

---

# Information-Theoretic Considerations in Batch Reinforcement Learning

---

Jinglin Chen<sup>1</sup> Nan Jiang<sup>1</sup>

## Abstract

Value-function approximation methods that operate in batch mode have foundational importance to reinforcement learning (RL). Finite sample guarantees for these methods often crucially rely on two types of assumptions: (1) mild distribution shift, and (2) representation conditions that are stronger than realizability. However, the necessity (“why do we need them?”) and the naturalness (“when do they hold?”) of such assumptions have largely eluded the literature. In this paper, we revisit these assumptions and provide theoretical results towards answering the above questions, and make steps towards a deeper understanding of value-function approximation.

## 1. Introduction and Related Work

We are concerned with value-function approximation in batch-mode reinforcement learning, which is related to and sometimes known as Approximate Dynamic Programming (ADP; Bertsekas & Tsitsiklis, 1996). Such methods take sample transition data as input<sup>1</sup> and approximate the optimal value-function  $Q^*$  from a restricted class that encodes one’s prior knowledge and inductive biases. They provide an important foundation for RL’s empirical success today, as many popular deep RL algorithms find their prototypes in this literature. For example, when DQN (Mnih et al., 2015) is run on off-policy data, and the target network is updated slowly, it can be viewed as the stochastic approximation of its batch analog, Fitted Q-Iteration (Ernst et al., 2005), with a neural net as the function approximator (Riedmiller, 2005; Yang et al., 2019).

Given the importance of these methods, the question of

---

<sup>1</sup>University of Illinois at Urbana-Champaign, Urbana, Illinois, USA. Correspondence to: Nan Jiang <nanjiang@illinois.edu>.

*Proceedings of the 36<sup>th</sup> International Conference on Machine Learning*, Long Beach, California, PMLR 97, 2019. Copyright 2019 by the author(s).

<sup>1</sup>In this paper, we restrict ourselves to the so-called *one-path* setting and do not allow multiple samples from the same state (Sutton & Barto, 1998; Maillard et al., 2010), which is only feasible in certain simulated environments and allows algorithms to succeed with realizability as the only representation condition.

*when they work* is central to our understanding of RL. Existing works that analyze error propagation and finite sample behavior of ADP methods (Munos, 2003; Szepesvári & Munos, 2005; Antos et al., 2008; Munos & Szepesvári, 2008; Tosatto et al., 2017) have provided us with a decent understanding: To guarantee sample-efficient learning of near-optimal policies, we often need assumptions from the following two categories.

**Mild distribution shift** Many ADP methods can run completely off-policy and they do the best with whatever data available.<sup>2</sup> Therefore, it is necessary that the data have sufficient coverage over the state (and action) space.

**Representation condition** Since the ultimate goal is to find  $Q^*$ , we would expect that the function class we work with contains it (or at least a close approximation). While such realizability-type assumptions are sufficient for supervised learning, reinforcement learning faces the additional difficulties of delayed consequences and the lack of labels, and existing analyses often make stronger assumptions on the function class, such as (approximate) closedness under Bellman update (Szepesvári & Munos, 2005).

While the above assumptions make intuitive sense, and finite sample bounds have been proved when they hold, their necessity (“*can we prove similar results without making these assumptions?*”) and naturalness (“*do they actually hold in interesting problems?*”) have largely eluded the literature. In this paper, we revisit these assumptions and provide theoretical results towards answering the above questions. Below is a highlight of our results:

1. To prepare for later discussions, we provide an analysis of representative ADP algorithms (FQI and its variant) under a simplified and minimal setup (Section 3). As a side-product, our results improve upon prior analyses in the dependence of error rate on sample size.
2. We formally justify the necessity of mild distribution shift via an information-theoretic lower bound (Section 4.1). Our setup rules out trivial and uninteresting failure mode due to an adversarial choice of data: Even

---

<sup>2</sup>Even when they are on-policy or combined with a standard exploration module (e.g.,  $\epsilon$ -greedy), most often they fail in problems where exploration is difficult (e.g., combination lock; see Kakade, 2003) and rely on the benignness of data to succeed.

with the most favorable data distribution, polynomial sample complexity is not achievable if the MDP dynamics are not restricted.

3. We conjecture an information-theoretic lower bound against realizability alone as the representation condition (Conjecture 8, Section 5.1). While we are not able to prove the conjecture, important steps are made, as two very general proof styles are shown to be destined to fail, one of which is due to Sutton & Barto (2018) and has been used to prove a closely related result.
4. As another side-product, we prove that *model-based RL* can enjoy polynomial sample complexity with realizability alone (Corollary 6). If Conjecture 8 is true, we have a formal separation showing the gap between batch model-based vs value-based RL with function approximation (see the analog in the online exploration setting in Sun et al. (2019)).

Throughout the paper, we make novel connections to two subareas of RL: state abstractions (Whitt, 1978; Li et al., 2006) and PAC exploration under function approximation (Krishnamurthy et al., 2016; Jiang et al., 2017). In particular, we are able to utilize some of their results in our proofs (Sections 4.1 and 5.1), and find examples from these areas where the assumptions of interest hold (Sections 4.2 and 5.2). This suggests that the results in these other areas may be beneficial to the research in ADP, and we hope this work can inspire researchers from different subareas of RL to exchange ideas more often.

## 2. Preliminaries

### 2.1. Markov Decision Processes (MDPs)

Let  $M = (\mathcal{S}, \mathcal{A}, P, R, \gamma, \eta_1)$  be an MDP, where  $\mathcal{S}$  is the finite (but can be arbitrarily large) state space,  $\mathcal{A}$  is the finite action space,  $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  is the transition function ( $\Delta(\cdot)$  is the probability simplex),  $R : \mathcal{S} \times \mathcal{A} \rightarrow [0, R_{\max}]$  is the reward function,  $\gamma \in [0, 1)$  is the discount factor, and  $\eta_1$  is the initial distribution over states.

A (stochastic) policy  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$  prescribes a distribution over actions for each state. Fixing a start state  $s$ , the policy  $\pi$  induces a random trajectory  $s_1, a_1, r_1, s_2, a_2, r_2, \dots$ , where  $s_1 = s$ ,  $a_1 \sim \pi(s_1)$ ,  $r_1 = R(s_1, a_1)$ ,  $s_2 \sim P(s_1, a_1)$ ,  $a_2 \sim \pi(s_2)$ , etc. The goal is to find  $\pi$  that maximizes the expected return  $v^\pi := \mathbb{E}[\sum_{h=1}^{\infty} \gamma^{h-1} r_h | s_1 \sim \eta_1, \pi]$ . It will also be useful to define the value function  $V^\pi(s) := \mathbb{E}[\sum_{h=1}^{\infty} \gamma^{h-1} r_h | s_1 = s, \pi]$  and Q-value function  $Q^\pi(s, a) := \mathbb{E}[\sum_{h=1}^{\infty} \gamma^{h-1} r_h | s_1 = s, a_1 = a, a_{2:\infty} \sim \pi]$ , and these functions take values in  $[0, V_{\max}]$  with  $V_{\max} := R_{\max}/(1 - \gamma)$ .

There exists a deterministic policy<sup>3</sup>  $\pi^*$  that maximizes  $V^\pi(s)$  for all  $s \in \mathcal{S}$  simultaneously, and hence also maximizes  $v^\pi$  as  $v^\pi = \mathbb{E}_{s_1 \sim \eta_1}[V^\pi(s_1)]$ . Let  $V^*$  and  $Q^*$  be the shorthand for  $V^{\pi^*}$  and  $Q^{\pi^*}$  respectively. It is well known that  $\pi^*(s) = \pi_{Q^*}(s) := \arg \max_{a \in \mathcal{A}} Q^*(s, a)$ , and  $Q^*$  satisfies the *Bellman equation*  $Q^* = \mathcal{T}Q^*$ , where  $\mathcal{T} : \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \rightarrow \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$  is the *Bellman update operator*:  $\forall f \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ ,

$$(\mathcal{T}f)(s, a) := R(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)}[V_f(s')], \quad (1)$$

where  $V_f(s') := \max_{a' \in \mathcal{A}} f(s', a')$ .

**Additional notations** Let  $\eta_h^\pi$  be the marginal distribution of  $s_h$  under  $\pi$ , that is,  $\eta_h^\pi(s) := \Pr[s_h = s | s_1 \sim \eta_1, \pi]$ . For  $g : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ ,  $\nu \in \Delta(\mathcal{S} \times \mathcal{A})$ , and  $p \geq 1$ , define the shorthand  $\|g\|_{p, \nu} := (\mathbb{E}_{(s, a) \sim \nu}[|g(s, a)|^p])^{1/p}$ , which is a semi-norm. Furthermore, for any object that is a function of/distribution over  $\mathcal{S}$  (or  $\mathcal{S} \times \mathcal{A}$ ), we will treat it as a vector whenever convenient. We add a subscript to the value functions or Bellman update operators, e.g.,  $V_M^*$ , when it is necessary to clarify the MDP in which the object is defined.

### 2.2. Batch Value-Function Approximation

This paper is concerned with *batch-mode* RL with value-function approximation. As a typical setup, the agent does not have direct access to the MDP and instead is given the following inputs:

- A batch dataset  $D$  consisting of  $(s, a, r, s')$  tuples, where  $r = R(s, a)$  and  $s' \sim P(s, a)$ . For simplicity, we assume that  $(s, a)$  is generated i.i.d. from the *data distribution*  $\mu \in \Delta(\mathcal{S} \times \mathcal{A})$ .<sup>4</sup>
- A class of candidate value-functions,  $\mathcal{F} \subset (\mathcal{S} \times \mathcal{A} \rightarrow [0, V_{\max}])$ , which (approximately) captures  $Q^*$ ; such a property is often called *realizability*. We discuss additional assumptions on  $\mathcal{F}$  later. As a further simplification, we focus on finite but exponentially large  $\mathcal{F}$  and discuss how to handle infinite classes when appropriate.

The learning goal is to compute a near-optimal policy from the data, often via finding  $f \in \mathcal{F}$  that approximates  $Q^*$  and outputting  $\pi_f$ , the greedy policy w.r.t.  $f$ . A representative algorithm for this setting is Fitted Q-Iteration (FQI) (Ernst et al., 2005; Szepesvári, 2010).<sup>5</sup> The algorithm initializes

<sup>3</sup>A deterministic policy puts all the probability mass on a single action in each state. With a slight abuse of notation, we sometimes also treat the type of such policies as  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ .

<sup>4</sup>The agent may or may not have knowledge of  $\mu$ . Most existing algorithms are agnostic to such knowledge.

<sup>5</sup>Batch value-based algorithms can often be categorized into approximate value iteration (e.g., FQI) and approximate policy iteration (e.g., LSPI (Lagoudakis & Parr, 2003)). We focus on the former due to its simplicity and do not discuss the latter as its

$f_0 \in \mathcal{F}$  arbitrarily, and iteratively computes  $f_k$  as follows: in iteration  $k$ , the algorithm converts the dataset  $D$  into a regression dataset, with  $(s, a)$  being the input and  $r + \gamma V_{f_{k-1}}(s')$  as the output. It then minimizes the squared loss regression objective over  $\mathcal{F}$ , and the minimizer becomes  $f_k$ . More formally,  $f_k := \widehat{\mathcal{T}}_{\mathcal{F}} f_{k-1}$ , where

$$\widehat{\mathcal{T}}_{\mathcal{F}} f' := \arg \min_{f \in \mathcal{F}} \mathcal{L}_D(f; f') \quad (2)$$

$$\mathcal{L}_D(f; f') := \frac{1}{|D|} \sum_{(s,a,r,s') \in D} (f(s, a) - r - \gamma V_{f'}(s'))^2.$$

FQI may oscillate and a fixed point solution may not exist in general (Gordon, 1995). Nevertheless, under conditions which we will specify later, finite sample guarantees for FQI can still be obtained even if the process does not converge.

### 2.3. State Abstractions

A state abstraction  $\phi$  maps  $\mathcal{S}$  to a finite and potentially much smaller abstract state space,  $\mathcal{S}_\phi$ . Naturally,  $\phi$  is often a many-to-one mapping, inducing an equivalence notion over  $\mathcal{S}$  which encodes one’s prior knowledge of equivalent or similar states. A typical use of abstractions in the batch learning setting is to construct a tabular (or *certainty-equivalent*) model from a dataset  $\{(\phi(s), a, r, \phi(s'))\}$ , and compute the optimal policy in the resulting abstract model. There is a long history of studying abstractions, mostly focusing on their approximation guarantees (Whitt, 1978).

We note, however, that there is a direct connection between FQI and certainty-equivalence with abstractions. In particular, value iteration in the model estimated with abstraction  $\phi$  is *exactly equivalent* to FQI with  $\mathcal{F}$  being the class of piece-wise constant functions under  $\phi$ .<sup>6</sup> As such, the characterization of approximation errors in the two bodies of literature are closely related to each other. We will discuss further connections in the rest of this paper.

## 3. Bellman Error Minimization in Batch Reinforcement Learning

In this section, we give a complete analysis of FQI and a related algorithm, with the main results being two sample complexity bounds. Many of the insights and results in this section have either explicitly appeared in or been implicitly hinted by prior work (especially Szepesvári & Munos, 2005; Antos et al., 2008), and we include them because (1) the discussions in the rest of the paper are largely based on these results, and (2) our analyses simplify prior results with-

guarantees often rely on similar but more complicated assumptions (Lazaric et al., 2012). Moreover, our lower bounds are information-theoretic and algorithm-independent.

<sup>6</sup>This result is known anecdotally (see e.g., Pires & Szepesvári, 2016) and we include details in Appendix E for completeness.

out trivializing them, making the high-level insights more accessible. We also improve the results in some aspects.

### 3.1. Sample-Based Bellman Error Minimization

We start by deriving FQI from a slightly unusual perspective due to the aforementioned prior work, which motivates major assumptions in FQI analysis and introduces concepts that are important for later discussions.

Recall that the goal of value-based RL is to find  $f \in \mathcal{F}$  such that  $f \approx \mathcal{T}f$ , that is,  $\|f - \mathcal{T}f\| = 0$  where  $\|\cdot\|$  is some appropriate norm. For example, if  $\mu$  is a distribution supported on the entire  $\mathcal{S} \times \mathcal{A}$ , then  $\|f - \mathcal{T}f\|_{2,\mu}^2 = 0$  would guarantee that  $f = Q^*$ . While such an  $f$  can be found in principle by minimizing  $\|f - \mathcal{T}f\|_{2,\mu}^2$  over  $f \in \mathcal{F}$ , calculating  $\|f - \mathcal{T}f\|$  requires knowledge of the transition dynamics (recall Eq.(1)), which is unknown in the learning setting. Instead, we have access to the dataset  $D = \{(s, a, r, s')\}$ , and it may be tempting to minimize the following objective that is purely a function of data: (Recall  $\mathcal{L}_D$  in Eq.(2))

$$\mathcal{L}_D(f; f) := \frac{1}{|D|} \sum_{(s,a,r,s') \in D} (f(s, a) - r - \gamma V_f(s'))^2.$$

Unfortunately, even with the infinite amount of data, the above objective is still different from the actual Bellman error  $\|f - \mathcal{T}f\|_{2,\mu}^2$  that we wish to minimize. In particular, define  $\mathcal{L}_\mu(\cdot; \cdot) := \mathbb{E}[\mathcal{L}_D(\cdot; \cdot)]$ , where the expectation is w.r.t. the random draw of the dataset  $D$ . We have  $\mathcal{L}_\mu(f; f) =$

$$\|f - \mathcal{T}f\|_{2,\mu}^2 + \gamma^2 \mathbb{E}_{(s,a) \sim \mu} [\mathbb{V}_{s' \sim P(s,a)} [V_f(s')]]. \quad (3)$$

In words,  $\mathcal{L}_\mu(f; f)$  adds a conditional variance term to the desired objective, which incorrectly penalizes functions that have a large variance w.r.t. random state transitions.

**The minimax algorithm**<sup>7</sup> One way to fix the issue is to estimate the conditional variance term in Eq. (3) and subtracting it from  $\mathcal{L}_D(f; f)$ . In fact, it is easy to verify that  $\gamma^2 \mathbb{E}_{(s,a) \sim \mu} [\mathbb{V}_{s' \sim P(s,a)} [V_f(s')]]$  is the Bayes optimal error of the regression problem

$$(s, a) \mapsto r + \gamma V_f(s'). \quad (4)$$

One can estimate it by empirical risk minimization over a rich function class, and the estimate is consistent as long as the function class realizes the Bayes optimal regressor and has bounded statistical complexity. Following this idea, we assume access to another function class  $\mathcal{G} \subset (\mathcal{S} \times \mathcal{A} \rightarrow [0, V_{\max}])$  for solving the regression problem in Eq.(4). The estimated Bayes optimal error is

$$\inf_{g \in \mathcal{G}} \mathcal{L}_D(g; f). \quad (5)$$

<sup>7</sup>Also known under the name “modified Bellman Residual Minimization” (Antos et al., 2008).

A good approximation to  $\|f - \mathcal{T}f\|_{2,\mu}^2$  from data is then  $\sup_{g \in \mathcal{G}} \mathcal{L}_D(f; f) - \mathcal{L}_D(g; f)$ . This suggests that we can simply run the following optimization problem to find  $f \in \mathcal{F}$  that approximates  $Q^*$ :

$$\inf_{f \in \mathcal{F}} \sup_{g \in \mathcal{G}} \mathcal{L}_D(f; f) - \mathcal{L}_D(g; f). \quad (6)$$

Later in this section, we will provide a finite sample analysis of the above minimax algorithm, but before that, we will show that FQI can be viewed as its approximation.

**FQI as an approximation to Eq.(6)** FQI has a close connection to the above program and can be viewed as its approximation, when  $\mathcal{G}$  is chosen to be  $\mathcal{F}$ . Formally,

**Proposition 1.** *Let  $\hat{f}, \hat{g}$  be the solution to Eq.(6) when  $\mathcal{G} = \mathcal{F}$ .*

- If  $\mathcal{L}_D(\hat{f}; \hat{f}) - \mathcal{L}_D(\hat{g}; \hat{f}) = 0$ ,  $\hat{f}$  is a fixed point for FQI.
- Conversely, if  $f_k = f_{k-1}$  holds for some  $k$  in FQI, then  $\hat{f} = \hat{g} = f_k$  is a solution to Eq.(6).
- If  $\mathcal{L}_D(\hat{f}; \hat{f}) - \mathcal{L}_D(\hat{g}; \hat{f}) > 0$ , FQI oscillates and no fixed point exists.

The proof is deferred to Appendix A. The proposition states that the minimax algorithm is more stable than FQI, and when FQI reaches a fixed point, the solutions of the two algorithms coincide. In fact, Dai et al. (2018) derives a closely related algorithm using Fenchel dual and shows that the algorithm is always convergent.

### 3.2. Analysis of FQI and Its Minimax Variant

We provide finite sample guarantees to the two algorithms introduced above; closely related analyses have appeared in prior works (see Section 1 for references), and our version provides a cleaner analysis under simplification assumptions, improves the error rate as a function of sample size, and prepares us for later discussions.

To state the guarantees, we need to introduce the two assumptions that are core to this paper. The first assumption handles distribution shift, and we precede it with the definition of *admissible distributions*.

**Definition 1** (Admissible distributions). *We say a distribution  $\nu \in \Delta(\mathcal{S} \times \mathcal{A})$  is admissible in MDP  $M = (\mathcal{S}, \mathcal{A}, P, R, \gamma, \eta_1)$ , if there exists  $h \geq 0$ , and a (potentially non-stationary and stochastic) policy  $\pi$ , such that  $\nu(s, a) = \Pr[s_h = s, a_h = a | s_1 \sim \eta_1, \pi]$ .*

Intuitively, a distribution is admissible if it can be generated in the MDP by following some policy for a number of timesteps. The following assumption on *concentratability* asserts that all admissible distributions are not “far away” from the data distribution  $\mu$ . The original definition is due to Munos (2003).

**Assumption 1** (Concentratability coefficient). *We assume that there exists  $C < \infty$  s.t. for any admissible  $\nu$ ,*

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \frac{\nu(s, a)}{\mu(s, a)} \leq C.$$

The real (and implicit) assumption here is that  $C$  is manageable large, as our sample complexity bounds scale linearly with  $C$ . Prior works have used more sophisticated definitions (Farahmand et al., 2010).<sup>8</sup> The technicalities introduced are largely orthogonal to the discussions in this paper, so we choose to adopt a much simplified version. Despite the simplification, we will see natural examples that yield small  $C$  under our definition in Section 4. We will also discuss how to relax it using the structure of  $\mathcal{F}$  at the end of the paper.

Next, we introduce the assumption on the representation power of  $\mathcal{F}$  and  $\mathcal{G}$ .

**Assumption 2** (Realizability).  $Q^* \in \mathcal{F}$ .

(When this holds approximately, we measure violation by  $\epsilon_{\mathcal{F}} := \inf_{f \in \mathcal{F}} \|f - \mathcal{T}f\|_{2,\mu}^2$ .)

**Assumption 3** (Completeness).  $\forall f \in \mathcal{F}, \mathcal{T}f \in \mathcal{G}$ .

(When this holds approximately, we measure violation by  $\epsilon_{\mathcal{F},\mathcal{G}} := \sup_{f \in \mathcal{F}} \inf_{g \in \mathcal{G}} \|g - \mathcal{T}f\|_{2,\mu}^2$ .)

These assumptions lead to finite sample guarantees for both the minimax algorithm and FQI. For FQI, since  $\mathcal{G} = \mathcal{F}$ , Assumption 3 essentially states that  $\mathcal{F}$  is closed under operator  $\mathcal{T}$ , hence “completeness”.<sup>9</sup> The assumption is natural from how we derive the minimax algorithm in Sec 3.1, as Eq.(5) is only a consistent estimate of the Bayes optimal error of Eq.(4) if  $\mathcal{G}$  realizes the Bayes optimal regressor, which is  $\mathcal{T}f$ .

A few remarks in order:

1. When  $\mathcal{F} = \mathcal{G}$  is finite, completeness implies realizability.<sup>10</sup> However, completeness is stronger and much less desired than realizability: realizability is monotone in  $\mathcal{F}$  (adding functions to  $\mathcal{F}$  never hurts realizability), while completeness is not (adding functions to  $\mathcal{F}$  may break completeness).
2. While we focus on completeness, it is not the only condition that leads to guarantees for ADP algorithms. We discuss alternative assumptions in Section 6.

<sup>8</sup>This often comes at the cost of their bound being not *a priori*, i.e., having a dependence on the randomness of data, initialization, and tie-breaking in optimization.

<sup>9</sup>In the literature, the violation of completeness when  $\mathcal{F} = \mathcal{G}$ ,  $\epsilon_{\mathcal{F},\mathcal{F}}$ , is called *inherent Bellman error*.

<sup>10</sup>This is because  $\mathcal{T}^k f$  never repeats itself, as its  $\ell_\infty$  distance to  $Q^*$  shrinks exponentially with a rate of  $\gamma$  due to contraction.

Now we are ready to state the sample complexity results. In Appendices C and D we provide more general error bounds (Theorems 11 and 17) that handle the approximate case where  $\epsilon_{\mathcal{F}}$  and  $\epsilon_{\mathcal{F},\mathcal{G}}$  are not zero and iteration  $k$  is finite. To keep the main text focused and accessible, we only present their sample complexity corollaries in the exact case.

**Theorem 2** (Sample complexity of FQI). *Given a dataset  $D = \{(s, a, r, s')\}$  with sample size  $|D| = n$  and  $\mathcal{F}$  that satisfies completeness (Assumption 3 when  $\mathcal{G} = \mathcal{F}$ ), w.p.  $\geq 1 - \delta$ , the output policy of FQI after  $k$  iterations,  $\pi_{f,k}$ , satisfies  $v^* - v^{\pi_{f,k}} \leq \epsilon \cdot V_{\max}$  when  $k \rightarrow \infty$  and<sup>11</sup>*

$$n = O\left(\frac{C \ln \frac{|\mathcal{F}|}{\delta}}{\epsilon^2(1-\gamma)^4}\right).$$

**Theorem 3** (Sample complexity of the minimax variant). *Given a dataset  $D = \{(s, a, r, s')\}$  with sample size  $|D| = n$  and  $\mathcal{F}, \mathcal{G}$  that satisfy realizability (Assumption 2) and completeness (Assumption 3) respectively, w.p.  $\geq 1 - \delta$ , the output policy of the minimax algorithm (Eq.(6)),  $\pi_{f^*}$ , satisfies  $v^* - v^{\pi_{f^*}} \leq \epsilon \cdot V_{\max}$ , if  $n = O\left(\frac{C \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{\epsilon^2(1-\gamma)^4}\right)$ .*

Our results show that the suboptimality  $\epsilon$  decreases in the rate of  $n^{-1/2}$  when realizability and completeness hold exactly, and the more general error bounds (Theorems 11 and 17) degrade gracefully from the exact case as  $\epsilon_{\mathcal{F},\mathcal{F}}$  (or  $\epsilon_{\mathcal{F}}$  and  $\epsilon_{\mathcal{F},\mathcal{G}}$ ) increases. This is obtained via the use of Bernstein’s inequality to achieve fast rate in least square regression. While results similar to Theorems 2 and 11 exist (Farahmand 2011, Chapter 5; see also Lazaric et al. (2012); Pires & Szepesvári (2012); Farahmand et al. (2016)), according to our knowledge, fast rate for the minimax algorithm has not been established before: for example, Antos et al. (2008); Munos & Szepesvári (2008) obtain an error rate of  $n^{-1/4}$  in closely related settings, but their rates do not improve to  $n^{-1/2}$  in the absence of approximation.<sup>12</sup> The major limitation of our result is the assumption of finite  $\mathcal{F}$  and  $\mathcal{G}$  due to our minimal setup, and we refer readers to Yang et al. (2019) for a recent analysis that specializes in ReLU networks.<sup>13</sup>

We do not discuss the proofs in further details since the improvement in error rate is a side-product and this section is mainly meant to simplify prior analyses and provide a basis for subsequent discussions. Interested readers are invited to consult Appendices C and D where we provide sketched outlines as well as detailed proofs.

<sup>11</sup>Only absolute constants are suppressed in Big-Oh notations.

<sup>12</sup>Note however that they handle infinite function classes. In fact, Munos & Szepesvári (2008, pg.831) have discussed the possibility of an  $n^{-1/2}$  result, which we obtain here. See the beginning of Appendix C for further discussions.

<sup>13</sup>Their analysis modifies the FQI algorithm and samples fresh data in each iteration, dodging some of the technical difficulties due to reusing the same batch of data, which we handle here.

## 4. On Concentratability

In this section, we establish the necessity of Assumption 1 and show natural examples where concentratability is low. While it is easy to construct a counterexample of missing data<sup>14</sup> against removing Assumption 1, such a counterexample only reflects a trivial failure mode due to an adversarial choice of data. What we show is a deeper and nontrivial failure mode: Even with the *most favorable* data distribution, polynomial sample complexity is precluded if we put no restriction on MDP dynamics. This result improves our understanding on concentratability, and shows that this assumption is not only about the data distribution, but also (and perhaps more) about the environment and the state distributions induced therein.

### 4.1. Lower Bound

To show that low concentratability is necessary, we prove a hardness result, where both realizability and completeness hold, and an algorithm has the freedom to choose *any* data distribution  $\mu$  that is favorable, yet *no algorithm* can achieve  $\text{poly}(|\mathcal{A}|, \frac{1}{1-\gamma}, \ln |\mathcal{F}|, \ln |\mathcal{G}|, \frac{1}{\epsilon}, \frac{1}{\delta})$  sample complexity. Crucially, the concentratability coefficient of any data distribution on the worst-case MDP is always exponential in horizon, so the lower bound does not conflict with the upper bounds in Section 3, as the exponential sample complexity would have been explained away by the dependence on  $C$ .

**Theorem 4.** *There exists a family of MDPs  $\mathcal{M}$  (they share the same  $S, \mathcal{A}, \gamma$ ),  $\mathcal{F}$  that realizes the  $Q^*$  of every MDP in the family, and  $\mathcal{G}$  that realizes  $\mathcal{T}_{M'}f$  for any  $M' \in \mathcal{M}$  and any  $f \in \mathcal{F}$ , such that: for any data distribution and any batch algorithm with  $(\mathcal{F}, \mathcal{G})$  as input, an adversary can choose an MDP from the family, such that the sample complexity for the algorithm to find an  $\epsilon$ -optimal policy cannot be  $\text{poly}(|\mathcal{A}|, \frac{1}{1-\gamma}, \ln |\mathcal{F}|, \ln |\mathcal{G}|, \frac{1}{\epsilon}, \frac{1}{\delta})$ .*

*Proof.* We construct  $\mathcal{M}$ , a family of hard MDPs, and prove the theorem via the combination of two arguments:

1. All algorithms are subject to an exponential lower bound (w.r.t. the horizon) even if (a) they have compact  $\mathcal{F}$  and  $\mathcal{G}$  that satisfy realizability and completeness as inputs, and (b) they can perform *exploration* during data collection.
2. Since the MDPs in the construction share the same deterministic transition dynamics, the combination of any data distribution and any batch RL algorithm is a *special case* of an exploration algorithm.

We first provide argument (1), which reuses the construction by Krishnamurthy et al. (2016). Let each instance of  $\mathcal{M}$  be a complete tree with branching factor  $|\mathcal{A}|$  and depth

<sup>14</sup> That is,  $\mu$  puts 0 probability on important states and actions.

$H = \lfloor 1/(1 - \gamma) \rfloor$ . Transitions are deterministic, and only leaf nodes have non-zero rewards. All leaves give  $\text{Ber}(1/2)$  rewards, except for one that gives  $\text{Ber}(1/2 + \epsilon)$ . Changing the position of this optimal leaf yields a family of  $|\mathcal{A}|^H$  MDPs, and in order to achieve a suboptimality that is a constant fraction of  $\epsilon$ , the algorithm is required to identify this optimal leaf.<sup>15</sup> In fact, the problem is equivalent to the hard instances of best arm identification with  $|\mathcal{A}|^H$  arms, so even if an algorithm can perform active exploration, the sample complexity is still  $\Omega(|\mathcal{A}|^H \ln(1/\delta)/\epsilon^2)$  (see Krishnamurthy et al. (2016) for details, who use standard techniques from Auer et al. (2002)).

Now we provide  $\mathcal{F}$  and  $\mathcal{G}$  that (1) satisfy Assumptions 2 and 3, (2) do not provide any information other than the fact that the problem is in  $\mathcal{M}$ , and (3) have “small” logarithmic sizes so that  $\ln |\mathcal{F}|$  and  $\ln |\mathcal{G}|$  cannot explain away the exponential sample complexity. Let  $\mathcal{F} = \{Q_{M'}^* : M' \in \mathcal{M}\}$ , where the subscript specifies the MDP with respect to which we compute  $Q^*$ . Let  $\mathcal{G} = \{\mathcal{T}_{M'}, Q_{M''}^* : M', M'' \in \mathcal{M}\}$ . Such  $\mathcal{F}$  and  $\mathcal{G}$  satisfy realizability and completeness by definition, and have statistical complexities  $\ln |\mathcal{F}| = H \ln |\mathcal{A}|$  and  $\ln |\mathcal{G}| \leq 2H \ln |\mathcal{A}|$ , respectively. With this, we conclude that any *exploration* algorithm cannot obtain  $\text{poly}(|\mathcal{A}|, \frac{1}{1-\gamma}, \ln |\mathcal{F}|, \ln |\mathcal{G}|, \frac{1}{\epsilon})$  sample complexity.

We complete the proof with the second argument. Note that all the MDPs in  $\mathcal{M}$  only differ in leaf rewards and share the same deterministic transition dynamics. Therefore, a learner with the ability to actively explore can mimic the combination of *any data distribution*  $\mu \in \Delta(\mathcal{S} \times \mathcal{A})$  and *any batch RL algorithm*, by (1) collecting data from  $\mu$  (which is always doable due to known and deterministic transitions), and (2) running the batch algorithm after data is collected. This completes the proof.  $\square$

## 4.2. Natural Examples

We have shown that polynomial learning is precluded if no restriction is put on the MDP dynamics, even if data is chosen in a favorable manner. The next question is, is low concentratability common, or at least found in interesting problems? In general, even if the data distribution  $\mu$  is uniform over the state-action space, the worst-case  $C$  might still scale with  $|\mathcal{S} \times \mathcal{A}|$ , which can be too large in challenging RL problems for the guarantees to be any meaningful. To this end, Munos (2007) has provided several carefully constructed tabular examples, demonstrating that  $C$  does not always scale badly. However, are there more general problem families that capture RL scenarios found in empirical work, yet always yield a bounded  $C$ ?

**Example in problems with rich observations** We find an-

<sup>15</sup>All leaf rewards are discounted by only a constant when  $\gamma \rightarrow 1$ , as  $\gamma^{1/(1-\gamma)} \rightarrow e^{-1}$ .

swers to the above problem in recent development of PAC exploration in rich-observation problems (Krishnamurthy et al., 2016; Jiang et al., 2017; Dann et al., 2018), where a general low-rank condition (a.k.a. *Bellman rank* (Jiang et al., 2017)) has been identified that enables sample-efficient exploration under function approximation. One of the prominent examples where such a condition holds is inspired by “visual gridworld” environments in empirical RL research (see e.g., Johnson et al., 2016): the dynamics are defined over a small number of hidden states (e.g., grids), and the agent receives high dimensional observations that are generated i.i.d. from the hidden states (e.g., raw-pixel images as observations). Below we show that in these environments, there always exists a data distribution that yields small  $C$  for batch learning, and such a distribution can be naturally generated as a mixture of admissible distributions. We include an informal statement below, deferring the precise version and the proof to Appendix B.

**Proposition 5 (Informal).** *Let  $M$  be a reactive POMDP as defined in Jiang et al. (2017), where the underlying hidden state space  $\mathcal{Z}$  is finite but the (Markov) observation space  $\mathcal{S}$  can be arbitrarily large. There always exists a state-action distribution  $\mu$  such that  $C = |\mathcal{Z} \times \mathcal{A}|$  satisfies Assumption 1. Furthermore,  $\mu$  can be obtained by taking a probability mixture of several admissible distributions.*

Similar results can be established for other structures studied by Jiang et al. (2017) (e.g., large MDPs with low-rank transitions), which we omit here. These results suggest that *Bellman rank is the counterpart for concentratability coefficient in the online exploration setting*. Further implications and how to leverage this connection to improve the definition of concentratability will be discussed in Section 6.

## 5. On Completeness

### 5.1. Towards an Information-Theoretic Lower Bound in the Absence of Completeness

We would also like to establish the necessity of completeness by showing that, there exist hard MDPs that cannot be efficiently learned with value-function approximation, even under low concentratability and realizability (Assumptions 1 and 2).<sup>16</sup> In fact, *algorithm-specific* hardness results have been known for a long time (see e.g., Van Roy, 1994; Gordon, 1995; Tsitsiklis & Van Roy, 1997), where ADP algorithms are shown to diverge even in MDPs with a small number of states, when the algorithm is forced to work with a restricted class of functions.<sup>17</sup> Unfortunately, such

<sup>16</sup>Note that the existence of such a lower bound would not imply that completeness is indispensable. Rather it simply states that realizability alone is insufficient, and we need stronger conditions on  $\mathcal{F}$ , for which completeness is a candidate.

<sup>17</sup>Interested readers can consult Agrawal (2018). See also Dann et al. (2018, Theorem 45) for a more plain example.

hardness results are insufficient to confirm the fundamental difficulty of the problem, and it is important to seek *information-theoretic* lower bounds.

While we are not able to obtain such a lower bound, what we find is that the counterexample (if it exists) must be highly nontrivial and probably need ideas that are not present in standard statistical learning theory (SLT) and RL literature. More concretely, we show that two general proof styles are destined to fail in such a task, as polynomial sample complexity can be achieved information-theoretically.

**Exponential-sized model family will not work** Standard lower bounds in SLT often start with the construction of a family of problem instances that has an exponential size (Yu, 1997).<sup>18</sup> We show that this will simply never work, which is a direct corollary of Theorem 3:

**Corollary 6** (Batch model-based RL only needs realizability). *Let  $D = \{(s, a, r, s')\}$  be a dataset with sample size  $|D| = n$ ,  $C$  as defined in Assumption 1, and  $\mathcal{M}$  a model class that realizes the true MDP  $M$ , i.e.,  $M \in \mathcal{M}$ . There exists an (information-theoretic) algorithm that takes  $\mathcal{M}$  as input and return an  $(\epsilon V_{\max})$ -optimal policy w.p.  $\geq 1 - \delta$ , if  $n = O\left(\frac{C \ln \frac{|\mathcal{M}|}{\delta}}{\epsilon^2(1-\gamma)^4}\right)$ .*

*Proof.* We use the same idea as the proof of Theorem 4: Let  $\mathcal{F} = \{Q_{M'}^* : M' \in \mathcal{M}\}$ , and  $\mathcal{G} = \{\mathcal{T}_{M'} Q_{M''}^* : M', M'' \in \mathcal{M}\}$ . Note that  $\ln |\mathcal{F}| \leq \ln |\mathcal{M}|$ , and  $\ln |\mathcal{G}| \leq 2 \ln |\mathcal{M}|$ .  $(\mathcal{F}, \mathcal{G})$  satisfy both realizability and completeness, so we apply the minimax algorithm (Eq.(6)) and the guarantee in Theorem 3 immediately holds.  $\square$

Essentially, this result shows that batch model-based RL can succeed with realizability as the only representation condition for the model class, because we can reduce it to value-based learning and obtain completeness *for free*. This illustrates a significant barrier to an algorithm-independent lower bound, that in an information-theoretic setting, the learner can always specialize in the family of hard instances and have the freedom to choose its algorithm style, thus can be *model-based*. However, in the context of value-function approximation, it is obvious that we are assuming no prior knowledge of the model class and hence cannot run any model-based algorithm. How can we encode such a constraint mathematically?

**Tabular MDPs with a restricted value-function class will not work** Sutton & Barto (2018, Section 11.6) proposes a clever way to prevent the learner to be model-based for linear function approximation, and a closely related definition is recently given by Sun et al. (2019) that applies to

arbitrary function classes.

The idea is the following: Instead of providing the dataset  $D = \{(s, a, r, s')\}$  directly, we preprocess the data and mask the identity of  $s$  (and  $s'$ ). While  $s$  is not directly observable, the learner can query the evaluation of any  $f \in \mathcal{F}$  on  $s$  for any  $a \in \mathcal{A}$ . That is, we represent each state  $s$  by its *value profile*,  $\{f(s, a) : f \in \mathcal{F}, a \in \mathcal{A}\}$ . This definition agrees with intuition and can be used to express a wide range of popular algorithms, including FQI.

Using this definition, Sutton & Barto (2018) proves a result closely related to what we aim at here: they show that the Bellman error  $\|f - \mathcal{T}f\|$  is not learnable. In particular, there exist two MDPs (with finite and constant-sized state space) and a value function, such that (1) a value-based learner (who only has access to the value profiles of states) cannot distinguish between the data coming from the two MDPs, and (2) the Bellman error of the value function is different in the two MDPs.

While encouraging and promising, their constructions have a crucial caveat for our purpose, that the value function class is not realizable.<sup>19</sup> With further investigation, we sadly find that such a caveat is fundamental: no information-theoretic lower bound can be shown if realizability holds in naïve tabular constructions with a constant-sized state-action space and uniform data, hence value profile cannot be the *only* mechanism to induce hardness. In fact, we can prove a stronger result than we need here for  $\mathcal{S}$  and  $\mathcal{A}$  that are not necessarily constant-sized:

**Proposition 7.** *Let  $M$  be an MDP with a finite state space and  $\mathcal{F}$  a realizable function class. Given a dataset  $D = \{(s, a, r, s')\}$  where each  $(s, a)$  receives  $\Omega(|D|/|\mathcal{S} \times \mathcal{A}|)$  samples, there exists an algorithm that only operates on states via their value profiles yet enjoy  $\text{poly}(|\mathcal{S}|, |\mathcal{A}|, \frac{1}{1-\gamma}, \frac{1}{\epsilon}, \frac{1}{\delta})$  sample complexity.*

*Proof Sketch.* (See full proof in Appendix F.) If every  $s \in \mathcal{S}$  has a unique value profile, the state is perfectly decodable and thus one can simply compute the optimal policy of the certainty-equivalent model. If a set of states share exactly the same value profile—and w.l.o.g. let's consider 2 states,  $s_1$  and  $s_2$ —realizability implies that  $Q^*(s_1, a) = Q^*(s_2, a), \forall a \in \mathcal{A}$ . Now consider the algorithm that treat all states with the same value profile as the same state, which essentially uses a state abstraction that is  $Q^*$ -irrelevant (Li et al., 2006). It is known that certainty-equivalence with  $Q^*$ -irrelevant abstraction is consistent and enjoys polynomial sample complexity when each state-action pair receives enough data (Li, 2009; Hutter, 2014; Jiang et al., 2015; Abel et al., 2016; Jiang, 2018).  $\square$

<sup>18</sup>In fact, our Theorem 4 also follows this style, whose construction is due to Krishnamurthy et al. (2016); Jiang et al. (2017).

<sup>19</sup>They force two states who have different optimal values to share the same features for linear function approximation.

Given that we fail to obtain the lower bound, a conjecture is made below and we hope to resolve it in future work.

**Conjecture 8.** *There exists a family of MDPs  $\mathcal{M}$  that share the same  $\mathcal{S}$ ,  $\mathcal{A}$ , and  $\gamma$ , such that: any algorithm with  $\mathcal{F} = \{Q_{M'}^* : M' \in \mathcal{M}\}$  as input that can only access states via value profiles cannot have  $\text{poly}(\frac{1}{1-\gamma}, C, \ln |\mathcal{F}|, \frac{1}{\epsilon}, \frac{1}{\delta})$  sample complexity.*

## 5.2. Connection to Bisimulation

As the last piece of technical result of this paper, we show that when  $\mathcal{F}$  is a space of piece-wise constant functions under a partition induced by state abstraction  $\phi$ , the notion of completeness (Assumption 3,  $\mathcal{F} = \mathcal{G}$ ) is exactly equivalent to a long-studied type of abstractions, known as *bisimulation* (Whitt, 1978; Even-Dar & Mansour, 2003; Ravindran, 2004; Li et al., 2006).

**Definition 2** (Bisimulation). *An abstraction  $\phi : \mathcal{S} \rightarrow \mathcal{S}_\phi$  is a bisimulation in an MDP  $M$ , if  $\forall s_1, s_2$  where  $\phi(s_1) = \phi(s_2)$  (i.e., they are aggregated),  $R(s_1, a) = R(s_2, a)$  and  $\sum_{s \in \phi^{-1}(x)} P(s|s_1, a) = \sum_{s \in \phi^{-1}(x)} P(s|s_2, a)$  for all  $a \in \mathcal{A}$ ,  $x \in \mathcal{S}_\phi$ .*

**Definition 3** (Piece-wise constant function class). *Given an abstraction  $\phi$ , define  $\mathcal{F}^\phi \subset (\mathcal{S} \times \mathcal{A} \rightarrow [0, V_{\max}])$  as the set of all functions  $f$  that are piece-wise constant under  $\phi$ . That is,  $\forall s_1, s_2 \in \mathcal{S}$  where  $\phi(s_1) = \phi(s_2)$ , we have  $f(s_1, a) = f(s_2, a)$ ,  $\forall a \in \mathcal{A}$ .*

**Proposition 9.**  $\phi$  is bisimulation  $\Leftrightarrow \mathcal{F}^\phi$  satisfies completeness (Assumption 3 with  $\mathcal{F} = \mathcal{G} = \mathcal{F}^\phi$ ).

The “ $\Rightarrow$ ” part is trivial, but the “ $\Leftarrow$ ” part is less obvious. The proof shows that if  $\phi$  is not a bisimulation, we can find  $f \in \mathcal{F}^\phi$  either to witness the reward error or the transition error, and in the latter case, the choice of  $f$  achieves the maximum discrepancy in an integral probability metric (Müller, 1997) interpretation of the bisimulation condition on transition dynamics. Details are provided in Appendix E, where we prove a stronger result that relates the approximation error of bisimulation to the violation of completeness.

## 6. Discussions and Related Work

In this paper, we examine the common assumptions that enable finite sample guarantees for value-function approximation methods. Concretely, we provide an information-theoretic lower bound in Section 4.1, showing that not constraining the concentratability coefficient  $C$  immediately precludes sample-efficient learning even with benign data. We also introduce a general family of problems of interest in empirical RL that yield low concentratability (Section 4.2).

In comparison, the necessity of completeness is still a mystery, and our investigation in Section 5.1 mostly shows the highly nontrivial nature of the lower bound (assuming it ex-

ists) as we eliminate two general proof styles. We hope these negative results can guide the search for novel constructions that reflect the fundamental difficulties of reinforcement learning in the function approximation setting.

We conclude the paper with some discussions.

**Alternative assumptions to completeness** As we note in Section 5.1, even if Conjecture 8 is true, it would not imply that completeness is absolutely necessary, as other assumptions may also break the lower bound. Furthermore, additional assumptions are not necessarily made on the value-function class (e.g., that  $\widehat{\mathcal{T}}_{\mathcal{F}}$  being a contraction (Gordon, 1995; Szepesvári & Smart, 2004; Lizotte, 2011; Pires & Szepesvári, 2016)), and can instead take the form of requiring another function class to realize other objects of interest, such as state distributions (Chen et al., 2018; Liu et al., 2018). Regardless, all of these approaches face the same fundamental question on the necessity of the additional/stronger assumptions being made, to which our Conjecture 8 is an important piece if not the final answer. We hope to resolve this important open question in the future.

**Related work that has not been covered** The conjectured insufficiency of realizability (Conjecture 8) is related to various undesirable phenomena in learning with bootstrapped targets, which has been of constant interest to RL researchers (Sutton, 2015; Van Hasselt et al., 2018; Lu et al., 2018). As far as we know, all existing efforts that investigate this issue are algorithm-specific (apart from Sutton & Barto (2018, Section 11.6) and the references therein, which has been discussed in Section 5.1), and our information-theoretic perspective is novel.

**Relaxation of Assumption 1 using the structure of  $\mathcal{F}$**  The concentratability coefficient  $C$  is defined as a function of the MDP, even in its most complicated version (Farahmand et al., 2010). In Section 4.2 we discover a connection to *Bellman rank* (Jiang et al., 2017), which can be viewed as its counterpart for online exploration. Interestingly, Bellman rank depends both on the environmental dynamics *and the function class  $\mathcal{F}$* , and in some cases, the latter dependence is crucial to obtaining low-rankness (e.g., for Linear Quadratic Regulators; see their Proposition 5). Similarly, we may improve the definition of concentratability and make it more widely applicable by incorporating  $\mathcal{F}$  into the definition. In Appendix G, we discuss some preliminary ideas based on the theoretical results in this paper.

## Acknowledgements

We gratefully thank the constructive comments from Alekh Agarwal and Anonymous Reviewer #3.



## References

- Abel, D., Hershkowitz, D. E., and Littman, M. L. Near optimal behavior via approximate state abstraction. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning-Volume 48*, pp. 2915–2923. JMLR. org, 2016.
- Agrawal, S. *IEOR 8100: Reinforcement Learning. Lecture 4: Approximate Dynamic Programming*. Columbia University, 2018. <https://ieor8100.github.io/rl/docs/Lecture%204%20-%20approximate%20DP.pdf>.
- Antos, A., Szepesvári, C., and Munos, R. Learning near-optimal policies with bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 71(1):89–129, 2008.
- Auer, P., Cesa-Bianchi, N., and Fischer, P. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- Bertsekas, D. P. and Tsitsiklis, J. N. *Neuro-Dynamic Programming*. Athena Scientific, Belmont, MA, 1996.
- Chen, Y., Li, L., and Wang, M. Scalable bilinear  $\pi$  learning using state and action features. *arXiv preprint arXiv:1804.10328*, 2018.
- Dai, B., Shaw, A., Li, L., Xiao, L., He, N., Liu, Z., Chen, J., and Song, L. Sbeed: Convergent reinforcement learning with nonlinear function approximation. In *International Conference on Machine Learning*, pp. 1133–1142, 2018.
- Dann, C., Jiang, N., Krishnamurthy, A., Agarwal, A., Langford, J., and Schapire, R. E. On Oracle-Efficient PAC RL with Rich Observations. In *Advances in Neural Information Processing Systems*, pp. 1429–1439, 2018.
- Ernst, D., Geurts, P., and Wehenkel, L. Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research*, 6:503–556, 2005.
- Even-Dar, E. and Mansour, Y. Approximate equivalence of Markov decision processes. In *Learning Theory and Kernel Machines*, pp. 581–594. 2003.
- Farahmand, A.-m. Regularization in reinforcement learning. 2011.
- Farahmand, A.-m., Szepesvári, C., and Munos, R. Error Propagation for Approximate Policy and Value Iteration. In *Advances in Neural Information Processing Systems*, pp. 568–576, 2010.
- Farahmand, A.-m., Ghavamzadeh, M., Szepesvári, C., and Mannor, S. Regularized policy iteration with nonparametric function spaces. *The Journal of Machine Learning Research*, 17(1):4809–4874, 2016.
- Gordon, G. J. Stable function approximation in dynamic programming. In *Proceedings of the twelfth international conference on machine learning*, pp. 261–268, 1995.
- Hutter, M. Extreme state aggregation beyond mdps. In *International Conference on Algorithmic Learning Theory*, pp. 185–199. Springer, 2014.
- Jiang, N. *CS 598: Notes on State Abstractions*. University of Illinois at Urbana-Champaign, 2018. <http://nanjiang.cs.illinois.edu/files/cs598/note4.pdf>.
- Jiang, N., Kulesza, A., and Singh, S. Abstraction Selection in Model-based Reinforcement Learning. In *Proceedings of the 32nd International Conference on Machine Learning*, pp. 179–188, 2015.
- Jiang, N., Krishnamurthy, A., Agarwal, A., Langford, J., and Schapire, R. E. Contextual Decision Processes with low Bellman rank are PAC-learnable. In *International Conference on Machine Learning*, 2017.
- Johnson, M., Hofmann, K., Hutton, T., and Bignell, D. The malmo platform for artificial intelligence experimentation. In *International joint conference on artificial intelligence (IJCAI)*, pp. 4246, 2016.
- Kakade, S. and Langford, J. Approximately Optimal Approximate Reinforcement Learning. In *Proceedings of the 19th International Conference on Machine Learning*, volume 2, pp. 267–274, 2002.
- Kakade, S. M. *On the sample complexity of reinforcement learning*. PhD thesis, University of College London, 2003.
- Krishnamurthy, A., Agarwal, A., and Langford, J. PAC reinforcement learning with rich observations. In *Advances in Neural Information Processing Systems*, pp. 1840–1848, 2016.
- Lagoudakis, M. G. and Parr, R. Least-squares policy iteration. *The Journal of Machine Learning Research*, 4: 1107–1149, 2003.
- Lazaric, A., Ghavamzadeh, M., and Munos, R. Finite-sample analysis of least-squares policy iteration. *The Journal of Machine Learning Research*, 13(1):3041–3074, 2012.
- Li, L. *A unifying framework for computational reinforcement learning theory*. PhD thesis, Rutgers, The State University of New Jersey, 2009.
- Li, L., Walsh, T. J., and Littman, M. L. Towards a unified theory of state abstraction for MDPs. In *Proceedings of the 9th International Symposium on Artificial Intelligence and Mathematics*, pp. 531–539, 2006.

- Liu, Q., Li, L., Tang, Z., and Zhou, D. Breaking the curse of horizon: Infinite-horizon off-policy estimation. In *Advances in Neural Information Processing Systems*, pp. 5361–5371, 2018.
- Lizotte, D. J. Convergent fitted value iteration with linear function approximation. In *Advances in Neural Information Processing Systems*, pp. 2537–2545, 2011.
- Lu, T., Schuurmans, D., and Boutilier, C. Non-delusional q-learning and value-iteration. In *Advances in Neural Information Processing Systems*, pp. 9971–9981, 2018.
- Maillard, O.-A., Munos, R., Lazaric, A., and Ghavamzadeh, M. Finite-sample analysis of bellman residual minimization. In *Proceedings of 2nd Asian Conference on Machine Learning*, pp. 299–314, 2010.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529–533, 2015.
- Müller, A. Integral probability metrics and their generating classes of functions. *Advances in Applied Probability*, 29(2):429–443, 1997.
- Munos, R. Error bounds for approximate policy iteration. In *ICML*, volume 3, pp. 560–567, 2003.
- Munos, R. Performance bounds in  $l_p$ -norm for approximate value iteration. *SIAM journal on control and optimization*, 46(2):541–561, 2007.
- Munos, R. and Szepesvári, C. Finite-time bounds for fitted value iteration. *Journal of Machine Learning Research*, 9(May):815–857, 2008.
- Pires, B. A. and Szepesvári, C. Statistical linear estimation with penalized estimators: an application to reinforcement learning. *arXiv preprint arXiv:1206.6444*, 2012.
- Pires, B. Á. and Szepesvári, C. Policy error bounds for model-based reinforcement learning with factored linear models. In *Conference on Learning Theory*, pp. 121–151, 2016.
- Ravindran, B. *An algebraic approach to abstraction in reinforcement learning*. PhD thesis, University of Massachusetts Amherst, 2004.
- Riedmiller, M. Neural fitted q iteration—first experiences with a data efficient neural reinforcement learning method. In *European Conference on Machine Learning*, pp. 317–328. Springer, 2005.
- Singh, S. and Yee, R. An upper bound on the loss from approximate optimal-value functions. *Machine Learning*, 16(3):227–233, 1994.
- Sun, W., Jiang, N., Krishnamurthy, A., Agarwal, A., and Langford, J. Model-based RL in Contextual Decision Processes: PAC bounds and Exponential Improvements over Model-free Approaches. In *Conference on Learning Theory*, 2019.
- Sutton, R. Introduction to reinforcement learning with function approximation. In *Tutorial at the Conference on Neural Information Processing Systems*, 2015.
- Sutton, R. S. and Barto, A. G. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, March 1998. ISBN 0-262-19398-1.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Szepesvári, C. Algorithms for reinforcement learning. *Synthesis lectures on artificial intelligence and machine learning*, 4(1):1–103, 2010.
- Szepesvári, C. and Munos, R. Finite time bounds for sampling based fitted value iteration. In *Proceedings of the 22nd international conference on Machine learning*, pp. 880–887. ACM, 2005.
- Szepesvári, C. and Smart, W. D. Interpolation-based q-learning. In *Proceedings of the twenty-first international conference on Machine learning*, pp. 100. ACM, 2004.
- Tosatto, S., Pirodda, M., D’Eramo, C., and Restelli, M. Boosted fitted q-iteration. In *Proceedings of the 34th International Conference on Machine Learning—Volume 70*, pp. 3434–3443. JMLR. org, 2017.
- Tsitsiklis, J. N. and Van Roy, B. An analysis of temporal-difference learning with function approximation. *IEEE TRANSACTIONS ON AUTOMATIC CONTROL*, 42(5), 1997.
- Van Hasselt, H., Doron, Y., Strub, F., Hessel, M., Sonnerat, N., and Modayil, J. Deep reinforcement learning and the deadly triad. *arXiv preprint arXiv:1812.02648*, 2018.
- Van Roy, B. *Feature-based methods for large scale dynamic programming*. PhD thesis, Massachusetts Institute of Technology, 1994.
- Whitt, W. Approximations of dynamic programs, I. *Mathematics of Operations Research*, 3(3):231–243, 1978.
- Yang, Z., Xie, Y., and Wang, Z. A Theoretical Analysis of Deep Q-Learning. *arXiv preprint arXiv:1901.00137*, 2019.
- Yu, B. Assouad, Fano, and Le Cam. In *Festschrift for Lucien Le Cam*, pp. 423–435. Springer, 1997.

## A. Proof of Proposition 1

**Claim 1:** Since  $\mathcal{G} = \mathcal{F}$ , we have that,  $\forall f \in \mathcal{F}$ ,

$$\arg \max_{g \in \mathcal{G}} (\mathcal{L}_D(f; f) - \mathcal{L}_D(g; f)) = \arg \max_{g \in \mathcal{G}} -\mathcal{L}_D(g; f) = \arg \min_{g \in \mathcal{G}} \mathcal{L}_D(g; f) = \widehat{\mathcal{T}}_{\mathcal{G}} f = \widehat{\mathcal{T}}_{\mathcal{F}} f.$$

Therefore,  $\mathcal{L}_D(\hat{g}; \hat{f}) = \mathcal{L}_D(\widehat{\mathcal{T}}_{\mathcal{F}} \hat{f}; \hat{f})$ , and the condition  $\mathcal{L}_D(\hat{f}; \hat{f}) - \mathcal{L}_D(\hat{g}; \hat{f}) = 0$  gives us that  $\mathcal{L}_D(\hat{f}; \hat{f}) - \mathcal{L}_D(\widehat{\mathcal{T}}_{\mathcal{F}} \hat{f}; \hat{f}) = 0$ . From the definition, we know that  $\widehat{\mathcal{T}}_{\mathcal{F}} \hat{f} = \arg \min_{f \in \mathcal{F}} \mathcal{L}_D(f; \hat{f})$ . Hence  $\mathcal{L}_D(\hat{f}; \hat{f}) = \mathcal{L}_D(\widehat{\mathcal{T}}_{\mathcal{F}} \hat{f}; \hat{f}) = \min_{f \in \mathcal{F}} \mathcal{L}_D(f; \hat{f})$ , which means  $\hat{f} = \arg \min_{f \in \mathcal{F}} \mathcal{L}_D(f; \hat{f})$  and  $\hat{f}$  is a fixed point for FQI.

**Claim 2:** Since  $\mathcal{G} = \mathcal{F}$ , for any  $f \in \mathcal{F}$ , we can always choose  $g = f$ . Therefore, for any  $f \in \mathcal{F}$ ,  $\sup_{g \in \mathcal{G}} \mathcal{L}_D(f; f) - \mathcal{L}_D(g; f) \geq 0$ , which further means  $\inf_{f \in \mathcal{F}} \sup_{g \in \mathcal{G}} \mathcal{L}_D(f; f) - \mathcal{L}_D(g; f) \geq 0$ , and the value of the optimization problem is non-negative. If we have  $f_k = f_{k-1}$  for some  $k$  in FQI, then we know that  $f_{k-1} \in \mathcal{F}$ ,  $f_k \in \mathcal{F} = \mathcal{G}$  and  $\mathcal{L}_D(f_{k-1}; f_{k-1}) - \mathcal{L}_D(f_k; f_{k-1}) = 0$ . This tells us that  $f = f_{k-1}$  and  $g = f_k$  achieve the optimal value, so  $\hat{f} = \hat{g} = f_k$  is a solution to Eq.(6).

**Claim 3:** Prove by contradiction. If FQI does not oscillate and a fixed point of FQI is  $f_{k-1}(= f_k)$ , the previous result gives us that  $\hat{f} = \hat{g} = f_k$  is a solution to Eq.(6), with the minimax objective value being  $\mathcal{L}_D(f_k; f_k) - \mathcal{L}_D(f_k; f_k) = 0$ . Contradiction.

## B. Example of Low Concentratability in Rich-Observation Problems

**Definition 4** (Reactive POMDPs (Jiang et al., 2017)). *A reactive POMDP is a decision process specified by a finite hidden state space  $\mathcal{Z}$ , an (arbitrarily large) observation space  $\mathcal{S}$ , an action space  $\mathcal{A}$ , hidden state dynamics  $\Gamma : \mathcal{Z} \times \mathcal{A} \rightarrow \Delta(\mathcal{Z})$ , an initial hidden state distribution  $\Gamma_1 \in \Delta(\mathcal{Z})$ , an emission process  $P : \mathcal{Z} \rightarrow \Delta(\mathcal{S})$ , a reward function  $R : \mathcal{X} \times \mathcal{A} \rightarrow \Delta([0, 1])$ , and a discount factor  $\gamma \in [0, 1)$ . A trajectory is generated as  $z_1 \sim \Gamma_1$ ,  $s_1 \sim P(\cdot|z_1)$ ,  $r_1 \sim R(s_1, a_1)$ ,  $z_2 \sim \Gamma(z_1, a_1)$ ,  $s_2 \sim P(\cdot|z_2)$ ,  $\dots$ , where the hidden states  $z_h$ 's are not observable to the agent. Moreover, the  $Q^*$  function of this POMDP is assumed to only depend on the last observation  $s_h$ , hence ‘‘reactive’’ POMDPs. We make a further simplification by assuming that the observations are indeed Markov (which implies reactive  $Q^*$ ).*

**Proposition 10** (Formal version of Proposition 5). *Let the environment be a reactive POMDP as defined above, where the underlying hidden state space  $\mathcal{Z}$  is finite. The (Markov) observation space  $\mathcal{S}$  is finite but can be arbitrarily large. Assume that the number of admissible distributions is finite,<sup>20</sup> there exists a distribution  $\mu_{\mathcal{S}} \in \Delta(\mathcal{S})$  that can be expressed as a mixture of admissible distributions (more accurately, their marginals over states), such that  $C \leq |\mathcal{Z} \times \mathcal{A}|$  when  $\mu := \mu_{\mathcal{S}} \times \text{Unif}(\mathcal{A})$  is used as the data distribution (recall the definition of  $C$  in Assumption 1).*

*Proof.* The proof contains two parts: the first part shows that a certain matrix consisting of admissible distributions has low rank, and the second part exploits the low-rankness to construct the mixture distribution described in the proposition statement and shows that it yields low concentratability coefficient  $C$ .

By definition, an admissible (state-action) distribution takes the form of  $\eta_h^\pi \in \Delta(\mathcal{S} \times \mathcal{A})$ , that is the distribution over state-action pairs induced by rolling into time step  $h$  with policy  $\pi$ . Let  $\nu_h^\pi(s)$  denote the corresponding marginal probability over states, which we call an *admissible state distribution*. Note that  $\eta_h^\pi(s, a) = \nu_h^\pi(s)\pi(a|s)$ .

Let there be a total of  $N$  admissible state distributions (we assumed  $N$  to be finite). Order them in an arbitrary manner and let the  $i$ -th admissible state distribution be  $\nu_{h_i}^{\pi_i}$ , for  $i = 1, \dots, N$ . Stacking these distributions as a matrix:

$$A_{\mathcal{S}} := \begin{bmatrix} \nu_{h_1}^{\pi_1}(s^1) & \dots & \nu_{h_1}^{\pi_1}(s^{|\mathcal{S}|}) \\ \vdots & \ddots & \vdots \\ \nu_{h_N}^{\pi_N}(s^1) & \dots & \nu_{h_N}^{\pi_N}(s^{|\mathcal{S}|}) \end{bmatrix},$$

where each row is indexed by an admissible state distribution and each column is indexed by a state, and  $\mathcal{S} := \{s^1, \dots, s^{|\mathcal{S}|}\}$ .

<sup>20</sup>This assumption is only introduced to get around of some technical subtleties, and the resulting upper bound on  $C$  has no dependence on the number of admissible distributions.

In reactive POMDPs, we can also define admissible distributions over hidden states  $\mathcal{Z}$ . For any  $z \in \mathcal{Z}$  and  $a \in \mathcal{A}$ , with abuse of notation, we use  $\eta_h^\pi(z, a)$  and  $\nu_h^\pi(z)$  to denote the distribution over hidden states (and actions) at step  $h$  induced by  $\pi$ . For any  $\pi$ ,  $h$ , and  $s$ , the distribution over observations can be decomposed as  $\nu_h^\pi(s) = \sum_{z \in \mathcal{Z}} P(s|z) \nu_h^\pi(z)$ , where  $P(s|z)$  is the emission process and is independent of the policy or the timestep. Therefore, we have

$$A_{\mathcal{S}} = \begin{bmatrix} \nu_{h_1}^{\pi_1}(z^1) & \cdots & \nu_{h_1}^{\pi_1}(z^{|\mathcal{Z}|}) \\ \vdots & \ddots & \vdots \\ \nu_{h_N}^{\pi_N}(z^1) & \cdots & \nu_{h_N}^{\pi_N}(z^{|\mathcal{Z}|}) \end{bmatrix} \begin{bmatrix} P(s^1|z^1) & \cdots & P(s^{|\mathcal{S}|}|z^1) \\ \vdots & \ddots & \vdots \\ P(s^1|z^{|\mathcal{Z}|}) & \cdots & P(s^{|\mathcal{S}|}|z^{|\mathcal{Z}|}) \end{bmatrix} := A_{\mathcal{Z}} P_{\mathcal{S}|\mathcal{Z}}.$$

From the above, we conclude that  $r := \text{rank}(A_{\mathcal{Z}}) \leq |\mathcal{Z}|$ .

In the rest of the proof we describe how to construct the mixture distribution and show that it yields low concentratability coefficient  $C$ . First, we factorize  $A_{\mathcal{Z}}$  as the product of two matrices with full column rank and full row rank, respectively:

$$A_{\mathcal{Z}} = B_{\mathcal{Z}} C_{\mathcal{Z}}.$$

We know that  $\text{rank}(B_{\mathcal{Z}}) = \text{rank}(A_{\mathcal{Z}}) = r \leq |\mathcal{Z}|$ .

Now let's focus on  $B_{\mathcal{Z}} := [b_1 \ \cdots \ b_M]^\top$ , where  $b_i^\top$  is its  $i$ -th row. Let  $D_{\mathcal{Z}}$  consists of  $r$  rows from  $B_{\mathcal{Z}}$  that maximize the absolute value of determinant (i.e., the spanned volume). That is

$$D_{\mathcal{Z}} := \begin{bmatrix} b_{i_1}^\top \\ \vdots \\ b_{i_r}^\top \end{bmatrix}, \quad \text{where } (i_1, \dots, i_r) := \arg \max_{i'_1, \dots, i'_r \in \{1, \dots, N\}} \left| \det \begin{bmatrix} b_{i'_1}^\top \\ \vdots \\ b_{i'_r}^\top \end{bmatrix} \right|.$$

Since  $D_{\mathcal{Z}}$  maximizes the absolute value of the determinant,  $|\det D_{\mathcal{Z}}| > 0$  and  $D_{\mathcal{Z}}$  is a full-rank square matrix. As a result, any row  $b_i^\top$  of  $B_{\mathcal{Z}}$  is a linear combination of rows in  $D_{\mathcal{Z}}$ . So there exists  $\alpha_1, \dots, \alpha_r \in \mathbb{R}$ , such that  $b_i = \sum_{j=1}^r \alpha_j d_j$  where  $d_j := b_{i_j}$ . We claim that  $|\alpha_j| \leq 1$  always holds.

This can be proved by contradiction. Assume that  $|\alpha_{j_0}| > 1$ , then consider the matrix

$$E_{\mathcal{Z}} = \begin{bmatrix} d_1^\top \\ \vdots \\ d_{j_0-1}^\top \\ b_i^\top \\ d_{j_0+1}^\top \\ \vdots \\ d_r^\top \end{bmatrix} = \begin{bmatrix} 1 & \cdots & 0 & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 1 & 0 & 0 & \cdots & 0 \\ \alpha_1 & \cdots & \alpha_{j_0-1} & \alpha_{j_0} & \alpha_{j_0+1} & \cdots & \alpha_r \\ 0 & \cdots & 0 & 0 & 1 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & 0 & \cdots & 1 \end{bmatrix} D_{\mathcal{Z}} := T_{\mathcal{Z}} D_{\mathcal{Z}}.$$

This matrix essentially replaces the  $j_0$ -th row of  $D_{\mathcal{Z}}$  with  $b_i^\top$ . Since  $D_{\mathcal{Z}}$  is volume maximizing, the volume of  $E_{\mathcal{Z}}$  should not increase. Calculating the determinant, however, we get  $|\det E_{\mathcal{Z}}| = |\det T_{\mathcal{Z}} \det D_{\mathcal{Z}}| = |\alpha_{j_0}| |\det D_{\mathcal{Z}}| > |\det D_{\mathcal{Z}}|$ , which causes a contradiction.

Finally, we construct the data distribution as a mixture of admissible distributions. Let  $\mu(s) = \frac{1}{r} \sum_{j=1}^r \nu_{h_{i_j}}^{\pi_{i_j}}(s)$  and  $\mu(a|s) = 1/|\mathcal{A}|$ . It is easy to check that  $\mu(s, a)$  is a valid distribution. Then for any  $i \in \{1, \dots, N\}$ ,

$$\frac{\nu_{h_i}^{\pi_i}(s, a)}{\mu(s, a)} = \frac{\nu_{h_i}^{\pi_i}(s)}{\mu(s)} \frac{\nu(a|s)}{\mu(a|s)} \leq \frac{\nu_{h_i}^{\pi_i}(s)}{\mu(s)} \frac{1}{1/|\mathcal{A}|}.$$

Now recall that for any  $i$ , there exists  $|\alpha_j| \leq 1, j = 1, \dots, r$ , such that  $b_i = \sum_{j=1}^r \alpha_j b_{i_j}$ . Since  $A_{\mathcal{S}} = B_{\mathcal{Z}} C_{\mathcal{Z}} P_{\mathcal{S}|\mathcal{Z}}$ , comparing the  $i$ -th row of both sides, we have

$$\nu_{h_i}^{\pi_i}(s) = \sum_{j=1}^r \alpha_j b_{i_j}^\top C_{\mathcal{Z}} P_{\mathcal{S}|\mathcal{Z}} = \sum_{j=1}^r \alpha_j \nu_{h_{i_j}}^{\pi_{i_j}}(s) \leq \sum_{j=1}^r |\alpha_j| \nu_{h_{i_j}}^{\pi_{i_j}}(s).$$

The inequality follows from the non-negativity of probabilities. Hence,

$$\frac{\nu_{h_i}^{\pi_i}(s, a)}{\mu(s, a)} \leq \frac{\sum_{j=1}^r |\alpha_j| \nu_{h_{i_j}}^{\pi_{i_j}}(s)}{\frac{1}{r} \sum_{j=1}^r \nu_{h_{i_j}}^{\pi_{i_j}}(s)} |\mathcal{A}| \leq \frac{\sum_{j=1}^r \nu_{h_{i_j}}^{\pi_{i_j}}(s)}{\frac{1}{r} \sum_{j=1}^r \nu_{h_{i_j}}^{\pi_{i_j}}(s)} |\mathcal{A}| = r |\mathcal{A}| \leq |\mathcal{Z}| \times |\mathcal{A}|. \quad \square$$

### C. Analysis of FQI

We state the more general error bound for FQI when Assumption 3 only holds approximately; Theorem 2 is a direct corollary of this result. Note that although our bound contains a slow-rate term ( $n^{-1/4}$ ), it is multiplied by  $\sqrt[4]{\epsilon_{\mathcal{F}, \mathcal{F}}}$  and becomes small when  $\epsilon_{\mathcal{F}, \mathcal{F}}$  is small. Furthermore, a closer examination of the bound reveals that the slow-rate term is *always* a geometric mean of the fast-rate term and the approximation error term, so the slow-rate term never dominates the bound. The bound for the minimax algorithm (Theorem 17) is in a similar situation, which distinguishes our bound from prior results for this algorithm that contains a “real” and dominating slow-rate term (Antos et al., 2008; Munos & Szepesvári, 2008).

**Theorem 11** (Error bound for FQI). *Given a dataset  $D = \{(s, a, r, s')\}$  with sample size  $|D| = n$ ,  $\mathcal{F}$  that satisfies approximate completeness (Assumption 3) with error  $\epsilon_{\mathcal{F}, \mathcal{F}}$ , with probability at least  $1 - \delta$ , the output policy of FQI after  $k$  iterations,  $\pi_{f_k}$ , satisfies<sup>21</sup>*

$$v^* - v^{\pi_{f_k}} \leq O \left( \frac{V_{\max}}{(1-\gamma)^2} \left( \sqrt{\frac{C \ln \frac{|\mathcal{F}|}{\delta}}{n}} + \sqrt[4]{\frac{C \ln \frac{|\mathcal{F}|}{\delta}}{n} \epsilon_{\mathcal{F}, \mathcal{F}}} \right) \right) + \frac{2(\sqrt{C \epsilon_{\mathcal{F}, \mathcal{F}}} + \gamma^k (1-\gamma) V_{\max})}{(1-\gamma)^2}.$$

To prove the theorem, we first define some useful notations for the proof and prove a few helper lemmas. Some of these notations/lemmas will also be helpful for the later analysis of the minimax algorithm, and we will reuse them.

**Additional Notations** We use  $\eta_h^\pi \times \pi'$  to denote the joint distribution over  $(s, a)$ , where  $s \sim \eta_h^\pi$  and  $a \sim \pi'(s)$ . For any  $\nu \in \Delta(\mathcal{S} \times \mathcal{A})$ , define  $P(\nu)$  as a distribution over states such that  $s' \sim P(\nu) \Leftrightarrow (s, a) \sim \nu, s' \sim P(s, a)$ .

The first lemma is the direct consequence of concentratability (recall Assumption 1).

**Lemma 12.** *Let  $\mu$  be any admissible distribution.  $\|\cdot\|_{2, \nu} \leq \sqrt{C} \|\cdot\|_{2, \mu}$ .*

*Proof.* For any function  $g : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , we have

$$\begin{aligned} \|g\|_{2, \nu} &= \left( \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} |g(s, a)|^2 \nu(s, a) \right)^{1/2} \\ &\leq \left( \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} |g(s, a)|^2 C \mu(s, a) \right)^{1/2} \\ &= \sqrt{C} \left( \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} |g(s, a)|^2 \mu(s, a) \right)^{1/2} = \sqrt{C} \|g\|_{2, \mu}. \quad \square \end{aligned}$$

The next lemma relates the suboptimality of a policy greedy w.r.t. a function  $f$  to  $\|f - Q^*\|$ .

**Lemma 13.** *Let  $f : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  and  $\hat{\pi} = \pi_f$  be the policy of interest, we have*

$$v^* - v^{\hat{\pi}} \leq \sum_{h=1}^{\infty} \gamma^{h-1} \left( \|Q^* - f\|_{2, \eta_h^{\hat{\pi}} \times \pi^*} + \|Q^* - f\|_{2, \eta_h^{\hat{\pi}} \times \hat{\pi}} \right).$$

<sup>21</sup>Big-Oh notations in this paper only suppress absolute constants.

*Proof.*

$$\begin{aligned}
 v^* - v^{\hat{\pi}} &= \sum_{h=1}^{\infty} \gamma^{h-1} \mathbb{E}_{s \sim \eta_h^{\hat{\pi}}} [V^*(s) - Q^*(s, \hat{\pi})] && \text{(see e.g., Kakade \& Langford (2002, Lemma 6.1))} \\
 &\leq \sum_{h=1}^{\infty} \gamma^{h-1} \mathbb{E}_{s \sim \eta_h^{\hat{\pi}}} [Q^*(s, \pi^*) - f(s, \pi^*) + f(s, \hat{\pi}) - Q^*(s, \hat{\pi})] \\
 &\leq \sum_{h=1}^{\infty} \gamma^{h-1} \left( \|Q^* - f\|_{1, \eta_h^{\hat{\pi}} \times \pi^*} + \|Q^* - f\|_{1, \eta_h^{\hat{\pi}} \times \hat{\pi}} \right) \\
 &\leq \sum_{h=1}^{\infty} \gamma^{h-1} \left( \|Q^* - f\|_{2, \eta_h^{\hat{\pi}} \times \pi^*} + \|Q^* - f\|_{2, \eta_h^{\hat{\pi}} \times \hat{\pi}} \right). \quad \square
 \end{aligned}$$

The following lemma, vaguely speaking, shows that max operator is a non-expansion in the function approximation setting.

**Lemma 14.** *Assume  $f, f' : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  and define  $\pi_{f, f'}(s) := \arg \max_{a \in \mathcal{A}} \max\{f(s, a), f'(s, a)\}$ . Then we have  $\forall \nu \in \Delta(\mathcal{S} \times \mathcal{A})$ ,*

$$\|V_f - V_{f'}\|_{2, P(\nu)} \leq \|f - f'\|_{2, P(\nu) \times \pi_{f, f'}}.$$

*Proof.*

$$\begin{aligned}
 \|V_f - V_{f'}\|_{2, P(\nu)}^2 &= \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} \sum_{s' \in \mathcal{S}} P(s' | s, a) (\max_{a \in \mathcal{A}} f(s', a) - \max_{a' \in \mathcal{A}} f'(s', a'))^2 \\
 &\leq \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} \sum_{s' \in \mathcal{S}} P(s' | s, a) (f(s', \pi_{f, f'}) - f'(s', \pi_{f, f'}))^2 = \|f - f'\|_{2, P(\nu) \times \pi_{f, f'}}^2. \quad \square
 \end{aligned}$$

With the help of Lemma 14, we are able to upper bound  $\|f - Q^*\|$  using the Bellman error  $\|f - \mathcal{T}f\|$  under  $\ell_2$  norm. The more coarse-grained version w.r.t.  $\ell_\infty$  norm has been proved by Singh & Yee (1994).

**Lemma 15.** *For an exploratory distribution  $\mu \in \Delta(\mathcal{S} \times \mathcal{A})$ , any distribution  $\nu \in \Delta(\mathcal{S} \times \mathcal{A})$ , policy  $\pi$ , and  $f, f' : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , we have*

$$\|f - Q^*\|_{2, \nu} \leq \sqrt{C} \|f - \mathcal{T}f'\|_{2, \mu} + \gamma \|f' - Q^*\|_{2, P(\nu) \times \pi_{f', Q^*}}$$

and

$$\|f - Q^*\|_{2, \nu} \leq \frac{\sqrt{C}}{1 - \gamma} \|f - \mathcal{T}f\|_{2, \mu}.$$

*Proof.* For any fixed distribution  $\nu$ , we have

$$\begin{aligned}
 \|f - Q^*\|_{2, \nu} &= \|f - \mathcal{T}f' + \mathcal{T}f' - Q^*\|_{2, \nu} \\
 &\leq \|f - \mathcal{T}f'\|_{2, \nu} + \|\mathcal{T}f' - \mathcal{T}Q^*\|_{2, \nu} \\
 &\leq \sqrt{C} \|f - \mathcal{T}f'\|_{2, \mu} + \gamma \|V_{f'} - V^*\|_{2, P(\nu)} \quad (*) \\
 &\leq \sqrt{C} \|f - \mathcal{T}f'\|_{2, \mu} + \gamma \|f' - Q^*\|_{2, P(\nu) \times \pi_{f', Q^*}}. \quad \text{(Lemma 14)}
 \end{aligned}$$

Step (\*) holds because:

$$\begin{aligned}
 \|\mathcal{T}f' - \mathcal{T}Q^*\|_{2, \nu}^2 &= \mathbb{E}_{(s, a) \sim \nu} \left[ ((\mathcal{T}f')(s, a) - (\mathcal{T}Q^*)(s, a))^2 \right] \\
 &= \mathbb{E}_{(s, a) \sim \nu} \left[ \left( \gamma \mathbb{E}_{s' \sim P(s, a)} [V_{f'}(s') - V^*(s')] \right)^2 \right] \\
 &\leq \gamma^2 \mathbb{E}_{(s, a) \sim \nu, s' \sim P(s, a)} \left[ (V_{f'}(s') - V^*(s'))^2 \right] \quad \text{(Jensen)} \\
 &= \gamma^2 \mathbb{E}_{s' \sim P(\nu)} \left[ (V_{f'}(s') - V^*(s'))^2 \right] \\
 &= \gamma^2 \|V_{f'} - V^*\|_{2, P(\nu)}^2.
 \end{aligned}$$

For the second term, let  $f' = f$  and  $\nu_0 = \arg \max_{\nu} \|f - Q^*\|_{2,\nu}$ , then we have

$$\begin{aligned} \|f - Q^*\|_{2,\nu_0} &\leq \sqrt{C} \|f - \mathcal{T}f\|_{2,\mu} + \gamma \|f - Q^*\|_{2,P(\nu_0) \times \pi_{f,Q^*}} \\ &\leq \sqrt{C} \|f - \mathcal{T}f\|_{2,\mu} + \gamma \|f - Q^*\|_{2,\nu_0} \end{aligned}$$

Therefore,  $\|f - Q^*\|_{2,\nu} \leq \|f - Q^*\|_{2,\nu_0} \leq \frac{\sqrt{C}}{1-\gamma} \|f - \mathcal{T}f\|_{2,\mu}$ .  $\square$

Finally, a concentration result that yields fast rate when completeness holds.

**Lemma 16.** *Given the MDP  $M = (\mathcal{S}, \mathcal{A}, P, R, \gamma, \eta_1)$ , we assume that the  $Q$ -function classes  $\mathcal{F}$  and  $\mathcal{G}$  are finite but can be exponentially large.  $\mathcal{G}$  approximately realizes  $\mathcal{T}\mathcal{F}$  ( $\forall f \in \mathcal{F}$ , let  $g_f^* = \arg \min_{g \in \mathcal{G}} \|g - \mathcal{T}f\|_{2,\mu}$ , then  $\|g_f^* - \mathcal{T}f\|_{2,\mu}^2 \leq \epsilon_{\mathcal{F},\mathcal{G}}$ ). The dataset  $D$  is generated from  $M$  as follows:  $(s, a) \sim \mu$ ,  $r = R(s, a)$ ,  $s' \sim P(s, a)$ . We have that  $\forall f \in \mathcal{F}$ , with probability at least  $1 - \delta$ ,*

$$\mathcal{L}_\mu(\widehat{\mathcal{T}}_{\mathcal{G}}f; f) - \mathcal{L}_\mu(g_f^*; f) \leq \frac{56V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{3n} + \sqrt{\frac{32V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{n} \epsilon_{\mathcal{F},\mathcal{G}}}$$

*Proof.* Fix  $f \in \mathcal{F}$  and  $g \in \mathcal{G}$ , define

$$X(g, f, g_f^*) := (g(s, a) - r - \gamma V_f(s'))^2 - (g_f^*(s, a) - r - \gamma V_f(s'))^2.$$

Plugging each  $(s, a, r, s') \in D$  into  $X(g, f, g_f^*)$ , we get i.i.d. variables  $X_1(g, f, g_f^*), X_2(g, f, g_f^*), \dots, X_n(g, f, g_f^*)$ . It is easy to see that

$$\frac{1}{n} \sum_{i=1}^n X_i(g, f, g_f^*) = \mathcal{L}_D(g; f) - \mathcal{L}_D(g_f^*; f).$$

Then we bound variance of  $X$ :

$$\begin{aligned} \mathbb{V}[X(g, f, g_f^*)] &\leq \mathbb{E}[X(g, f, g_f^*)^2] \\ &= \mathbb{E} \left[ \left( (g(s, a) - r - \gamma V_f(s'))^2 - (g_f^*(s, a) - r - \gamma V_f(s'))^2 \right)^2 \right] \\ &= \mathbb{E} \left[ (g(s, a) - g_f^*(s, a))^2 (g(s, a) + g_f^*(s, a) - 2r - 2\gamma V_f(s'))^2 \right] \\ &\leq 4V_{\max}^2 \mathbb{E} \left[ (g(s, a) - g_f^*(s, a))^2 \right] \\ &= 4V_{\max}^2 \|g - g_f^*\|_{2,\mu}^2 \tag{7} \\ &\leq 8V_{\max}^2 (\mathbb{E}[X(g, f, g_f^*)] + 2\epsilon_{\mathcal{F},\mathcal{G}}). \tag{*} \end{aligned}$$

Step (\*) holds because

$$\begin{aligned} &\|g - g_f^*\|_{2,\mu}^2 \\ &\leq 2 (\|g - \mathcal{T}f\|_{2,\mu}^2 + \|\mathcal{T}f - g_f^*\|_{2,\mu}^2) \tag{((a+b)^2 \leq 2a^2 + 2b^2)} \\ &\leq 2 (\|g - \mathcal{T}f\|_{2,\mu}^2 - \|\mathcal{T}f - g_f^*\|_{2,\mu}^2 + 2\|\mathcal{T}f - g_f^*\|_{2,\mu}^2) \\ &= 2 [(\mathcal{L}_\mu(g; f) - \mathcal{L}_\mu(\mathcal{T}f; f)) - (\mathcal{L}_\mu(g_f^*; f) - \mathcal{L}_\mu(\mathcal{T}f; f)) + 2\|\mathcal{T}f - g_f^*\|_{2,\mu}^2] \\ &= 2 (\mathbb{E}[X(g, f, g_f^*)] + 2\|\mathcal{T}f - g_f^*\|_{2,\mu}^2) \\ &\leq 2(\mathbb{E}[X(g, f, g_f^*)] + 2\epsilon_{\mathcal{F},\mathcal{G}}) \end{aligned}$$

Next, we apply (one-sided) Bernstein's inequality and union bound over all  $f \in \mathcal{F}$  and  $g \in \mathcal{G}$ . With probability at least

$1 - \delta$ , we have

$$\begin{aligned}
 & \mathbb{E}[X(g, f, g_f^*)] - \frac{1}{n} \sum_{i=1}^n X_i(f, f, g_f^*) \\
 & \leq \sqrt{\frac{2\mathbb{V}[X(g, f, g_f^*)] \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{n}} + \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{3n} \\
 & = \sqrt{\frac{16V_{\max}^2 \left( \mathbb{E}[X(g, f, g_f^*)] + 2\epsilon_{\mathcal{F}, \mathcal{G}} \right) \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{n}} + \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{3n}. \tag{8}
 \end{aligned}$$

Since  $\widehat{\mathcal{T}}_{\mathcal{G}} f$  minimizes  $\mathcal{L}_D(\cdot; f)$ , it also minimizes  $\frac{1}{n} \sum_{i=1}^n X_i(\cdot, f, g_f^*)$ . This is because the two objectives only differ by a constant  $\mathcal{L}_D(g_f^*; f)$ . Hence,

$$\frac{1}{n} \sum_{i=1}^n X_i(\widehat{\mathcal{T}}_{\mathcal{G}} f, f, g_f^*) \leq \frac{1}{n} \sum_{i=1}^n X_i(g_f^*, f, g_f^*) = 0.$$

Then,

$$\mathbb{E}[X(\widehat{\mathcal{T}}_{\mathcal{G}} f, f, g_f^*)] \leq \sqrt{\frac{16V_{\max}^2 \left( \mathbb{E}[X(\widehat{\mathcal{T}}_{\mathcal{G}} f, f, g_f^*)] + 2\epsilon_{\mathcal{F}, \mathcal{G}} \right) \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{n}} + \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{3n}.$$

Solving for the quadratic formula,

$$\begin{aligned}
 \mathbb{E}[X(\widehat{\mathcal{T}}_{\mathcal{G}} f, f, g_f^*)] & \leq \sqrt{48 \left( \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{3n} \right)^2 + \frac{32V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{n} \epsilon_{\mathcal{F}, \mathcal{G}} + \frac{28V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{3n}} \\
 & \leq \frac{(28 + 16\sqrt{3})V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{3n} + \sqrt{\frac{32V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{n} \epsilon_{\mathcal{F}, \mathcal{G}}} \quad (\sqrt{a+b} \leq \sqrt{a} + \sqrt{b} \text{ and } \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta} > 0) \\
 & \leq \frac{56V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{3n} + \sqrt{\frac{32V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{n} \epsilon_{\mathcal{F}, \mathcal{G}}}
 \end{aligned}$$

Noticing that  $\mathbb{E}[X(\widehat{\mathcal{T}}_{\mathcal{G}} f, f, g_f^*)] = \mathcal{L}_{\mu}(\widehat{\mathcal{T}}_{\mathcal{G}} f; f) - \mathcal{L}_{\mu}(g_f^*; f)$ , we complete the proof.  $\square$

Now we are ready to prove the main theorem.

*Proof of Theorem 11.* Firstly, we can let  $f = f_k$  and  $f' = f_{k-1}$  in Lemma 15. This gives us that

$$\|f_k - Q^*\|_{2, \nu} \leq \sqrt{C} \|f_k - \mathcal{T}f_{k-1}\|_{2, \mu} + \gamma \|f_{k-1} - Q^*\|_{2, P(\nu) \times \pi_{f_{k-1}, Q^*}}.$$

Note that we can apply the same analysis on  $P(\nu) \times \pi_{f_{k-1}, Q^*}$  and expand the inequality  $k$  times. It then suffices to upper bound  $\|f_k - \mathcal{T}f_{k-1}\|_{2, \mu}$ .

$$\begin{aligned}
 & \|f_k - \mathcal{T}f_{k-1}\|_{2, \mu}^2 \\
 & = \mathcal{L}_{\mu}(f_k; f_{k-1}) - \mathcal{L}_{\mu}(\mathcal{T}f_{k-1}; f_{k-1}) \quad (\mathcal{L} \text{ squared loss} + \mathcal{T}f_{k-1} \text{ Bayes optimal}) \\
 & = [\mathcal{L}_{\mu}(f_k; f_{k-1}) - \mathcal{L}_{\mu}(g_{f_{k-1}}^*; f_{k-1})] + [\mathcal{L}_{\mu}(g_{f_{k-1}}^*; f_{k-1}) - \mathcal{L}_{\mu}(\mathcal{T}f_{k-1}; f_{k-1})] \\
 & \leq \epsilon_1 + \|g_{f_{k-1}}^* - \mathcal{T}f_{k-1}\|_{2, \mu}^2 \quad (\text{Let } \mathcal{G} = \mathcal{F} \text{ in Lemma 16} + \mathcal{L} \text{ squared loss} + \mathcal{T}f_{k-1} \text{ Bayes optimal}) \\
 & \leq \epsilon_1 + \epsilon_{\mathcal{F}, \mathcal{F}}. \quad (\text{The selection of } g_{f_{k-1}}^*)
 \end{aligned}$$

The inequality holds with probability at least  $1 - \delta$  and  $\epsilon_1 = \frac{56V_{\max}^2 \ln \frac{|\mathcal{F}|^2}{\delta}}{3n} + \sqrt{\frac{32V_{\max}^2 \ln \frac{|\mathcal{F}|^2}{\delta}}{n} \epsilon_{\mathcal{F}, \mathcal{F}}}$ .



Noticing that  $\epsilon_1$  and  $\epsilon_{\mathcal{F},\mathcal{F}}$  do not depend on  $k$ , and the inequality holds simultaneously for different  $k$ , we have that

$$\|f_k - Q^*\|_{2,\nu} \leq \frac{1 - \gamma^k}{1 - \gamma} \sqrt{C(\epsilon_1 + \epsilon_{\mathcal{F},\mathcal{F}})} + \gamma^k V_{\max}.$$

Applying this to Lemma 13, we have that

$$\begin{aligned} v^* - v^{\pi_{f_k}} &\leq \frac{2}{1 - \gamma} \left( \frac{1 - \gamma^k}{1 - \gamma} \sqrt{C(\epsilon_1 + \epsilon_{\mathcal{F},\mathcal{F}})} + \gamma^k V_{\max} \right) \\ &\leq \frac{2}{(1 - \gamma)^2} \left( \sqrt{C\epsilon_1} + \sqrt{C\epsilon_{\mathcal{F},\mathcal{F}}} + \gamma^k (1 - \gamma) V_{\max} \right) \\ &\leq \frac{2}{(1 - \gamma)^2} \left( \sqrt{\frac{56CV_{\max}^2 \ln \frac{|\mathcal{F}|^2}{\delta}}{3n}} + \sqrt{\frac{32CV_{\max}^2 \ln \frac{|\mathcal{F}|^2}{\delta}}{n}} \epsilon_{\mathcal{F},\mathcal{F}} + \sqrt{C\epsilon_{\mathcal{F},\mathcal{F}}} + \gamma^k (1 - \gamma) V_{\max} \right). \end{aligned}$$

The proof is completed by simplifying the expression.  $\square$

## D. Analysis of the Minimax Algorithm

We state the more general error bound for the minimax algorithm when Assumptions 2 and 3 only hold approximately; Theorem 3 is a direct corollary of this result. See Appendix C for the interpretations and discussions of this result.

**Theorem 17** (Error bound for the minimax algorithm). *Given a dataset  $D = \{(s, a, r, s')\}$  with sample size  $|D| = n$ ,  $\mathcal{F}$  that satisfies approximate realizability with error  $\epsilon_{\mathcal{F}}$ , and  $\mathcal{G}$  that satisfies approximate completeness with error  $\epsilon_{\mathcal{F},\mathcal{G}}$ , with probability at least  $1 - \delta$ , the output policy of the minimax algorithm (Eq.(6)),  $\pi_{\hat{f}}$ , satisfies:*

$$v^* - v^{\pi_{\hat{f}}} \leq O \left( \frac{V_{\max} \sqrt{C}}{(1 - \gamma)^2} \left( \sqrt{\frac{\ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{n}} + \sqrt[4]{\frac{\ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta}}{n}} (\epsilon_{\mathcal{F}} + \epsilon_{\mathcal{F},\mathcal{G}}) \right) \right) + \frac{2\sqrt{2C}}{(1 - \gamma)^2} (\sqrt{\epsilon_{\mathcal{F}}} + \sqrt{2\epsilon_{\mathcal{F},\mathcal{G}}}).$$

We provide a sketched outline before diving into the detailed proof:

1. The objective in the minimax form is

$$\inf_{f \in \mathcal{F}} \sup_{g \in \mathcal{G}} (\mathcal{L}_D(f; f) - \mathcal{L}_D(g; f)) = \inf_{f \in \mathcal{F}} (\mathcal{L}_D(f; f) - \mathcal{L}_D(\hat{\mathcal{T}}_G f; f)).$$

2. We begin with dropping the dependence on function class  $\mathcal{G}$  by upper bounding the difference  $|\mathcal{L}_D(\hat{\mathcal{T}}_G f; f) - \mathcal{L}_D(\mathcal{T}f; f)|$ . This is Lemma 18 and can be separated into two substeps.

The first substep is to bound  $|\frac{1}{n} \sum_{i=1}^n X_i(\hat{\mathcal{T}}_G f, f, g_f^*)|$ , where  $X(\hat{\mathcal{T}}_G f, f, g_f^*) = (\hat{\mathcal{T}}_G f(s, a) - r - \gamma V_f(s'))^2 - (g_f^*(s, a) - r - \gamma V_f(s'))^2$ . This error is between the output of the algorithm and the best function in class  $\mathcal{G}$ .

The second substep is to bound  $|\frac{1}{n} \sum_{i=1}^n Y_i(g_f^*, f)|$ , where  $Y(g_f^*, f) := (g_f^*(s, a) - r - \gamma V_f(s'))^2 - ((\mathcal{T}f)(s, a) - r - \gamma V_f(s'))^2$ . This error is between the best function in class  $\mathcal{G}$  and the true Bellman update  $\mathcal{T}f$ .

In this way, we can change the objective in the minimax form to  $\inf_{f \in \mathcal{F}} (\mathcal{L}_D(f; f) - \mathcal{L}_D(\mathcal{T}f; f))$ , within a bounded error.

3. Then, we only need to consider the function class  $\mathcal{F}$ , since  $\inf_{f \in \mathcal{F}} (\mathcal{L}_D(f; f) - \mathcal{L}_D(\mathcal{T}f; f))$  is only related to  $\mathcal{F}$ . The proof can be finished by the following three substeps.

Firstly, by the optimality of  $\hat{f}$  and the previous error bounds, we can bound the difference between  $\frac{1}{n} \sum_{i=1}^n Z_i(\hat{f})$  and  $\frac{1}{n} \sum_{i=1}^n Z_i(f^*)$  by  $\epsilon_2$  in Lemma 18, where  $Z(f) = (f(s, a) - r - \gamma V_f(s'))^2 - ((\mathcal{T}f)(s, a) - r - \gamma V_f(s'))^2$  and  $\frac{1}{n} \sum_{i=1}^n Z_i(f) = \mathcal{L}_D(f; f) - \mathcal{L}_D(\mathcal{T}f; f)$ .

Secondly, by the property of  $f^*$ , we can bound  $\frac{1}{n} \sum_{i=1}^n Z_i(f^*)$  by  $\epsilon_3$  in Lemma 18.

These two substeps give us the bound of  $\frac{1}{n} \sum_{i=1}^n Z_i(\hat{f})$ .

Thirdly, applying Lemma 15 and Lemma 13, which is the similar steps in FQI, we obtain the desired result.

We start proving Theorem 17 by a concentration result.

**Lemma 18.** *Under the same assumption as Lemma 16, we have that  $\forall f \in \mathcal{F}$ , with probability at least  $1 - \delta$ ,*

$$\left| \mathcal{L}_D(\widehat{\mathcal{T}}_{\mathcal{G}} f; f) - \mathcal{L}_D(\mathcal{T} f; f) \right| \leq \frac{43V_{\max}^2 \ln \frac{4|\mathcal{F}||\mathcal{G}|}{\delta}}{n} + \sqrt{\frac{239V_{\max}^2 \ln \frac{4|\mathcal{F}||\mathcal{G}|}{\delta}}{n} \epsilon_{\mathcal{F},\mathcal{G}}} + \epsilon_{\mathcal{F},\mathcal{G}}.$$

*Proof.* We first apply (two-sided) Bernstein's inequality and union bound over all  $f \in \mathcal{F}$  and  $g \in \mathcal{G}$  (similar to Eq.(8) in Lemma 16). Define  $\delta' := \delta/4$ . With probability at least  $1 - 2\delta'$ , we have

$$\left| \frac{1}{n} \sum_{i=1}^n X_i(\widehat{\mathcal{T}}_{\mathcal{G}}, f, g_f^*) - \mathbb{E}[X(\widehat{\mathcal{T}}_{\mathcal{G}}, f, g_f^*)] \right| \leq \sqrt{\frac{16V_{\max}^2 \left( \mathbb{E}[X(\widehat{\mathcal{T}}_{\mathcal{G}}, f, g_f^*)] + 2\epsilon_{\mathcal{F},\mathcal{G}} \right) \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n}} + \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{3n},$$

which means that

$$\left| \frac{1}{n} \sum_{i=1}^n X_i(\widehat{\mathcal{T}}_{\mathcal{G}}, f, g_f^*) \right| \leq \left| \mathbb{E}[X(\widehat{\mathcal{T}}_{\mathcal{G}}, f, g_f^*)] \right| + \sqrt{\frac{16V_{\max}^2 \left( \mathbb{E}[X(\widehat{\mathcal{T}}_{\mathcal{G}}, f, g_f^*)] + 2\epsilon_{\mathcal{F},\mathcal{G}} \right) \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n}} + \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{3n}.$$

Noticing that  $\mathbb{E}[X(\widehat{\mathcal{T}}_{\mathcal{G}} f, f, g_f^*)] = \mathcal{L}_{\mu}(\widehat{\mathcal{T}}_{\mathcal{G}} f; f) - \mathcal{L}_{\mu}(g_f^*; f) = [\mathcal{L}_{\mu}(\widehat{\mathcal{T}}_{\mathcal{G}} f; f) - \mathcal{L}_{\mu}(\mathcal{T} f; f)] + [\mathcal{L}_{\mu}(\mathcal{T} f; f) - \mathcal{L}_{\mu}(g_f^*; f)] = \|\widehat{\mathcal{T}}_{\mathcal{G}} f - \mathcal{T} f\|_{2,\mu}^2 - \|g_f^* - \mathcal{T} f\|_{2,\mu}^2 \geq 0$ , and the results in Lemma 16 also holds (one-sided Bernstein's inequality is implied by the two-sided Bernstein's inequality), we have

$$0 \leq \mathcal{L}_{\mu}(\widehat{\mathcal{T}}_{\mathcal{G}} f; f) - \mathcal{L}_{\mu}(g_f^*; f) \leq \frac{56V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{3n} + \sqrt{\frac{32V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n} \epsilon_{\mathcal{F},\mathcal{G}}}.$$

Therefore, we have

$$\left| \mathbb{E}[X(\widehat{\mathcal{T}}_{\mathcal{G}} f, f, g_f^*)] \right| \leq \frac{56V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{3n} + \sqrt{\frac{32V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n} \epsilon_{\mathcal{F},\mathcal{G}}}.$$

Substituting this inequality into the bound of  $\left| \frac{1}{n} \sum_{i=1}^n X_i(\widehat{\mathcal{T}}_{\mathcal{G}}, f, g_f^*) \right|$ , we have that with probability at least  $1 - 2\delta'$ ,

$$\begin{aligned} & \left| \frac{1}{n} \sum_{i=1}^n X_i(\widehat{\mathcal{T}}_{\mathcal{G}}, f, g_f^*) \right| \\ & \leq \frac{56V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{3n} + \sqrt{\frac{32V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n} \epsilon_{\mathcal{F},\mathcal{G}}} \\ & \quad + \sqrt{\frac{16V_{\max}^2 \left( \frac{56V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{3n} + \sqrt{\frac{32V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n} \epsilon_{\mathcal{F},\mathcal{G}}} + 2\epsilon_{\mathcal{F},\mathcal{G}} \right) \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n}} + \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{3n} \quad (*) \\ & \leq \frac{60V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{3n} + \sqrt{\frac{32V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n} \epsilon_{\mathcal{F},\mathcal{G}}} + \sqrt{\frac{16V_{\max}^2 \left( \frac{80V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{3n} + 3\epsilon_{\mathcal{F},\mathcal{G}} \right) \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n}} \\ & \leq \frac{122V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{3n} + \sqrt{\frac{159V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n} \epsilon_{\mathcal{F},\mathcal{G}}} \quad (\sqrt{a+b} \leq \sqrt{a} + \sqrt{b} \text{ and } \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'} > 0) \end{aligned}$$

In Step (\*), we use  $\sqrt{\frac{32V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n} \epsilon_{\mathcal{F},\mathcal{G}}} \leq \frac{8V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n} + \epsilon_{\mathcal{F},\mathcal{G}}$ .

Then, define

$$Y(g, f) := (g(s, a) - r - \gamma V_f(s'))^2 - ((\mathcal{T}f)(s, a) - r - \gamma V_f(s'))^2.$$

Plugging each  $(s, a, r, s') \in D$  into  $Y_1(g_f^*, f)$ , we get i.i.d. variables  $Y_2(g_f^*, f), Y(g_f^*, f), \dots, Y_n(g_f^*, f)$ . Applying same derivations in Lemma 16, we can get similar bound as Inequality (7),

$$0 \leq \mathbb{V}[Y(g_f^*, f)] \leq 4V_{\max}^2 \|g_f^* - \mathcal{T}f\|_{2, \mu}^2 (= 4V_{\max}^2 \mathbb{E}[Y(g_f^*, f)]) \leq 4V_{\max}^2 \epsilon_{\mathcal{F}, \mathcal{G}}.$$

We can apply (two-sided) Bernstein's inequality and union bound over all  $f \in \mathcal{F}$  and  $g \in \mathcal{G}$ . With probability at least  $1 - 2\delta'$ , we have

$$\left| \frac{1}{n} \sum_{i=1}^n Y_i(g_f^*, f) - \mathbb{E}[Y(g_f^*, f)] \right| \leq \sqrt{\frac{8V_{\max}^2 \mathbb{E}[Y(g_f^*, f)] \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n}} + \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{3n},$$

which means that

$$\left| \frac{1}{n} \sum_{i=1}^n Y_i(g_f^*, f) \right| \leq \epsilon_{\mathcal{F}, \mathcal{G}} + \sqrt{\frac{8V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n}} \epsilon_{\mathcal{F}, \mathcal{G}} + \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{3n}.$$

Union bounding the results of  $\frac{1}{n} \sum_{i=1}^n X_i(\widehat{\mathcal{T}}_G f, f, g_f^*)$  and  $\frac{1}{n} \sum_{i=1}^n Y_i(g_f^*, f)$ , we have that with probability at least  $1 - 4\delta'$ ,

$$\begin{aligned} & \left| \mathcal{L}_D(\widehat{\mathcal{T}}_G f; f) - \mathcal{L}_D(\mathcal{T}f; f) \right| \\ &= \left| \frac{1}{n} \sum_{i=1}^n X_i(\widehat{\mathcal{T}}_G f, f, g_f^*) + \frac{1}{n} \sum_{i=1}^n Y_i(g_f^*, f) \right| \\ &\leq \frac{43V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n} + \sqrt{\frac{239V_{\max}^2 \ln \frac{|\mathcal{F}||\mathcal{G}|}{\delta'}}{n}} \epsilon_{\mathcal{F}, \mathcal{G}} + \epsilon_{\mathcal{F}, \mathcal{G}}. \end{aligned}$$

Noticing  $\delta' = \delta/4$ , we complete the proof.  $\square$

*Proof of Theorem 17.* Firstly, Lemma 15 gives us that

$$\|\hat{f} - Q^*\|_{2, \nu} \leq \frac{\sqrt{C}}{1 - \gamma} \|\hat{f} - \mathcal{T}\hat{f}\|_{2, \mu}.$$

It then suffices to upper bound  $\|\hat{f} - \mathcal{T}\hat{f}\|_{2, \mu}$ .

The objective of the minimax form minimization can be written as  $\inf_{f \in \mathcal{F}} \sup_{g \in \mathcal{G}} (\mathcal{L}_D(f; f) - \mathcal{L}_D(g; f))$ . We can find that,  $\forall f \in \mathcal{F}$ ,

$$\arg \max_{g \in \mathcal{G}} (\mathcal{L}_D(f; f) - \mathcal{L}_D(g; f)) = \arg \max_{g \in \mathcal{G}} -\mathcal{L}_D(g; f) = \arg \min_{g \in \mathcal{G}} \mathcal{L}_D(g; f) = \widehat{\mathcal{T}}_G f.$$

Define  $\delta' := \delta/2$ , Lemma 18 tells us that  $\forall f \in \mathcal{F}$ , we have that with probability at least  $1 - \delta'$ ,

$$\left| \mathcal{L}_D(\widehat{\mathcal{T}}_G f; f) - \mathcal{L}_D(\mathcal{T}f; f) \right| \leq \epsilon_2,$$

where

$$\epsilon_2 = \frac{43V_{\max}^2 \ln \frac{4|\mathcal{F}||\mathcal{G}|}{\delta'}}{n} + \sqrt{\frac{239V_{\max}^2 \ln \frac{4|\mathcal{F}||\mathcal{G}|}{\delta'}}{n}} \epsilon_{\mathcal{F}, \mathcal{G}} + \epsilon_{\mathcal{F}, \mathcal{G}}.$$

From the approximate realizability of  $\mathcal{F}$ , we know there exists  $f^* \in \mathcal{F}$ , s.t.  $\|f^* - \mathcal{T}f^*\|_{2, \mu}^2 \leq \epsilon_{\mathcal{F}}$ . Then by the optimality of  $\hat{f}$ , we have that  $\mathcal{L}_D(\hat{f}; \hat{f}) - \mathcal{L}_D(\widehat{\mathcal{T}}_G \hat{f}; \hat{f}) \leq \mathcal{L}_D(f^*; f^*) - \mathcal{L}_D(\widehat{\mathcal{T}}_G f^*; f^*)$ . Therefore, with probability at least  $1 - \delta'$ , we have that

$$\mathcal{L}_D(\hat{f}; \hat{f}) - \mathcal{L}_D(\mathcal{T}\hat{f}; \hat{f}) \leq \mathcal{L}_D(f^*; f^*) - \mathcal{L}_D(\mathcal{T}f^*; f^*) + 2\epsilon_2.$$

Define

$$Z(f) := (f(s, a) - r - \gamma V_f(s'))^2 - ((\mathcal{T}f)(s, a) - r - \gamma V_f(s'))^2.$$

Plugging each  $(s, a, r, s') \in D$  into  $Z(f)$ , we get i.i.d. variables  $Z_1(f), Z_2(f), \dots, Z_n(f)$ . Applying Ineq. (7) in Lemma 16, we get

$$\mathbb{V}[Z(f)] \leq 4V_{\max}^2 \|f - \mathcal{T}f\|_{2,\mu}^2 = 4V_{\max}^2 \mathbb{E}[Z(f)].$$

We can apply (one-sided) Bernstein's inequality and union bound over all  $f \in \mathcal{F}$ . With probability at least  $1 - \delta'$ , we have that  $\forall f \in \mathcal{F}$ ,

$$\frac{1}{n} \sum_{i=1}^n Z_i(f) - \mathbb{E}[Z(f)] \leq \sqrt{\frac{8V_{\max}^2 \mathbb{E}[Z(f)] \ln \frac{|\mathcal{F}|}{\delta'}}{n}} + \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta'}}{3n}.$$

Substituting  $f^*$  into the inequality and noticing  $\|f^* - \mathcal{T}f^*\|_{2,\mu}^2 \leq \epsilon_{\mathcal{F}}$ , we have

$$\frac{1}{n} \sum_{i=1}^n Z_i(f^*) \leq \epsilon_{\mathcal{F}} + \sqrt{\frac{8V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta'}}{n}} \epsilon_{\mathcal{F}} + \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta'}}{3n} := \epsilon_3.$$

Since  $\forall f \in \mathcal{F}$ ,  $\frac{1}{n} \sum_{i=1}^n Z_i(f) = \mathcal{L}_D(f; f) - \mathcal{L}_D(\mathcal{T}f; f)$ , with probability at least  $1 - 2\delta'$ , we have

$$\frac{1}{n} \sum_{i=1}^n Z_i(\hat{f}) = \mathcal{L}_D(\hat{f}; \hat{f}) - \mathcal{L}_D(\mathcal{T}\hat{f}; \hat{f}) \leq \mathcal{L}_D(f^*; f^*) - \mathcal{L}_D(\mathcal{T}f^*; f^*) + 2\epsilon_2 \leq 2\epsilon_2 + \epsilon_3.$$

Finally, we consider  $Z(\hat{f})$ . Our goal is to bound  $\|\hat{f} - \mathcal{T}\hat{f}\|_{2,\mu} = \sqrt{\mathbb{E}[Z(\hat{f})]}$ . Substituting  $\hat{f}$  into the concentration bound of  $Z(f)$ , we have

$$\mathbb{E}[Z(\hat{f})] - \frac{1}{n} \sum_{i=1}^n Z_i(\hat{f}) \leq \sqrt{\frac{8V_{\max}^2 \mathbb{E}[Z(\hat{f})] \ln \frac{|\mathcal{F}|}{\delta'}}{n}} + \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta'}}{3n}.$$

Substituting the upper bound of  $\frac{1}{n} \sum_{i=1}^n Z_i(\hat{f})$  into the equality, we have that, with probability at least  $1 - 2\delta'$ ,

$$\|\hat{f} - \mathcal{T}\hat{f}\|_{2,\mu}^2 = \mathbb{E}[Z(\hat{f})] \leq \sqrt{\frac{8V_{\max}^2 \mathbb{E}[Z(\hat{f})] \ln \frac{|\mathcal{F}|}{\delta'}}{n}} + \frac{4V_{\max}^2 \ln \frac{|\mathcal{F}|}{\delta'}}{3n} + 2\epsilon_2 + \epsilon_3.$$

Solving this quadratic formula and noticing that  $\delta = 2\delta'$ , we have that with probability at least  $1 - \delta$ ,

$$\|\hat{f} - \mathcal{T}\hat{f}\|_{2,\mu}^2 \leq \frac{16V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{3n} + 2\epsilon_2 + \epsilon_3 + \sqrt{\frac{8V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{n} \left( \frac{10V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{3n} + 2\epsilon_2 + \epsilon_3 \right)},$$

where

$$\epsilon_2 = \frac{43V_{\max}^2 \ln \frac{8|\mathcal{F}||\mathcal{G}|}{\delta}}{n} + \sqrt{\frac{239V_{\max}^2 \ln \frac{8|\mathcal{F}||\mathcal{G}|}{\delta}}{n}} \epsilon_{\mathcal{F},\mathcal{G}} + \epsilon_{\mathcal{F},\mathcal{G}},$$

and

$$\epsilon_3 = \epsilon_{\mathcal{F}} + \sqrt{\frac{8V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{n}} \epsilon_{\mathcal{F}} + \frac{4V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{3n}.$$

In this way, we obtain the bound for  $\|\hat{f} - \mathcal{T}\hat{f}\|_{2,\mu}$ , and further the bound for  $\|f - Q^*\|_{2,\mu}$  (by Lemma 15). Finally, applying the bound for  $\|f - Q^*\|_{2,\mu}$  to Lemma 13, we have that with probability at least  $1 - \delta$ ,

$$\begin{aligned}
 v^* - v^{\pi_f} &\leq \frac{2\sqrt{C}}{(1-\gamma)^2} \sqrt{\frac{16V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{3n} + 2\epsilon_2 + \epsilon_3} + \sqrt{\frac{8V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{n} \left( \frac{10V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{3n} + 2\epsilon_2 + \epsilon_3 \right)} \\
 &\leq \frac{2\sqrt{C}}{(1-\gamma)^2} \sqrt{\frac{16V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{3n} + 2\epsilon_2 + \epsilon_3 + \frac{2V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{n} + \left( \frac{10V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{3n} + 2\epsilon_2 + \epsilon_3 \right)} \\
 &= \frac{2\sqrt{C}}{(1-\gamma)^2} \sqrt{\frac{32V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{3n} + 4\epsilon_2 + 2\epsilon_3} \\
 &\leq \frac{2\sqrt{C}}{(1-\gamma)^2} \left( \sqrt{\frac{32V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{3n}} + \sqrt{4\epsilon_2} + \sqrt{2\epsilon_3} \right) \\
 &\leq \frac{2\sqrt{C}}{(1-\gamma)^2} \left( \sqrt{\frac{32V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{3n}} + \sqrt{\frac{172V_{\max}^2 \ln \frac{8|\mathcal{F}||\mathcal{G}|}{\delta'}}{n}} + \sqrt{\frac{3824V_{\max}^2 \ln \frac{8|\mathcal{F}||\mathcal{G}|}{\delta}}{n}} \epsilon_{\mathcal{F},\mathcal{G}} + 2\sqrt{\epsilon_{\mathcal{F},\mathcal{G}}} \right) \\
 &\quad + \frac{2\sqrt{C}}{(1-\gamma)^2} \left( \sqrt{2\epsilon_{\mathcal{F}}} + \sqrt[4]{\frac{32V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{n}} \epsilon_{\mathcal{F}} + \sqrt{\frac{8V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{3n}} \right) \\
 &\leq \frac{2\sqrt{C}}{(1-\gamma)^2} (\sqrt{2\epsilon_{\mathcal{F}}} + 2\sqrt{\epsilon_{\mathcal{F},\mathcal{G}}}) + \frac{2\sqrt{C}}{(1-\gamma)^2} \left( \sqrt{\frac{24V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{n}} + \sqrt{\frac{172V_{\max}^2 \ln \frac{8|\mathcal{F}||\mathcal{G}|}{\delta}}{n}} \right) \\
 &\quad + \frac{2\sqrt{C}}{(1-\gamma)^2} \left( \sqrt[4]{\frac{32V_{\max}^2 \ln \frac{2|\mathcal{F}|}{\delta}}{n}} \epsilon_{\mathcal{F}} + \sqrt[4]{\frac{3824V_{\max}^2 \ln \frac{8|\mathcal{F}||\mathcal{G}|}{\delta}}{n}} \epsilon_{\mathcal{F},\mathcal{G}} \right).
 \end{aligned}$$

The proof is completed by combining the terms and absorbing the constants using Big-Oh notation.  $\square$

## E. Proofs Related to State Abstractions

### E.1. Equivalence Between MBRL with State Abstractions and FQI with Piece-wise Constant Function Class

**Proposition 19.** *In model-based RL with abstraction  $\phi : \mathcal{S} \rightarrow \mathcal{S}_\phi$ , we estimate an abstract model  $\widehat{M}_\phi = (\mathcal{S}_\phi, \mathcal{A}, \widehat{P}_\phi, \widehat{R}_\phi, \gamma)$  and then perform planning. When value iteration is used as the planning algorithm, the procedure is exactly equivalent to FQI with  $\mathcal{F}^\phi$  as the function class.*

To prove the result, we first define a few notations.

**Definition 5 (Lifting).** *Given the MDP  $M = (\mathcal{S}, \mathcal{A}, P, R, \gamma)$  and the state abstraction  $\phi$  that operates on  $\mathcal{S}$ , for any function  $f$  that operates on  $\mathcal{S}_\phi \times \mathcal{A}$ , we use  $[f]_M$  to denote its lifted version, which is a function over  $\mathcal{S} \times \mathcal{A}$  and defined as  $[f]_M(s, a) := f(\phi(s), a)$ .*

*Similarly, we can also lift a state value function. For any function  $f$  that operates on  $\mathcal{S}_\phi$ , we also use  $[f]_M$  to denote its lifted version, which is a function over  $\mathcal{S}$  and defined as  $[f]_M(s) := f(\phi(s))$ . Lifting a real-valued function  $f$  over states can also be expressed in vector form:  $[f]_M = \Phi^\top f$ , where  $\Phi$  is an  $|\mathcal{S}_\phi| \times |\mathcal{S}|$  matrix with entries  $\phi(s, x) = \mathbb{I}[\phi(s) = x]$ .*

**Definition 6.** *For piece-wise constant state-action value function  $f$  and  $x \in \mathcal{S}_\phi$ , define  $[f]_\phi(x, a) = f(s, a)$  for any  $s \in \phi^{-1}(x)$ ; note that the notation  $[\cdot]_\phi$  can only be applied to functions that are piece-wise constant under  $\phi$ .*

*Proof of Proposition 19.* Let  $D = \{D_{s,a}\}_{(s,a) \in \mathcal{S} \times \mathcal{A}}$  where  $D_{s,a}$  is the collection of transition tuples that start with  $(s, a)$ . We also let  $\mathbf{e}_{\phi(s')}$  be the unit vector whose  $\phi(s')$ -th entry is 1 and all other entries are 0. Then, for any abstract state-action

pair  $(x, a) \in \mathcal{S}_\phi \times \mathcal{A}$ , the certainty-equivalence estimate of model parameters are:

$$\widehat{R}_\phi(x, a) = \frac{1}{|D_{x,a}|} \sum_{(r,s') \in D_{x,a}} r \quad \text{and} \quad \widehat{P}_\phi(x, a) = \frac{1}{|D_{x,a}|} \sum_{(r,s') \in D_{x,a}} \mathbf{e}_{\phi(s')}.$$

If we use value iteration as the planning algorithm, we will first initialize  $g_0 \in [0, R_{\max}]^{|\mathcal{S}_\phi \times \mathcal{A}|}$ . Then in each iteration, we let  $g_t = \widehat{\mathcal{T}}_{M_\phi} g_{t-1}$ . Expanding the operator  $\widehat{\mathcal{T}}_{M_\phi}$ , for  $x \in \mathcal{S}_\phi$  and  $a \in \mathcal{A}$ , we have

$$\begin{aligned} g_t(x, a) &= \widehat{R}_\phi(x, a) + \gamma \langle \widehat{P}_\phi(x, a), V_{g_{t-1}} \rangle \\ &= \frac{1}{|D_{x,a}|} \sum_{(r,s') \in D_{x,a}} (r + \gamma \langle \mathbf{e}_{\phi(s')}, V_{g_{t-1}} \rangle) \\ &= \frac{1}{|D_{x,a}|} \sum_{(r,s') \in D_{x,a}} (r + \gamma V_{g_{t-1}}(\phi(s'))) \end{aligned}$$

For the FQI with  $\mathcal{F}^\phi$ , we first initialize  $f_0$  as any function in  $\mathcal{F}^\phi = [0, R_{\max}]^{|\mathcal{S}_\phi \times \mathcal{A}|}$ . Then in each iteration, we let  $f_t = \widehat{\mathcal{T}}_{\mathcal{F}^\phi} f_{t-1}$ . From the definition of  $\widehat{\mathcal{T}}_{\mathcal{F}^\phi}$ , for  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$ , we have

$$f_t(s, a) = \arg \min_{f \in \mathcal{F}^\phi} \frac{1}{|D_{\phi(s),a}|} \sum_{(r,s') \in D_{\phi(s),a}} (f - r - \gamma V_{f_{t-1}}(s'))^2.$$

This is a regression problem and the solution is

$$\begin{aligned} f_t(s, a) &= \frac{1}{|D_{\phi(s),a}|} \sum_{(r,s') \in D_{\phi(s),a}} (r + \gamma \langle \mathbf{e}_{\phi(s')}, [V_{f_{t-1}}]_\phi \rangle) \\ &= \frac{1}{|D_{\phi(s),a}|} \sum_{(r,s') \in D_{\phi(s),a}} (r + \gamma V_{f_{t-1}}(s')) \end{aligned}$$

Therefore, if  $f_0 = [g_0]_M$ , the two algorithms give us that  $f_t = [g_t]_M$  for any  $t$ . This shows that model-based RL with abstraction  $\phi$  is exactly equivalent to FQI with  $\mathcal{F}^\phi$ .  $\square$

## E.2. Proof of Equivalence Between Bisimulation and Completeness for Piece-wise Constant Function Class

We first define approximate bisimulation, which is a generalization of Definition 2.

**Definition 7** (Approximate model-irrelevant). *Given the MDP  $M = (\mathcal{S}, \mathcal{A}, P, R, \gamma)$  and the state abstraction  $\phi : \mathcal{S} \rightarrow \mathcal{S}_\phi$ , we call  $\phi$  an  $(\epsilon_R, \epsilon_P)$ -approximate bisimulation if*

$$\max_{s_1, s_2: \phi(s_1) = \phi(s_2), a \in \mathcal{A}} |R(s_1, a) - R(s_2, a)| = \epsilon_R, \quad (9)$$

$$\max_{s_1, s_2: \phi(s_1) = \phi(s_2), a \in \mathcal{A}} \|\Phi P(s_1, a) - \Phi P(s_2, a)\|_1 = \epsilon_P, \quad (10)$$

where  $\Phi$  is as defined in Definition 5.

**Proposition 20** (Completeness=Bisimulation). *Suppose that  $\phi$  is an  $(\epsilon_R, \epsilon_P)$ -approximate  $Q^*$ -irrelevant abstraction, then we have*

$$\max \left\{ \frac{\epsilon_R}{2}, \frac{\gamma \epsilon_P V_{\max}}{4} \right\} \leq \sup_{f \in \mathcal{F}^\phi} \inf_{f' \in \mathcal{F}^\phi} \|f' - \mathcal{T}f\|_\infty \leq \frac{\epsilon_R}{2} + \frac{\gamma \epsilon_P V_{\max}}{4}.$$

Proposition 9 is a direct corollary of the above result when  $\epsilon_R, \epsilon_P, \sup_{f \in \mathcal{F}^\phi} \inf_{f' \in \mathcal{F}^\phi} \|f' - \mathcal{T}f\|_\infty$  are all 0's. In the approximate case, however, we use  $\epsilon_R$  and  $\epsilon_P$  to provide both upper and lower bounds of the violation of completeness, but do not obtain an equality relationship. This is purely an artifact that bisimulation considers rewards and transitions separately, whereas completeness considers both of them together in terms of the Bellman update operator  $\mathcal{T}$ , and cancellation between reward/transition errors may occur.

*Proof of Proposition 20.* We first prove the upper bound. For any fixed  $f \in \mathcal{F}^\phi$ , we show that there exists  $f'_1 \in \mathcal{F}^\phi$  such that  $\|f'_1 - \mathcal{T}f\|_\infty \leq \epsilon_R/2 + \gamma\epsilon_P V_{\max}/4$ . Therefore  $\inf_{f' \in \mathcal{F}^\phi} \|f' - \mathcal{T}f\|_\infty \leq \|f'_1 - \mathcal{T}f\|_\infty$  and hence is subject to the same upper bound.

Since  $f'_1$  is required to be piece-wise constant, it suffices to specify  $[f'_1]_\phi(x, a)$  for each  $x \in \mathcal{S}_\phi, a \in \mathcal{A}$ . Fixing any  $x, a$ , define

$$s_+ := \arg \max_{s \in \phi^{-1}(x), a \in \mathcal{A}} (\mathcal{T}f)(s, a), \quad s_- := \arg \min_{s \in \phi^{-1}(x), a \in \mathcal{A}} (\mathcal{T}f)(s, a), \quad (11)$$

and

$$[f'_1]_\phi(x, a) := \frac{1}{2} ((\mathcal{T}f)(s_+, a) + (\mathcal{T}f)(s_-, a)). \quad (12)$$

Note that  $f'_1 \in [0, V_{\max}]$  so  $f'_1 \in \mathcal{F}^\phi$ . It remains to upper bound  $\|f'_1 - \mathcal{T}f\|_\infty$ .

For any  $s \in \mathcal{S}, a \in \mathcal{A}$ , let  $x = \phi(s)$ , and  $s_+$  and  $s_-$  as defined in Eq.(11) for  $(x, a)$ ,

$$\begin{aligned} & f'_1(s, a) - (\mathcal{T}f)(s, a) \\ & \leq \frac{1}{2} ((\mathcal{T}f)(s_+, a) + (\mathcal{T}f)(s_-, a)) - (\mathcal{T}f)(s, a) && \text{(Eq.(11) and (12))} \\ & = \frac{1}{2} ((\mathcal{T}f)(s_+, a) - (\mathcal{T}f)(s, a)) \\ & = \frac{1}{2} (R(s_+, a) + \gamma \langle P(s_+, a), V_f \rangle - R(s, a) - \gamma \langle P(s, a), V_f \rangle) \\ & \leq \frac{1}{2} |R(s_+, a) - R(s, a)| + \frac{\gamma}{2} |\langle P(s_+, a) - P(s, a), V_f \rangle| \\ & \leq \frac{1}{2} \epsilon_R + \frac{\gamma}{2} |\langle \Phi P(s_+, a) - \Phi P(s, a), [V_f]_\phi \rangle| && (f \text{ is piece-wise constant and so is } V_f) \\ & = \frac{1}{2} \epsilon_R + \frac{\gamma}{2} |\langle \Phi P(s_+, a) - \Phi P(s, a), [V_f]_\phi - \frac{V_{\max}}{2} \cdot \mathbf{1} \rangle| && (*) \\ & \leq \frac{1}{2} \epsilon_R + \frac{\gamma}{2} \|\Phi P(s_+, a) - \Phi P(s, a)\|_1 \cdot V_{\max}/2 && \text{(H\"older's inequality)} \\ & \leq \frac{\epsilon_R}{2} + \frac{\gamma\epsilon_P V_{\max}}{4}. \end{aligned}$$

Here Step (\*) holds because  $\langle \Phi P(s_+, a) - \Phi P(s, a), \mathbf{1} \rangle = 0$ , as  $\Phi P(\cdot, \cdot)$  is always a valid distribution. The other direction follows exactly the same argument due to symmetry and is omitted, and together we conclude that  $\|f'_1 - \mathcal{T}f\|_\infty \leq \frac{\epsilon_R}{2} + \frac{\gamma\epsilon_P V_{\max}}{4}$ .

We then turn to the lower bound of the theorem statement. It suffices to show  $\exists f_R, f_P \in \mathcal{F}^\phi$ , such that  $\inf_{f' \in \mathcal{F}^\phi} \|f' - \mathcal{T}f_R\|_\infty \geq \epsilon_R/2$  and  $\inf_{f' \in \mathcal{F}^\phi} \|f' - \mathcal{T}f_P\|_\infty \geq \gamma\epsilon_P V_{\max}/4$ , respectively.

**Case of  $f_R$**  Let  $f_R := \mathbf{0} \in \mathcal{F}^\phi$ , so  $\inf_{f' \in \mathcal{F}^\phi} \|f' - \mathcal{T}f_R\|_\infty = \inf_{f' \in \mathcal{F}^\phi} \|f' - R\|_\infty$ , where  $R \in [0, R_{\max}]^{|\mathcal{S} \times \mathcal{A}|}$  is the reward function. It is obvious from the definition of  $\epsilon_R$  in Eq.(9) that  $\inf_{f' \in \mathcal{F}^\phi} \|f' - R\|_\infty = \epsilon_R/2$ , which proves the result.

**Case of  $f_P$**  Let  $s_1^P, s_2^P \in \mathcal{S}, a^P \in \mathcal{A}$  be the arguments that achieve the maximum in Eq.(10), i.e.,  $\phi(s_1^P) = \phi(s_2^P)$  and  $\|\Phi P(s_1^P, a^P) - \Phi P(s_2^P, a^P)\|_1 = \epsilon_P$ . We construct  $f_P \in \mathcal{F}^\phi$  as follows: Assume w.l.o.g. that  $R(s_1^P, a^P) \geq R(s_2^P, a^P)$ . For any  $x \in \mathcal{S}_\phi$  and  $a \in \mathcal{A}$ , define

$$[f_P]_\phi(x, a) = \mathbb{I}[P(x|s_1^P, a^P) > P(x|s_2^P, a^P)] \cdot V_{\max}.$$

Note that the RHS has no dependence on  $a$ , so  $V_f(s) = [f_P]_\phi(\phi(s), a)$  for any  $a \in \mathcal{A}$ . It is easy to verify that  $f_P \in \mathcal{F}^\phi$  as its value is either 0 or  $V_{\max}$ . Essentially  $f_P$  is designed such that  $[V_{f_P}]_\phi$  witnesses the  $\ell_1$  error (or total variation) between  $\Phi P(s_1^P, a^P)$  and  $\Phi P(s_2^P, a^P)$ . The inequality sign inside the indicator could be either " $>$ " or " $<$ ", and we choose it in consistence with the relationship between  $R(s_1^P, a^P)$  and  $R(s_2^P, a^P)$ , which guarantees that the reward error and the

transition error wouldn't cancel out with each other. Now consider the difference between two entries in  $(\mathcal{T}f_P)$ :

$$\begin{aligned}
 & (\mathcal{T}f_P)(s_1^P, a^P) - (\mathcal{T}f_P)(s_2^P, a^P) \\
 &= R(s_1^P, a^P) + \gamma \langle P(s_1^P, a^P), V_{f_P} \rangle - R(s_2^P, a^P) - \gamma \langle P(s_2^P, a^P), V_{f_P} \rangle \\
 &= (R(s_1^P, a^P) - R(s_2^P, a^P)) + \gamma \langle P(s_1^P, a^P) - P(s_2^P, a^P), V_{f_P} \rangle \quad (\text{Both terms are positive due to construction}) \\
 &= |R(s_1^P, a^P) - R(s_2^P, a^P)| + \gamma |\langle \Phi P(s_1^P, a^P) - \Phi P(s_2^P, a^P), [V_{f_P}]_\phi \rangle| \\
 &\geq 0 + \gamma \|\Phi P(s_1^P, a^P) - \Phi P(s_2^P, a^P)\|_{TV} V_{\max} \\
 &= \gamma \epsilon_P V_{\max} / 2. \tag{*}
 \end{aligned}$$

Step (\*) follows because  $[V_{f_P}]_\phi$  takes  $V_{\max}$  on the subset of  $\mathcal{S}_\phi$  where  $\Phi P(s_1^P, a^P)$  has a greater probability than  $\Phi P(s_2^P, a^P)$ , and 0 otherwise, so the dot product is equal to the total variation up to the scaling factor of  $V_{\max}$ .

Now  $\sup_{f \in \mathcal{F}^\phi} \inf_{f' \in \mathcal{F}^\phi} \|f' - \mathcal{T}f\|_\infty \geq \inf_{f' \in \mathcal{F}^\phi} \|f' - \mathcal{T}f_P\|_\infty$ , which is the approximation error of  $\mathcal{T}f_P$  in  $\mathcal{F}^\phi$ . Since  $\mathcal{T}f_P$  takes values that are  $\gamma \epsilon_P V_{\max} / 2$  apart for aggregated states  $s_1^P$  and  $s_2^P$  on action  $a^P$ , the approximation error is at least  $\gamma \epsilon_P V_{\max} / 4$ . This completes the proof.  $\square$

## F. Proof of Proposition 7

Firstly we introduce a standard result that bounds the loss of acting greedily with respect to an approximate Q-value function.

**Lemma 21.** (Singh & Yee, 1994) *For any  $f : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , let  $\pi_f$  be its greedy policy, then*

$$\|V^* - V^{\pi_f}\|_\infty \leq \frac{2\|f - Q^*\|_\infty}{1 - \gamma}.$$

Now we are ready to prove the proposition.

*Proof of Proposition 7.* We represent each state  $s$  by its value profile  $\{f(s, a) : f \in \mathcal{F}, a \in \mathcal{A}\}$ , estimate a tabular model, and output its optimal policy. Since a set of states may share exactly the same value profile, we are essentially using a state abstraction, denoted as  $\phi$ . For any two states  $s_1, s_2 \in \mathcal{S}$  that share the same value profile, the realizability assumption implies that  $Q^*(s_1, a) = Q^*(s_2, a), \forall a \in \mathcal{A}$ , so  $\phi$  is  $Q^*$ -irrelevant (Li et al., 2006). In the following, we will show that certainty-equivalence with  $Q^*$ -irrelevant abstraction is consistent and enjoys polynomial sample complexity if each state-action pair  $(s, a)$  receives  $\Omega(|D|/|\mathcal{S} \times \mathcal{A}|)$  data.

Let  $D_{s,a}$  be the collection of transition tuples that start with  $(s, a)$  and  $D_{x,a} := \sum_{s \in \phi^{-1}(x)} |D_{s,a}|$ . We first consider an abstract MDP  $M_\phi = (\mathcal{S}_\phi, \mathcal{A}, P_\phi, R_\phi, \gamma)$ , where  $\mathcal{S}_\phi$  is the abstract state space (isomorphic to the set of distinct value profiles),

$$R_\phi(x, a) = \frac{\sum_{s \in \phi^{-1}(x)} |D_{s,a}| R(s, a)}{|D_{\phi(s), a}|}, \quad P_\phi(x' | x, a) = \frac{\sum_{s \in \phi^{-1}(x)} |D_{s,a}| P(x' | s, a)}{|D_{\phi(s), a}|}, \quad \forall x, x' \in \mathcal{S}_\phi, a \in \mathcal{A}.$$

Recall the notations in Definitions 5, 6, and 7. We claim that  $[Q_{M_\phi}^*]_M = Q_M^*$ , where  $[Q_{M_\phi}^*]_M$  is the lifted version of  $Q_{M_\phi}^*$ : Since  $Q_M^*(s, a)$  is piece-wise constant under  $\phi$ , we let  $[Q_M^*]_\phi(x, a) = Q_M^*(s, a)$  for any  $s \in \phi^{-1}(x)$ . It suffices to show  $Q_{M_\phi}^* = [Q_M^*]_\phi$ , by showing that  $[Q_M^*]_\phi$  is the fixed point of  $\mathcal{T}_{M_\phi}$ . This is because, for any  $x \in \mathcal{S}_\phi, a \in \mathcal{A}$ ,

$$\begin{aligned}
 (\mathcal{T}_{M_\phi}[Q_M^*]_\phi)(x, a) &= R_\phi(x, a) + \gamma \langle P_\phi(x, a), [V_M^*]_\phi \rangle \\
 &= \sum_{s \in \phi^{-1}(x)} \frac{|D_{s,a}|}{|D_{\phi(s), a}|} (R(s, a) + \gamma \langle \Phi P(s, a), [V_M^*]_\phi \rangle) \\
 &= \sum_{s \in \phi^{-1}(x)} \frac{|D_{s,a}|}{|D_{\phi(s), a}|} (R(s, a) + \gamma \langle P(s, a), V_M^* \rangle) \\
 &= \sum_{s \in \phi^{-1}(x)} \frac{|D_{s,a}|}{|D_{\phi(s), a}|} [Q_M^*]_\phi(x, a) = [Q_M^*]_\phi(x, a).
 \end{aligned}$$



Then we consider the estimated model using the abstract representation  $\widehat{M}_\phi = (\mathcal{S}_\phi, \mathcal{A}, \widehat{P}_\phi, \widehat{R}_\phi, \gamma)$ . Let  $\mathbf{e}_{\phi(s')}$  be the unit vector whose  $\phi(s')$ -th entry is 1 and all other entries are 0, the parameters are

$$\widehat{R}_\phi(x, a) = \frac{1}{|D_{x,a}|} \sum_{(r,s') \in D_{x,a}} r \quad \text{and} \quad \widehat{P}_\phi(x, a) = \frac{1}{|D_{x,a}|} \sum_{(r,s') \in D_{x,a}} \mathbf{e}_{\phi(s')}, \quad \forall (x, a) \in \mathcal{S}_\phi \times \mathcal{A}.$$

Define the minimal samples received by any abstract state-action pair as  $n_\phi(D) := \min_{x \in \mathcal{S}_\phi, a \in \mathcal{A}} |D_{x,a}|$ . We can upper bound  $\left\| [Q_{M_\phi}^*]_M - [Q_{\widehat{M}_\phi}^*]_M \right\|_\infty$  by a function of  $n_\phi(D)$ .

Noticing the contraction property of  $\mathcal{T}_{\widehat{M}_\phi}$ , we have

$$\begin{aligned} \left\| [Q_{M_\phi}^*]_M - [Q_{\widehat{M}_\phi}^*]_M \right\|_\infty &= \left\| Q_{M_\phi}^* - Q_{\widehat{M}_\phi}^* \right\|_\infty \\ &\leq \frac{1}{1-\gamma} \left\| Q_{M_\phi}^* - \mathcal{T}_{\widehat{M}_\phi} Q_{M_\phi}^* \right\|_\infty \\ &= \frac{1}{1-\gamma} \left\| \mathcal{T}_{\widehat{M}_\phi} Q_{M_\phi}^* - \mathcal{T}_{M_\phi} Q_{M_\phi}^* \right\|_\infty. \end{aligned} \quad (*)$$

Step(\*) holds because

$$\left\| Q_{M_\phi}^* - Q_{\widehat{M}_\phi}^* \right\|_\infty \leq \left\| Q_{M_\phi}^* - \mathcal{T}_{\widehat{M}_\phi} Q_{M_\phi}^* \right\|_\infty + \left\| \mathcal{T}_{\widehat{M}_\phi} Q_{M_\phi}^* - \mathcal{T}_{\widehat{M}_\phi} Q_{\widehat{M}_\phi}^* \right\|_\infty \leq \left\| Q_{M_\phi}^* - \mathcal{T}_{\widehat{M}_\phi} Q_{M_\phi}^* \right\|_\infty + \gamma \left\| Q_{M_\phi}^* - Q_{\widehat{M}_\phi}^* \right\|_\infty.$$

Then we plug in the definition of  $\mathcal{T}_{\widehat{M}_\phi}$  and  $\mathcal{T}_{M_\phi}$ . For each  $(x, a) \in \mathcal{S}_\phi \times \mathcal{A}$ ,

$$\begin{aligned} &|(\mathcal{T}_{\widehat{M}_\phi} Q_{M_\phi}^*)(x, a) - (\mathcal{T}_{M_\phi} Q_{M_\phi}^*)(x, a)| \\ &= |\widehat{R}_\phi(x, a) + \gamma \langle \widehat{P}_\phi(x, a), V_{M_\phi}^* \rangle - R_\phi(x, a) - \gamma \langle P_\phi(x, a), V_{M_\phi}^* \rangle| \\ &= \left| \frac{1}{|D_{x,a}|} \sum_{s \in \phi^{-1}(x)} \sum_{(r,s') \in D_{s,a}} \left( r + \gamma V_{M_\phi}^*(\phi(s')) - R(s, a) - \gamma \langle P(s, a), [V_{M_\phi}^*]_M \rangle \right) \right|. \end{aligned}$$

If we view the nested sum as a flat sum, the expression is the sum of the differences between random variables  $r + \gamma V_{M_\phi}^*(s')$  and their expectation w.r.t. the randomness of  $(r, s')$ . Each sample is independent and bounded in  $[0, V_{\max}]$ , so Hoeffding's inequality applies: with probability at least  $1 - \delta/|\mathcal{S}_\phi \times \mathcal{A}|$ ,

$$\left| (\mathcal{T}_{\widehat{M}_\phi} Q_{M_\phi}^*)(x, a) - (\mathcal{T}_{M_\phi} Q_{M_\phi}^*)(x, a) \right| \leq V_{\max} \sqrt{\frac{1}{2n_\phi(D)} \ln \frac{2|\mathcal{S}_\phi \times \mathcal{A}|}{\delta}}.$$

Union bounding over all  $(x, a) \in \mathcal{S}_\phi \times \mathcal{A}$ , with probability at least  $1 - \delta$ , we get

$$\left\| \mathcal{T}_{\widehat{M}_\phi} Q_{M_\phi}^* - \mathcal{T}_{M_\phi} Q_{M_\phi}^* \right\|_\infty \leq V_{\max} \sqrt{\frac{1}{2n_\phi(D)} \ln \frac{2|\mathcal{S}_\phi \times \mathcal{A}|}{\delta}}.$$

Therefore, with probability at least  $1 - \delta$ , we have

$$\begin{aligned} \left\| Q_M^* - [Q_{\widehat{M}_\phi}^*]_M \right\|_\infty &\leq \left\| Q_M^* - [Q_{M_\phi}^*]_M \right\|_\infty + \left\| [Q_{M_\phi}^*]_M - [Q_{\widehat{M}_\phi}^*]_M \right\|_\infty \\ &\leq \frac{V_{\max}}{1-\gamma} \sqrt{\frac{1}{2n_\phi(D)} \ln \frac{2|\mathcal{S}_\phi \times \mathcal{A}|}{\delta}}. \end{aligned}$$

Finally, applying Lemma 21 with  $f = [Q_{\widehat{M}_\phi}^*]_M$ , we get

$$v_M^* - v_M^{\pi_{\widehat{M}_\phi}^*} \leq \left\| V_M^* - V_M^{\pi_{\widehat{M}_\phi}^*} \right\|_\infty = \left\| V_M^* - V_M^{\pi_{[Q_{\widehat{M}_\phi}^*]_M}^*} \right\|_\infty \leq \frac{2 \left\| [Q_{\widehat{M}_\phi}^*]_M - Q_M^* \right\|_\infty}{1-\gamma}.$$

This means that with probability at least  $1 - \delta$ , the output of certainty-equivalence with  $Q^*$ -irrelevant abstraction  $[\pi_{M\phi}^*]_M$  satisfies

$$v_M^* - v_M^{[\pi_{M\phi}^*]_M} \leq \frac{2V_{\max}}{(1-\gamma)^2} \sqrt{\frac{1}{2n_\phi(D)} \ln \frac{2|\mathcal{S}_\phi \times \mathcal{A}|}{\delta}} \leq \frac{2V_{\max}}{(1-\gamma)^2} \sqrt{\frac{1}{2n_\phi(D)} \ln \frac{2|\mathcal{S} \times \mathcal{A}|}{\delta}}.$$

We complete the proof by noticing that  $n_\phi(D) = \Omega(|D|/|\mathcal{S} \times \mathcal{A}|)$ , so to guarantee the above bound to be  $\epsilon$ , the necessary sample size  $|D|$  will be polynomial in all relevant parameters.  $\square$

## G. Possible Relaxation of Assumption 1

We illustrate the possibility of relaxing Assumption 1 using a simple example on state abstractions: when learning with abstractions, it is sufficient to have data that is relatively uniform over the abstract state space, even if some raw state receives no data. Due to the connection to FQI (Section 2.3), one would expect that with  $\mathcal{F}^\phi$  as the function class, concentratability coefficient can be upper bounded by the number of abstract states and incur no dependence on the raw state space, which is unfortunately not the case according to the current definition. It turns out that our analysis provides an easy fix to this issue: the proof of Theorem 2 (for FQI) only depends on Assumption 1 via

$$\|f - \mathcal{T}f'\|_{2,\nu} \leq \sqrt{C} \|f - \mathcal{T}f'\|_{2,\mu}, \forall f, f' \in \mathcal{F}. \quad (13)$$

And the proof of Theorem 3 (for the minimax algorithm) only depends on Assumption 1 via

$$\|f - \mathcal{T}f\|_{2,\nu} \leq \sqrt{C} \|f - \mathcal{T}f\|_{2,\mu}, \forall f \in \mathcal{F}. \quad (14)$$

If we define  $C$  through the above inequalities (which are strict relaxations of Assumption 1), Theorems 2 and 3 still hold under Eq.(13) and (14) respectively. Furthermore, when  $\mathcal{F} = \mathcal{F}^\phi$  and  $\phi$  is a bisimulation, we can easily verify that  $C$  can be upper bounded by the number of abstract state-action pairs with uniform data. One issue here is that Eq.(13) is specialized to the completeness assumption, and if we wish to work with alternative assumptions (as discussed earlier), we may need to relax the definition in a different manner. Another interesting observation is that Eq.(14) is less strict and much nicer than Eq.(13), but so far we have not been able to modify the FQI analysis to work under Eq.(14). It is unclear whether this is an artifact of proof techniques or a fundamental difference between FQI and the minimax algorithm. We leave the investigation of these issues to future work.