beta_1}{1-\gamma}$.

Proof of Lemma 2. Invoking the inequality (42) to bound the second term of (52), we reach

$$\|\widehat{\mathbf{V}}^\pi - \mathbf{V}^\pi\|_\infty \leq \frac{\gamma\beta_1}{(1-\gamma)N} \|\mathbf{V}^\pi\|_\infty + 8 \log 2 \sqrt{\frac{\beta_1}{N(1-\gamma)^3}} \|\mathbf{r}\|_\infty \leq 6 \sqrt{\frac{\beta_1}{N(1-\gamma)^3}} \|\mathbf{r}\|_\infty, \tag{53}$$

where the last inequality uses the elementary fact that $\|\mathbf{V}^\pi\|_\infty \leq \frac{1}{1-\gamma} \|\mathbf{r}\|_\infty$ and the assumption that $N > \frac{16e^2}{1-\gamma} \beta_1$. We complete the proof of Lemma 2.

Proof of Lemma 1. Finally, to establish Lemma 1, we observe that: for any fixed policy π , the vector $\mathbf{V}^{(l)}$ is independent of $\widehat{\mathbf{P}}_\pi$. The Bernstein inequality (e.g. (Agarwal et al., 2019, Lemma 6)) then reveals that with probability at least $1 - \delta$,

$$\left| (\widehat{\mathbf{P}}_\pi - \mathbf{P}_\pi) \mathbf{V}^{(l)} \right| \leq \sqrt{\frac{2 \log \left(\frac{4m|S|}{\delta} \right)}{N}} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^{(l)}]} + \frac{\|\mathbf{V}^{(l)}\|_\infty \log \left(\frac{4m|S|}{\delta} \right)}{N} \mathbf{1} \tag{54}$$

holds uniformly for all $0 \leq l \leq m$. This means that we can take $\beta_1 := 2 \log \left(\frac{4m|S|}{\delta} \right)$ with $m = \log(\frac{e}{1-\gamma})$ for this case. Combining this with the inequality (52), we derive the advertised instance-dependent bound

$$\|\widehat{\mathbf{V}}^\pi - \mathbf{V}^\pi\|_\infty \leq \frac{2\gamma \log \left(\frac{4|S| \log \frac{e}{1-\gamma}}{\delta} \right)}{(1-\gamma)N} \|\mathbf{V}^\pi\|_\infty + 4\gamma \sqrt{\frac{2 \log \left(\frac{4|S| \log \frac{e}{1-\gamma}}{\delta} \right)}{(1-\gamma)N}} \left\| (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \sqrt{\text{Var}_{\mathbf{P}_\pi}[\mathbf{V}^\pi]} \right\|_\infty.$$

Further, this taken collectively with (42) and the crude bound $\|\mathbf{V}^\pi\|_\infty \leq \frac{1}{1-\gamma} \|\mathbf{r}\|_\infty$ gives

$$\|\widehat{\mathbf{V}}^\pi - \mathbf{V}^\pi\|_\infty \leq \frac{2\gamma \log \left(\frac{4|S| \log \frac{e}{1-\gamma}}{\delta} \right)}{(1-\gamma)^2 N} \|\mathbf{r}\|_\infty + 8 \log 2 \sqrt{\frac{2 \log \left(\frac{4|S| \log \frac{e}{1-\gamma}}{\delta} \right)}{(1-\gamma)^3 N}} \|\mathbf{r}\|_\infty$$

$$\leq 6\sqrt{\frac{2\log\left(\frac{4|\mathcal{S}|\log\frac{e}{1-\gamma}}{\delta}\right)}{(1-\gamma)^3 N}}\|\mathbf{r}\|_\infty,$$

with the proviso that $N \geq \frac{32e^2}{1-\gamma}\log\left(\frac{4|\mathcal{S}|\log\frac{e}{1-\gamma}}{\delta}\right)$.

B.1.1 Proof of Lemma 8

To begin with, we make the observation that

$$\begin{aligned}\text{Var}_{\mathbf{P}_\pi}(\mathbf{V}) &= \mathbf{P}_\pi(\mathbf{V} \circ \mathbf{V}) - (\mathbf{P}_\pi \mathbf{V}) \circ (\mathbf{P}_\pi \mathbf{V}) \\ &= \mathbf{P}_\pi(\mathbf{V} \circ \mathbf{V}) - \frac{1}{\gamma^2}(\mathbf{V} - \mathbf{r}_\pi) \circ (\mathbf{V} - \mathbf{r}_\pi)\end{aligned}\tag{55}$$

$$\begin{aligned}&= \mathbf{P}_\pi(\mathbf{V} \circ \mathbf{V}) - \frac{1}{\gamma^2}\mathbf{V} \circ \mathbf{V} + \frac{2}{\gamma^2}\mathbf{V} \circ \mathbf{r}_\pi - \frac{1}{\gamma^2}\mathbf{r}_\pi \circ \mathbf{r}_\pi \\ &\leq \frac{1}{\gamma^2}(\gamma^2\mathbf{P}_\pi - \mathbf{I})(\mathbf{V} \circ \mathbf{V}) + \frac{2}{\gamma^2}\mathbf{V} \circ \mathbf{r}_\pi,\end{aligned}\tag{56}$$

where the identity (55) makes use of the relation $\mathbf{V} = \mathbf{r}_\pi + \gamma\mathbf{P}_\pi\mathbf{V}$. In addition, one can deduce that

$$\begin{aligned}\left\|(\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})}\right\|_\infty &= \frac{1}{1-\gamma}\left\|(1-\gamma)(\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})}\right\|_\infty \\ &\stackrel{(i)}{\leq} \frac{1}{\sqrt{1-\gamma}}\left\|\sqrt{(\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})}\right\|_\infty \\ &\stackrel{(ii)}{\leq} \frac{1}{\sqrt{1-\gamma}}\sqrt{\|2(\mathbf{I} - \gamma^2\mathbf{P}_\pi)^{-1}\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})\|_\infty}.\end{aligned}$$

Here, (i) comes from Jensen's inequality (so that $\mathbb{E}[\sqrt{v}] \leq \sqrt{\mathbb{E}[v]}$) recognizing that each row of $(1-\gamma)(\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}$ is a probability distribution, and Lemma 7(d), (ii) is an elementary fact established in (Agarwal et al., 2019, Lemma 4). Combining Lemma 7(e) and the inequality (56) further yields

$$\begin{aligned}\left\|(\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})}\right\|_\infty &\leq \frac{1}{\sqrt{1-\gamma}}\sqrt{\left\|2(\mathbf{I} - \gamma^2\mathbf{P}_\pi)^{-1}\left(\frac{1}{\gamma^2}(\gamma^2\mathbf{P}_\pi - \mathbf{I})(\mathbf{V} \circ \mathbf{V}) + \frac{2}{\gamma^2}\mathbf{V} \circ \mathbf{r}_\pi\right)\right\|_\infty} \\ &\leq \frac{1}{\gamma\sqrt{1-\gamma}}\sqrt{2\|\mathbf{V} \circ \mathbf{V}\|_\infty} + \frac{2}{\gamma\sqrt{1-\gamma}}\sqrt{\|(\mathbf{I} - \gamma^2\mathbf{P}_\pi)^{-1}(\mathbf{V} \circ \mathbf{r}_\pi)\|_\infty}.\end{aligned}\tag{57}$$

where the last step arises from the triangle inequality and $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$. This leaves us with two terms to deal with.

Regarding the first term of (57), we observe that $\|(\mathbf{V} \circ \mathbf{V})\|_\infty = \|\mathbf{V}\|_\infty^2$. When it comes to the second term of (57), it is seen that

$$\begin{aligned}\|(\mathbf{I} - \gamma^2\mathbf{P}_\pi)^{-1}(\mathbf{V} \circ \mathbf{r}_\pi)\|_\infty &\leq \|(\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}(\mathbf{V} \circ \mathbf{r}_\pi)\|_\infty \\ &\leq \|\mathbf{V}\|_\infty \|(\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\mathbf{r}_\pi\|_\infty \\ &= \|\mathbf{V}\|_\infty^2.\end{aligned}$$

Here, the first inequality holds true since $(\mathbf{I} - \gamma^2\mathbf{P}_\pi)^{-1} = \sum_{i=0}^{\infty}\gamma^{2i}\mathbf{P}_\pi^i \leq \sum_{i=0}^{\infty}\gamma^i\mathbf{P}_\pi^i = (\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}$, while the second line holds since \mathbf{V} , \mathbf{r} and $(\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}$ are all non-negative. Substitution into (57) thus yields

$$\left\|(\mathbf{I} - \gamma\mathbf{P}_\pi)^{-1}\sqrt{\text{Var}_{\mathbf{P}_\pi}(\mathbf{V})}\right\|_\infty \leq \frac{\sqrt{2}}{\gamma\sqrt{1-\gamma}}\|\mathbf{V}\|_\infty + \frac{2}{\gamma\sqrt{1-\gamma}}\|\mathbf{V}\|_\infty \leq \frac{4}{\gamma\sqrt{1-\gamma}}\|\mathbf{V}\|_\infty$$

as claimed.

B.2 Proof of Lemma 3

To establish Lemma 3, it suffices to check that \mathbf{V}^* and \mathbf{Q}^* satisfy the Bellman optimality equations underlying \mathcal{M}_{s,a,u^*} . Towards this end, we study the absorbing state-action pair (s, a) and other pairs separately. For notational simplicity, we shall let \mathbf{P}^{abs} and $r^{\text{abs}}(\cdot, \cdot)$ denote respectively the probability transition matrix and the reward function associated with \mathcal{M}_{s,a,u^*} .

First, we observe that, by construction,

$$r^{\text{abs}}(s, a) + \gamma(\mathbf{P}^{\text{abs}}\mathbf{V}^*)_{s,a} = u^* + \gamma V^*(s).$$

Recall that \mathbf{V}^* satisfies the Bellman optimality equation w.r.t. the original MDP, namely, $Q^*(s, a) = r(s, a) + \gamma(\mathbf{P}\mathbf{V}^*)_{s,a}$. This together with our choice of u^* gives

$$u^* + \gamma V^*(s) = Q^*(s, a) - \gamma V^*(s) + \gamma V^*(s) = Q^*(s, a).$$

Putting the above identities together, we arrive at

$$Q^*(s, a) = r^{\text{abs}}(s, a) + \gamma(\mathbf{P}^{\text{abs}}\mathbf{V}^*)_{s,a}. \quad (58)$$

Next, consider any state-action pair $(s', a') \neq (s, a)$. Recalling again the properties of \mathcal{M}_{s,a,u^*} , we reach

$$r^{\text{abs}}(s', a') + \gamma(\mathbf{P}^{\text{abs}}\mathbf{V}^*)_{s',a'} = r(s', a') + \gamma(\mathbf{P}\mathbf{V}^*)_{s',a'} = Q^*(s', a'). \quad (59)$$

Here the last identity is due to the Bellman equation w.r.t. the original MDP. Combining (58) and (59) implies that \mathbf{V}^* and \mathbf{Q}^* satisfy Bellman's optimality equations in \mathcal{M}_{s,a,u^*} , thus concluding the proof.

B.3 Proof of Lemma 4

Our first observation is that $\widehat{\mathbf{Q}}_{s,a,u}^*$ satisfies Lipschitz continuity w.r.t. u in the sense that

$$\|\widehat{\mathbf{Q}}_{s,a,u}^* - \widehat{\mathbf{Q}}_{s,a,u'}^*\|_\infty \leq \frac{1}{1-\gamma}|u - u'|. \quad (60)$$

The proof of this relation is identical to that of (Agarwal et al., 2019, Lemma 8); we omit here for brevity. In view of Lemma 3, if we set $\widehat{u}^* := r(s, a) + \gamma(\widehat{\mathbf{P}}\widehat{\mathbf{V}}^*)_{s,a} - \gamma\widehat{V}^*(s)$, then one has

$$\widehat{\mathbf{Q}}_{s,a,\widehat{u}^*}^* = \widehat{\mathbf{Q}}^* \quad \text{and} \quad \widehat{\mathbf{V}}_{s,a,\widehat{u}^*}^* = \widehat{\mathbf{V}}^*.$$

In addition, there exists a point u_0 in the epsilon-net $\mathcal{N}_{(1-\gamma)\omega/4}$ such that $|\widehat{u}^* - u_0| \leq (1-\gamma)\omega/4$, which combined with the Lipschitz continuity property (60) gives

$$\|\widehat{\mathbf{Q}}^* - \widehat{\mathbf{Q}}_{s,a,u_0}^*\|_\infty = \|\widehat{\mathbf{Q}}_{s,a,\widehat{u}^*}^* - \widehat{\mathbf{Q}}_{s,a,u_0}^*\|_\infty \leq \frac{1}{1-\gamma}|\widehat{u}^* - u_0| \leq \frac{\omega}{4}. \quad (61)$$

Additionally, for any $s' \in \mathcal{S}$ and any $a_1, a_2 \in \mathcal{A}$ with $a_1 \neq a_2$, we have the following decomposition

$$\begin{aligned} & \widehat{\mathbf{Q}}_{s,a,u_0}^*(s', a_1) - \widehat{\mathbf{Q}}_{s,a,u_0}^*(s', a_2) \\ &= \widehat{\mathbf{Q}}^*(s', a_1) - \widehat{\mathbf{Q}}^*(s', a_2) + \widehat{\mathbf{Q}}_{s,a,u_0}^*(s', a_1) - \widehat{\mathbf{Q}}^*(s', a_1) - \left(\widehat{\mathbf{Q}}_{s,a,u_0}^*(s', a_2) - \widehat{\mathbf{Q}}^*(s', a_2) \right) \\ &\geq \widehat{\mathbf{Q}}^*(s', a_1) - \widehat{\mathbf{Q}}^*(s', a_2) - 2\|\widehat{\mathbf{Q}}^* - \widehat{\mathbf{Q}}_{s,a,u_0}^*\|_\infty \\ &\geq \widehat{\mathbf{Q}}^*(s', a_1) - \widehat{\mathbf{Q}}^*(s', a_2) - \frac{\omega}{2}, \end{aligned} \quad (62)$$

where the last inequality invokes the inequality (61). Moreover, our separation condition defined in (31) requires that: for any $s' \in \mathcal{S}$ and any $a_2 \neq \widehat{\pi}^*(s')$, one has $\widehat{\mathbf{Q}}^*(s', \widehat{\pi}^*(s')) - \widehat{\mathbf{Q}}^*(s', a_2) \geq \omega$, which together with (62) reveals that

$$\widehat{\mathbf{Q}}_{s,a,u_0}^*(s', \widehat{\pi}^*(s')) - \widehat{\mathbf{Q}}_{s,a,u_0}^*(s', a_2) \geq \widehat{\mathbf{Q}}^*(s', \widehat{\pi}^*(s')) - \widehat{\mathbf{Q}}^*(s', a_2) - \frac{\omega}{2} \geq \frac{\omega}{2}. \quad (63)$$

Given that $\widehat{\pi}_{s,a,u_0}^*(s') := \arg \max_{a'} \widehat{\mathbf{Q}}_{s,a,u_0}^*(s', a')$, it is seen from (63) that

$$\widehat{\pi}_{s,a,u_0}^*(s') = \widehat{\pi}^*(s'),$$

which holds true for all $s' \in \mathcal{S}$. This concludes the proof.

B.4 Proof of Lemma 5

We start by bounding $\|\widehat{\mathbf{V}}^{\widehat{\pi}^*} - \mathbf{V}^{\widehat{\pi}^*}\|_\infty$. Recall the definition of the series $\{\mathbf{V}^{(l)}\}$ in (41). Throughout this proof, we shall write $\mathbf{V}_\pi^{(l)}$ instead in order to make apparent the dependency on the policy π .

For each state-action pair (s, a) , let us construct the epsilon-net $\mathcal{N}_{(1-\gamma)\omega/4}$ as in the expression (32). For every $u \in \mathcal{N}_{(1-\gamma)\omega/4}$, recall that $\widehat{\pi}_{s,a,u}^*$ is defined as the optimal policy with respect to the (s, a) -absorbing MDP $\widehat{\mathcal{M}}_{s,a,u}$. By construction, the set of policies $\widehat{\pi}_{s,a,u}^*$ ($u \in \mathcal{N}_{(1-\gamma)\omega/4}$) is independent of $\widehat{\mathbf{P}}_{s,a}$. The Bernstein inequality (e.g. (Agarwal et al., 2019, Lemma 6)) taken together with the union bound thus guarantees that with probability at least $1 - \delta$,

$$\left| (\widehat{\mathbf{P}} - \mathbf{P})_{s,a} \mathbf{V}_{\widehat{\pi}_{s,a,u}^*}^{(l)} \right| \leq \sqrt{\frac{\beta_1}{N}} \sqrt{(\text{Var}_{\mathbf{P}}[\mathbf{V}_{\widehat{\pi}_{s,a,u}^*}^{(l)}])_{s,a}} + \frac{\|\mathbf{V}_{\widehat{\pi}_{s,a,u}^*}^{(l)}\|_\infty \beta_1}{N} \quad (64)$$

holds uniformly over all $0 \leq l \leq \log \frac{e}{1-\gamma}$, $u \in \mathcal{N}_{(1-\gamma)\omega/4}$, $(s, a) \in \mathcal{S} \times \mathcal{A}$. Here, β_1 is given by

$$\beta_1 := 2 \log \left(\frac{4 \log \left(\frac{e}{1-\gamma} \right) |\mathcal{N}_{(1-\gamma)\omega/4}| |\mathcal{S}| |\mathcal{A}|}{\delta} \right) \leq 2 \log \left(\frac{32}{(1-\gamma)^2 \omega \delta} |\mathcal{S}| |\mathcal{A}| \log \left(\frac{e}{1-\gamma} \right) \right),$$

where we have used the fact $|\mathcal{N}_{(1-\gamma)\omega/4}| \leq \frac{8}{(1-\gamma)^2 \omega}$. In addition, for any $0 < \omega < 1$, Lemma 4 guarantees that for each state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$, there exists a point $u_0 \in \mathcal{N}_{(1-\gamma)\omega/4}$ such that $\widehat{\pi}^* = \widehat{\pi}_{s,a,u_0}^*$. Invoking this important fact, we obtain

$$\begin{aligned} \left| (\widehat{\mathbf{P}} - \mathbf{P})_{s,a} \mathbf{V}_{\widehat{\pi}^*}^{(l)} \right| &= \left| (\widehat{\mathbf{P}} - \mathbf{P})_{s,a} \mathbf{V}_{\widehat{\pi}_{s,a,u_0}^*}^{(l)} \right| \\ &\leq \sqrt{\frac{\beta_1}{N}} \sqrt{(\text{Var}_{\mathbf{P}}[\mathbf{V}_{\widehat{\pi}_{s,a,u_0}^*}^{(l)}])_{s,a}} + \frac{\|\mathbf{V}_{\widehat{\pi}_{s,a,u_0}^*}^{(l)}\|_\infty \beta_1}{N} \\ &= \sqrt{\frac{\beta_1}{N}} \sqrt{(\text{Var}_{\mathbf{P}}[\mathbf{V}_{\widehat{\pi}^*}^{(l)}])_{s,a}} + \frac{\|\mathbf{V}_{\widehat{\pi}^*}^{(l)}\|_\infty \beta_1}{N}. \end{aligned}$$

The above inequality further allows us to deduce that, with probability $1 - \delta$,

$$\begin{aligned} \left| (\widehat{\mathbf{P}}_{\widehat{\pi}^*} - \mathbf{P}_{\widehat{\pi}^*}) \mathbf{V}_{\widehat{\pi}^*}^{(l)} \right| &= \left| \Pi^{\widehat{\pi}^*} (\widehat{\mathbf{P}} - \mathbf{P}) \mathbf{V}_{\widehat{\pi}^*}^{(l)} \right| \\ &\leq \sqrt{\frac{\beta_1}{N}} \sqrt{\Pi^{\widehat{\pi}^*} \text{Var}_{\mathbf{P}}[\mathbf{V}_{\widehat{\pi}^*}^{(l)}]} + \frac{\|\mathbf{V}_{\widehat{\pi}^*}^{(l)}\|_\infty \beta_1}{N} \mathbf{1} \\ &= \sqrt{\frac{\beta_1}{N}} \sqrt{\text{Var}_{\mathbf{P}_{\widehat{\pi}^*}}[\mathbf{V}_{\widehat{\pi}^*}^{(l)}]} + \frac{\|\mathbf{V}_{\widehat{\pi}^*}^{(l)}\|_\infty \beta_1}{N} \mathbf{1}. \end{aligned}$$

The above derivation validates the assumption required for Lemma 2. As a result, if $N > \frac{16e^2}{1-\gamma} \beta_1$ and $0 < \omega < 1$, then Lemma 2 leads to the advertised bound

$$\begin{aligned} \|\widehat{\mathbf{V}}^{\widehat{\pi}^*} - \mathbf{V}^{\widehat{\pi}^*}\|_\infty &\leq \frac{6}{1-\gamma} \sqrt{\frac{\beta_1}{N(1-\gamma)}} \leq \frac{6}{1-\gamma} \sqrt{\frac{2 \log \left(\frac{32}{(1-\gamma)^2 \omega \delta} |\mathcal{S}| |\mathcal{A}| \log \left(\frac{e}{1-\gamma} \right) \right)}{N(1-\gamma)}} \\ &\leq 6 \sqrt{\frac{2 \log \left(\frac{32 |\mathcal{S}| |\mathcal{A}|}{(1-\gamma)^3 \omega \delta} \right)}{N(1-\gamma)^3}}. \end{aligned} \quad (65)$$

Finally, we move on to the term $\mathbf{V}^* - \widehat{\mathbf{V}}^{\pi^*}$. Given that π^* is independent of $\widehat{\mathbf{P}}$, invoke Lemma 1 to reach

$$\|\widehat{\mathbf{V}}^{\pi^*} - \mathbf{V}^*\|_\infty = \|\widehat{\mathbf{V}}^{\pi^*} - \mathbf{V}^{\pi^*}\|_\infty \leq 6\sqrt{2} \sqrt{\frac{\log \left(\frac{4|\mathcal{S}|}{\delta} \right) + \log \log \left(\frac{e}{1-\gamma} \right)}{N(1-\gamma)^3}}$$

$$\leq 6\sqrt{\frac{2\log\left(\frac{32|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^2\omega\delta}\right)}{N(1-\gamma)^3}} \quad (66)$$

with probability at least $1 - \delta$. This together with (24) and (65) immediately establishes Lemma 5.

B.5 Proof of Lemma 6

The proofs for Q_p^* and \widehat{Q}_p^* are exactly the same; for the sake of conciseness, we shall only provide the proof for Q_p^* . Here we aim to prove a more general result than Lemma 6, namely, with probability at least $1 - \delta$,

$$\forall s \in \mathcal{S} \text{ and } a_1, a_2 \in \mathcal{A} \text{ with } a_1 \neq a_2 : \quad |Q_p^*(s, a_1) - Q_p^*(s, a_2)| > \frac{\xi\delta(1-\gamma)}{4|\mathcal{S}||\mathcal{A}|^2}.$$

Consider any state s and any actions $a_1 \neq a_2$. In what follows, we allow $r_p(s, a_1) = \tau$ to vary, while *freezing* the values of all other rewards $\{r_p(\tilde{s}, a) \mid (\tilde{s}, a) \neq (s, a_1)\}$. To streamline notation, we define

- $\mathbf{r}_\tau = [r_\tau(s, a)]_{(s,a) \in \mathcal{S} \times \mathcal{A}}$: the reward vector obeying

$$r_\tau(s, a_1) = \tau \quad \text{and} \quad r_\tau(\tilde{s}, a) = r_p(\tilde{s}, a) \quad \text{for all } (\tilde{s}, a) \neq (s, a_1);$$

- Q_τ^* : the optimal Q-function when the reward vector is \mathbf{r}_τ ;
- V_τ^* : the optimal value function when the reward vector is \mathbf{r}_τ ;
- π_τ^* : the optimal policy when the reward vector is \mathbf{r}_τ .

Additionally, we claim for the moment that there exists a phase transition boundary τ_{th} such that

$$\pi_\tau^*(s) \neq a_1, \quad \text{for all } \tau < \tau_{\text{th}}; \quad (67a)$$

$$\pi_\tau^*(s) = a_1, \quad \text{for all } \tau > \tau_{\text{th}}. \quad (67b)$$

The proof of this claim is deferred to the end of this section. To establish Lemma 6, the idea is to control the size of the set

$$\mathcal{I}_{0,\omega} := \{\tau \mid |Q_\tau^*(s, a_1) - Q_\tau^*(s, a_2)| < \omega\} \quad (68)$$

for some $\omega > 0$ to be specified shortly. As motivated by (67), we further break down this set into two parts $\mathcal{I}_{0,\omega} = \mathcal{I}_{1,\omega} \cup \mathcal{I}_{2,\omega}$, where

$$\mathcal{I}_{1,\omega} := \{\tau \mid \tau < \tau_{\text{th}}, |Q_\tau^*(s, a_1) - Q_\tau^*(s, a_2)| < \omega\}, \quad (69a)$$

$$\mathcal{I}_{2,\omega} := \{\tau \mid \tau \geq \tau_{\text{th}}, |Q_\tau^*(s, a_1) - Q_\tau^*(s, a_2)| < \omega\}. \quad (69b)$$

In what follows, we first control the size of each set separately, and then demonstrate that the probability of these events happening is very small.

Step 1. We begin with $\mathcal{I}_{1,\omega}$ associated with the range $\tau < \tau_{\text{th}}$. In this case, the value function V_τ^* does not vary with τ , since the reward $r_\tau(s, a_1) = \tau$ is never active when calculating V_τ^* (by virtue of (67a)). Thus, the Bellman equation allows us to write

$$Q_\tau^*(s, a_1) = \tau + B_1 \quad \text{and} \quad Q_\tau^*(s, a_2) = B_2$$

for some quantities B_1 and B_2 , where neither B_1 nor B_2 depends on the value of τ . Armed with this observation, we can easily show that: for any $\omega > 0$, the interval $\mathcal{I}_{1,\omega}$ (cf. (69a)) obeys

$$\mathcal{I}_{1,\omega} \subseteq \{\tau \mid |\tau + B_1 - B_2| < \omega\},$$

and hence has length (or Lebesgue measure) at most 2ω .

Step 2. We then move on to $\mathcal{I}_{2,\omega}$ associated with the range $\tau > \tau_{\text{th}}$ in which case $\pi_\tau^*(s) = a_1$. Towards this, we first make some useful observations.

- To begin with, given the relation $\mathbf{Q}_\tau^* = \mathbf{r}_\tau + \gamma \mathbf{P}\mathbf{V}_\tau^*$, it is easily seen that for any $\tau_2 > \tau_1 > \tau_{\text{th}}$,

$$\mathbf{0} \leq \mathbf{Q}_{\tau_2}^* - \mathbf{Q}_{\tau_1}^* \leq \mathbf{r}_{\tau_2} - \mathbf{r}_{\tau_1} + \gamma \|\mathbf{V}_{\tau_2}^* - \mathbf{V}_{\tau_1}^*\|_\infty. \quad (70)$$

In addition, for any state-action pair $(\tilde{s}, a) \neq (s, a_1)$, by construction we have $r_{\tau_2}(\tilde{s}, a) - r_{\tau_1}(\tilde{s}, a) = 0$, which together with (70) indicates that

$$\forall (\tilde{s}, a) \neq (s, a_1) : \quad 0 \leq Q_{\tau_2}^*(\tilde{s}, a) - Q_{\tau_1}^*(\tilde{s}, a) \leq \gamma \|\mathbf{V}_{\tau_2}^* - \mathbf{V}_{\tau_1}^*\|_\infty. \quad (71)$$

- Next, observe that for any $\tau_2 > \tau_1 > \tau_{\text{th}}$,

$$\forall s' \in \mathcal{S} : \quad 0 \leq V_{\tau_2}^*(s') - V_{\tau_1}^*(s') = \max_a Q_{\tau_2}^*(s', a) - \max_a Q_{\tau_1}^*(s', a) \leq \|\mathbf{Q}_{\tau_2}^* - \mathbf{Q}_{\tau_1}^*\|_\infty \quad (72)$$

and hence $\|\mathbf{Q}_{\tau_2}^* - \mathbf{Q}_{\tau_1}^*\|_\infty \geq \|\mathbf{V}_{\tau_2}^* - \mathbf{V}_{\tau_1}^*\|_\infty$. This combined with (71) and the fact $\gamma < 1$ implies that

$$Q_{\tau_2}^*(s, a_1) - Q_{\tau_1}^*(s, a_1) \geq \|\mathbf{V}_{\tau_2}^* - \mathbf{V}_{\tau_1}^*\|_\infty, \quad (73)$$

which together with the facts $V_{\tau_1}^*(s) = Q_{\tau_1}^*(s, a_1)$ and $V_{\tau_2}^*(s) = Q_{\tau_2}^*(s, a_1)$ (by virtue of (67b)) yields

$$\begin{aligned} V_{\tau_2}^*(s) - V_{\tau_1}^*(s) &= Q_{\tau_2}^*(s, a_1) - Q_{\tau_1}^*(s, a_1) \geq \|\mathbf{V}_{\tau_2}^* - \mathbf{V}_{\tau_1}^*\|_\infty \\ \implies Q_{\tau_2}^*(s, a_1) - Q_{\tau_1}^*(s, a_1) &= \|\mathbf{V}_{\tau_2}^* - \mathbf{V}_{\tau_1}^*\|_\infty. \end{aligned}$$

Invoke the Bellman equation to further derive

$$\begin{aligned} Q_{\tau_2}^*(s, a_1) - Q_{\tau_1}^*(s, a_1) &= \|\mathbf{V}_{\tau_2}^* - \mathbf{V}_{\tau_1}^*\|_\infty = \|\mathbf{r}_{\tau_2} + \gamma \mathbf{P}\mathbf{V}_{\tau_2}^* - \mathbf{r}_{\tau_1} - \gamma \mathbf{P}\mathbf{V}_{\tau_1}^*\|_\infty \\ &\geq \|\mathbf{r}_{\tau_2} - \mathbf{r}_{\tau_1}\|_\infty = \tau_2 - \tau_1, \end{aligned} \quad (74)$$

where the last inequality holds since $\mathbf{r}_{\tau_2} + \gamma \mathbf{P}\mathbf{V}_{\tau_2}^* - \mathbf{r}_{\tau_1} - \gamma \mathbf{P}\mathbf{V}_{\tau_1}^* \geq \mathbf{r}_{\tau_2} - \mathbf{r}_{\tau_1} \geq \mathbf{0}$ (due to the monotonicity properties $\mathbf{r}_{\tau_2} \geq \mathbf{r}_{\tau_1}$ and $\mathbf{V}_{\tau_2}^* \geq \mathbf{V}_{\tau_1}^*$), and the last identity follows from the definition of \mathbf{r}_τ .

With the above two properties (71) and (74) in mind, we are ready to locate $\mathcal{I}_{2,\omega}$ by showing that

$$\mathcal{I}_{2,\omega} \subseteq \left[\tau_{\text{th}}, \tau_{\text{th}} + \frac{\omega}{1-\gamma} \right]. \quad (75)$$

Given that $Q_{\tau_{\text{th}}}^*(s, a_1) \geq Q_{\tau_{\text{th}}}^*(s, a_2)$ (in view of (67)), we have for any $\tau \geq \tau_{\text{th}}$ and any $a_2 \neq a_1$ that

$$\begin{aligned} Q_\tau^*(s, a_1) - Q_\tau^*(s, a_2) &\geq (Q_\tau^*(s, a_1) - Q_{\tau_{\text{th}}}^*(s, a_1)) - (Q_\tau^*(s, a_2) - Q_{\tau_{\text{th}}}^*(s, a_2)) \\ &\geq (1-\gamma)(Q_\tau^*(s, a_1) - Q_{\tau_{\text{th}}}^*(s, a_1)). \end{aligned} \quad (76)$$

Here, the last inequality holds since

$$Q_\tau^*(s, a_2) - Q_{\tau_{\text{th}}}^*(s, a_2) \stackrel{(i)}{\leq} \gamma \|\mathbf{V}_\tau^* - \mathbf{V}_{\tau_{\text{th}}}^*\|_\infty \stackrel{(ii)}{\leq} \gamma (Q_\tau^*(s, a_1) - Q_{\tau_{\text{th}}}^*(s, a_1)),$$

where (i) follows from (71) and (ii) is due to (73). As a result, for any $\tau > \tau_{\text{th}} + \frac{\omega}{1-\gamma}$, one can invoke (76) and (74) to see that

$$\forall a_2 \neq a_1 : \quad Q_\tau^*(s, a_1) - Q_\tau^*(s, a_2) \geq (1-\gamma)(Q_\tau^*(s, a_1) - Q_{\tau_{\text{th}}}^*(s, a_1)) \geq (1-\gamma)(\tau - \tau_{\text{th}}) > \omega,$$

which necessarily implies that such a τ does not lie within the interval $\mathcal{I}_{2,\omega}$ as defined in (69b). This establishes the claimed relation (75).

Step 3. Putting the results in the above two steps together, we see the set $\mathcal{I}_{0,\omega}$ (cf. (68)) has total length (or Lebesgue measure) at most $\frac{3\omega}{1-\gamma}$. Given that $r_p(s, a) = r(s, a) + \zeta(s, a)$ with $\zeta(s, a) \sim \text{Unif}(0, \xi)$, one has

$$\mathbb{P} \left\{ |Q_p^*(s, a_1) - Q_p^*(s, a_2)| < \omega \right\} \leq \mathbb{P} \{r(s, a) + \zeta(s, a) \in \mathcal{I}_{0,\omega}\} \leq \frac{3\omega}{\xi(1-\gamma)}.$$

By setting $\omega = \frac{\delta(1-\gamma)\xi}{3|\mathcal{S}||\mathcal{A}|^2}$, we arrive at

$$\mathbb{P} \left\{ |Q_p^*(s, a_1) - Q_p^*(s, a_2)| < \frac{\delta(1-\gamma)\xi}{3|\mathcal{S}||\mathcal{A}|^2} \right\} \leq \frac{\delta}{|\mathcal{S}||\mathcal{A}|^2}.$$

Finally, taking the union bound over all s, a_1, a_2 , we conclude that

$$\mathbb{P} \left\{ \exists s, a_1 \neq a_2 : |Q_p^*(s, a_1) - Q_p^*(s, a_2)| < \frac{\delta(1-\gamma)\xi}{3|\mathcal{S}||\mathcal{A}|^2} \right\} \leq \delta,$$

thus establishing Lemma 6 as long as the claim (67) is valid.

Proof of the claim (67). To establish the claim, it suffices to take

$$\tau_{\text{th}} = \sup \{u \mid \pi_\tau^*(s) \neq a_1 \text{ for all } \tau < u\}. \quad (77)$$

It thus suffices to verify (67b) for our choice (77). Towards this, suppose instead that there exist some $\tau_3 < \tau_{\text{th}} \leq \tau_2 < \tau_1$ such that

$$\pi_{\tau_3}^*(s) \neq a_1, \quad \pi_{\tau_2}^*(s) = a_1, \quad \text{and} \quad \pi_{\tau_1}^*(s) \neq a_1.$$

It is straightforward to see that $V_{\tau_1}^* = V_{\tau_3}^*$, since in both cases, the reward $r_\tau(s, a_1)$ does not enter the calculation of the optimal value function (while the rewards in other state-action pairs are identical in both cases). In view of the monotonicity of the value function w.r.t. the reward vector, we have

$$V_{\tau_1}^* = V_{\tau_2}^* = V_{\tau_3}^*.$$

However, this contradicts our assumption that a_1 is the optimal action for state s at τ_2 but not at τ_3 , since enlarging τ_2 to τ_3 otherwise will enlarge the optimal value function $V_{\tau_3}^*$. We have thus established (67). \square

References

- Agarwal, A., Kakade, S., and Yang, L. F. (2019). Model-based reinforcement learning with a generative model is minimax optimal. *arXiv preprint arXiv:1906.03804*.
- Azar, M. G., Munos, R., and Kappen, B. (2012). On the sample complexity of reinforcement learning with a generative model. *arXiv preprint arXiv:1206.6461*.
- Azar, M. G., Munos, R., and Kappen, H. J. (2013). Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, 91(3):325–349.
- Azar, M. G., Osband, I., and Munos, R. (2017). Minimax regret bounds for reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 263–272. JMLR.org.
- Beck, C. L. and Srikant, R. (2012). Error bounds for constant step-size Q-learning. *Systems & control letters*, 61(12):1203–1208.
- Bellman, R. (1952). On the theory of dynamic programming. *Proceedings of the National Academy of Sciences of the United States of America*, 38(8):716.
- Bertsekas, D. P. (2017). *Dynamic programming and optimal control (4th edition)*. Athena Scientific.

- Bhandari, J., Russo, D., and Singal, R. (2018). A finite time analysis of temporal difference learning with linear function approximation. In *Conference On Learning Theory*, pages 1691–1692.
- Bradtke, S. J. and Barto, A. G. (1996). Linear least-squares algorithms for temporal difference learning. *Machine learning*, 22(1-3):33–57.
- Cai, Q., Yang, Z., Lee, J. D., and Wang, Z. (2019). Neural temporal-difference learning converges to global optima. In *Advances in Neural Information Processing Systems*, pages 11312–11322.
- Chen, Y., Fan, J., Ma, C., and Wang, K. (2019). Spectral method and regularized mle are both optimal for top-k ranking. *Annals of statistics*, 47(4):2204.
- Chen, Z., Maguluri, S. T., Shakkottai, S., and Shanmugam, K. (2020). Finite-sample analysis of stochastic approximation using smooth convex envelopes. *arXiv preprint arXiv:2002.00874*.
- Dalal, G., Szörényi, B., Thoppe, G., and Mannor, S. (2018). Finite sample analyses for TD(0) with function approximation. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Dann, C. and Brunskill, E. (2015). Sample complexity of episodic fixed-horizon reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 2818–2826.
- Duan, Y. and Wang, M. (2020). Minimax-optimal off-policy evaluation with linear function approximation. *arXiv preprint arXiv:2002.09516*.
- El Karoui, N. (2015). On the impact of predictor geometry on the performance on high-dimensional ridge-regularized generalized robust regression estimators. *Probability Theory and Related Fields*, pages 1–81.
- Even-Dar, E. and Mansour, Y. (2003). Learning rates for Q-learning. *Journal of machine learning Research*, 5(Dec):1–25.
- Fan, J., Wang, Z., Xie, Y., and Yang, Z. (2019). A theoretical analysis of deep Q-learning. *arXiv preprint arXiv:1901.00137*.
- Gupta, H., Srikant, R., and Ying, L. (2019). Finite-time performance bounds and adaptive learning rate selection for two time-scale reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 4706–4715.
- Jaakkola, T., Jordan, M. I., and Singh, S. P. (1994). Convergence of stochastic iterative dynamic programming algorithms. In *Advances in neural information processing systems*, pages 703–710.
- Jiang, N. and Agarwal, A. (2018). Open problem: The dependence of sample complexity lower bounds on planning horizon. In *Conference On Learning Theory*, pages 3395–3398.
- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. (2018). Is Q-learning provably efficient? In *Advances in Neural Information Processing Systems*, pages 4863–4873.
- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. (2019). Provably efficient reinforcement learning with linear function approximation. *arXiv preprint arXiv:1907.05388*.
- Kakade, S. (2003). *On the sample complexity of reinforcement learning*. PhD thesis, University of London.
- Kaledin, M., Moulines, E., Naumov, A., Tadic, V., and Wai, H.-T. (2020). Finite time analysis of linear two-timescale stochastic approximation with Markovian noise. *arXiv preprint arXiv:2002.01268*.
- Kearns, M., Mansour, Y., and Ng, A. Y. (2002). A sparse sampling algorithm for near-optimal planning in large Markov decision processes. *Machine learning*, 49(2-3):193–208.
- Kearns, M. J. and Singh, S. P. (1999). Finite-sample convergence rates for Q-learning and indirect algorithms. In *Advances in neural information processing systems*, pages 996–1002.
- Khamaru, K., Pananjady, A., Ruan, F., Wainwright, M. J., and Jordan, M. I. (2020). Is temporal difference learning optimal? an instance-dependent analysis. *arXiv preprint arXiv:2003.07337*.

- Lakshminarayanan, C. and Szepesvari, C. (2018). Linear stochastic approximation: How far does constant step-size and iterate averaging go? In *International Conference on Artificial Intelligence and Statistics*, pages 1347–1355.
- Lattimore, T. and Hutter, M. (2012). PAC bounds for discounted MDPs. In *International Conference on Algorithmic Learning Theory*, pages 320–334. Springer.
- Li, G., Wei, Y., Chi, Y., Gu, Y., and Chen, Y. (2020). Sample complexity of asynchronous Q-learning: Sharper analysis and variance reduction. *arXiv preprint arXiv:2006.03041*.
- Ma, C., Wang, K., Chi, Y., and Chen, Y. (2020). Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion and blind deconvolution. *Foundations of Computational Mathematics*, 20(3):451–632.
- Mou, W., Li, C. J., Wainwright, M. J., Bartlett, P. L., and Jordan, M. I. (2020). On linear stochastic approximation: Fine-grained Polyak-Ruppert and non-asymptotic concentration. *arXiv preprint arXiv:2004.04719*.
- Pananjady, A. and Wainwright, M. J. (2019). Value function estimation in Markov reward processes: Instance-dependent ℓ_∞ -bounds for policy evaluation. *arXiv preprint arXiv:1909.08749*.
- Puterman, M. L. (2014). *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons.
- Qu, G. and Wierman, A. (2020). Finite-time analysis of asynchronous stochastic approximation and Q-learning. *arXiv preprint arXiv:2002.00260*.
- Shah, D. and Xie, Q. (2018). Q-learning with nearest neighbors. In *Advances in Neural Information Processing Systems*, pages 3111–3121.
- Sidford, A., Wang, M., Wu, X., Yang, L., and Ye, Y. (2018a). Near-optimal time and sample complexities for solving Markov decision processes with a generative model. In *Advances in Neural Information Processing Systems*, pages 5186–5196.
- Sidford, A., Wang, M., Wu, X., and Ye, Y. (2018b). Variance reduced value iteration and faster algorithms for solving Markov decision processes. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 770–787. SIAM.
- Srikant, R. and Ying, L. (2019). Finite-time error bounds for linear stochastic approximation and TD learning. In *Conference on Learning Theory*, pages 2803–2830.
- Strehl, A. L., Li, L., Wiewiora, E., Langford, J., and Littman, M. L. (2006). PAC model-free reinforcement learning. In *Proceedings of the 23rd international conference on Machine learning*, pages 881–888.
- Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.
- Szepesvári, C. (1998). The asymptotic convergence-rate of Q-learning. In *Advances in Neural Information Processing Systems*, pages 1064–1070.
- Szepesvári, C. (2010). Algorithms for reinforcement learning. *Synthesis lectures on artificial intelligence and machine learning*, 4(1):1–103.
- Tsitsiklis, J. and Van Roy, B. (1997). An analysis of temporal-difference learning with function approximation. *IEEE Transactions on Automatic Control*, 42(5):674–690.
- Tsitsiklis, J. N. (1994). Asynchronous stochastic approximation and Q-learning. *Machine learning*, 16(3):185–202.
- Tu, S. and Recht, B. (2018). The gap between model-based and model-free methods on the linear quadratic regulator: An asymptotic viewpoint. *arXiv preprint arXiv:1812.03565*.

- Vershynin, R. (2018). *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press.
- Wainwright, M. J. (2019a). Stochastic approximation with cone-contractive operators: Sharp ℓ_∞ -bounds for Q-learning. *arXiv preprint arXiv:1905.06265*.
- Wainwright, M. J. (2019b). Variance-reduced Q-learning is minimax optimal. *arXiv preprint arXiv:1906.04697*.
- Wang, M. (2019). Randomized linear programming solves the Markov decision problem in nearly linear (sometimes sublinear) time. *Mathematics of Operations Research*.
- Wang, R., Du, S. S., Yang, L. F., and Kakade, S. M. (2020). Is long horizon reinforcement learning more difficult than short horizon reinforcement learning? *arXiv preprint arXiv:2005.00527*.
- Xu, P. and Gu, Q. (2020). A finite-time analysis of Q-learning with neural network function approximation. *accepted to International Conference on Machine Learning*.
- Xu, T., Zou, S., and Liang, Y. (2019). Two time-scale off-policy TD learning: Non-asymptotic analysis over Markovian samples. In *Advances in Neural Information Processing Systems*, pages 10633–10643.
- Yang, L. and Wang, M. (2019). Sample-optimal parametric Q-learning using linearly additive features. In *International Conference on Machine Learning*, pages 6995–7004.