# Pessimistic Nonlinear Least-Squares Value Iteration for Offline Reinforcement Learning

**Anonymous Author(s)**
Affiliation
Address
`email`

## Abstract

Offline reinforcement learning, where the agent aims to learn the optimal policy based on the data collected by a behavior policy, has attracted increasing attention in recent years. While offline RL with linear function approximation has been extensively studied with optimal results achieved under certain assumptions, the theoretical understanding of offline RL with non-linear function approximation is still limited. Specifically, most existing works on offline RL with non-linear function approximation either have a poor dependency on the function class complexity or require an inefficient planning phase. In this paper, we propose an oracle-efficient algorithm PNLSVI for offline RL with non-linear function approximation. Our algorithmic design comprises three innovative components: (1) a variance-based weighted regression scheme that can be applied to a wide range of function classes, (2) a subroutine for variance estimation, and (3) a planning phase that utilizes a pessimistic value iteration approach. Our algorithm enjoys a regret bound that has a tight dependency on the function class complexity and achieves minimax optimal problem-dependent regret when specialized to linear function approximation. Our theoretical analysis introduces a new coverage assumption for nonlinear Q function, bridging the minimum-eigenvalue assumption and the uncertainty measure widely used in online nonlinear RL. To the best of our knowledge, this is the first statistically optimal algorithm for nonlinear offline RL.

## 1 Introduction

Offline reinforcement learning (RL), also known as batch RL, is a learning paradigm where an agent learns to make decisions based on a set of pre-collected data, instead of interacting with the environment in real-time like online RL. The goal of offline RL is to learn a policy that performs well in a given task, based on historical data that was collected from an unknown environment. Recent years have witnessed significant progress in developing offline RL algorithms that can leverage large amounts of data to learn effective policies. These algorithms often incorporate powerful function approximation techniques, such as deep neural networks, to generalize across large state-action spaces. They have achieved excellent performances in a wide range of domains, including the games of Go and chess (Silver et al., 2017; Schrittwieser et al., 2020), robotics (Gu et al., 2017; Levine et al., 2018), and control systems (Degrave et al., 2022).

Several studies have studied the theoretical guarantees of tabular offline RL and proved near-optimal sample complexities in this setting (Xie et al., 2021b; Shi et al., 2022; Li et al., 2022). However, these algorithms cannot handle numerous real-world applications with large state spaces. Consequently, a significant body of research has shifted its focus to offline RL with function approximation. For example, several works have analyzed the sample efficiency of offline RL with linear function approximation under different MDP models, including linear MDPs and their variants (Jin et al., 2021b; Zanette et al., 2021; Min et al., 2021; Yin et al., 2022a). To handle nonlinear function class, a recent line of research considered offline RL with general function approximation (Chen and Jiang,

2019; Xie et al., 2021a; Zhan et al., 2022). While these algorithms have sample efficiency guarantees, they often require an inefficient planning phase or have a poor dependency on the function class complexity. For example, Xie et al. (2021a) proposed an information-theoretic algorithm that requires solving an optimization problem over all potential policy and corresponding version space, which includes all functions with lower Bellman error. To overcome this limitation, Xie et al. (2021a) proposed a practical implementation, as a cost, the algorithm have a poor dependency on the function class complexity. Recently, (Yin et al., 2022b) studied the general differentiable function class and propose a computation efficient algorithm (PFQL). However, their result also have an addition dependence on the dimension $d$ of the parameter.

Therefore, a natural question arises:

*Can we design a computationally efficient algorithm that achieves the minimax optimality with respect to the complexity of nonlinear function class?*

We give an affirmative answer to the above question in this work. Our contributions are listed as follows:

- We propose a pessimism-based algorithm PNLSVI designed for nonlinear function approximation, which strictly generalizes the existing pessimism-based algorithms for both linear and differentiable function approximation (Xiong et al., 2022; Yin et al., 2022b). Our algorithm is oracle-efficient, i.e., it is computationally efficient when there exists an efficient regression oracle and bonus oracle for the function class (e.g., generalized linear function class).

- We prove a data-dependent regret bound with the widely used $D^2$-divergence in online nonlinear RL regime, which is optimal with respect to the function class complexity. Our analysis closes the gap to optimality for differentiable function approximation, which was previously an open problem (Yin et al., 2022b).

- We introduce a novel uniform coverage assumption for general function approximation that is generalized over the assumption in Yin et al. (2022b). Our assumption bridges between the minimum-eigenvalue assumption used in linear models and the generalized dimension for nonlinear function class, offering new insights into the function approximation problem in RL.

**Notation:** In this work, we use lowercase letters to denote scalars and use lower and uppercase boldface letters to denote vectors and matrices respectively. For a vector $\mathbf{x} \in \mathbb{R}^d$ and matrix $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$, we denote by $\|\mathbf{x}\|_2$ the Euclidean norm and $\|\mathbf{x}\|_{\boldsymbol{\Sigma}} = \sqrt{\mathbf{x}^\top \boldsymbol{\Sigma} \mathbf{x}}$. For two sequences $\{a_n\}$ and $\{b_n\}$, we write $a_n = O(b_n)$ if there exists an absolute constant $C$ such that $a_n \leq C b_n$, and we write $a_n = \Omega(b_n)$ if there exists an absolute constant $C$ such that $a_n \geq C b_n$. We use $\widetilde{O}(\cdot)$ and $\widetilde{\Omega}(\cdot)$ to further hide the logarithmic factors. For any $a \leq b \in \mathbb{R}$, $x \in \mathbb{R}$, let $[x]_{[a,b]}$ denote the truncate function $a \cdot \mathbb{1}(x \leq a) + x \cdot \mathbb{1}(a \leq x \leq b) + b \cdot \mathbb{1}(b \leq x)$, where $\mathbb{1}(\cdot)$ is the indicator function. For a positive integer $n$, we use $[n] = \{1, 2, .., n\}$ to denote the set of integers from 1 to $n$.

## 2  Related Work

**RL with function approximation.**  As one of the simplest function approximation classes, linear representation in RL has been extensively studied in recent years (Jiang et al., 2017; Dann et al., 2018; Yang and Wang, 2019; Jin et al., 2020; Wang et al., 2020c; Du et al., 2019; Sun et al., 2019; Zanette et al., 2020a,b; Weisz et al., 2021; Yang and Wang, 2020; Modi et al., 2020; Ayoub et al., 2020; Zhou et al., 2021; He et al., 2021). Several assumptions on the linear structure of the underlying MDPs have been made in these works, ranging from the *linear MDP* assumption (Yang and Wang, 2019; Jin et al., 2020; Hu et al., 2022; He et al., 2022; Agarwal et al., 2022) to the *low Bellman-rank* assumption (Jiang et al., 2017) and the *low inherent Bellman error* assumption (Zanette et al., 2020b). Extending the previous theoretical guarantees to more general problem classes, RL with nonlinear function classes has garnered increased attention in recent years (Wang et al., 2020b; Jin et al., 2021a; Foster et al., 2021; Du et al., 2021; Agarwal and Zhang, 2022; Agarwal et al., 2022). Various complexity measures of function classes have been studied including Bellman rank (Jiang et al., 2017), Bellman-Eluder dimension (Jin et al., 2021a), Decision-Estimation Coefficient (Foster et al., 2021) and generalized Eluder dimension (Agarwal et al., 2022). Among these works, the setting in our paper is most related to Agarwal et al. (2022) where $D^2$-divergence (Gentile et al., 2022) was introduced in RL to indicate the uncertainty of a sample with respect to a particular sample batch.

**Offline tabular RL.** There is a line of works integrating the principle of pessimism to develop statistically efficient algorithms for offline tabular RL setting (Rashidinejad et al., 2021; Yin and Wang, 2021; Xie et al., 2021b; Shi et al., 2022; Li et al., 2022). More specifically, Xie et al. (2021b) utilized the variance of transition noise and proposed a nearly optimal algorithm based on pessimism and Bernstein-type bonus. Subsequently, Li et al. (2022) proposed a model-based approach that achieves minimax-optimal sample complexity without burn-in cost for tabular MDPs. Shi et al. (2022) also contributed by proposing the first nearly minimax-optimal model-free offline RL algorithm.

**Offline RL with linear function approximation.** Jin et al. (2021b) presented the initial theoretical results on offline linear MDPs. They introduced a pessimism-principled algorithmic framework for offline RL and proposed an algorithm based on LSVI (Jin et al., 2020). Min et al. (2021) subsequently considered offline policy evaluation (OPE) in linear MDPs, assuming independence between data samples across time steps to obtain tighter confidence sets and proposed an algorithm with optimal $d$ dependence. Yin et al. (2022a) took one step further and considered the policy optimization in linear MDPs, which implicitly requires the same independence assumption. Zanette et al. (2021) proposed an actor-critic-based algorithm that establishes pessimism principle by directly perturbing the parameter vectors in a linear function approximation framework. Recently, Xiong et al. (2022) proposed a novel uncertainty decomposition technique via a reference function, which leads to a minimax-optimal sample complexity bound for offline linear MDPs without additional assumptions.

**Offline RL with general function approximation.** Chen and Jiang (2019) critically examined the assumptions underlying value-function approximation methods and established an information-theoretic lower bound. Xie et al. (2021a) introduced the concept of Bellman-consistent pessimism, which enables sample-efficient guarantees by relying solely on the Bellman-completeness assumption. Uehara and Sun (2021) focused on model-based offline RL with function approximation under partial coverage, demonstrating that realizability in the function class and partial coverage are sufficient for policy learning. Zhan et al. (2022) proposed an algorithm that achieves polynomial sample complexity under the realizability and single-policy concentrability assumptions. Nguyen-Tang and Arora (2023) proposed a method of random perturbations and pessimism for neural function approximation. For differentiable function classes, Yin et al. (2022b) made advancements by improving the sample complexity with respect to the stage $H$. However, their result had an additional dependence on the dimension $d$ of the parameter space, whereas in linear function approximation, the dependence is typically on $\sqrt{d}$.

## 3  Preliminaries

In our work, we consider the inhomogeneous episodic Markov Decision Processes (MDP), which can be denoted by a tuple of $\mathcal{M}\big(\mathcal{S}, \mathcal{A}, H, \{r_h\}_{h=1}^H, \{\mathbb{P}_h\}_{h=1}^H\big)$. In specific, $\mathcal{S}$ is the state space, $\mathcal{A}$ is the finite action space, $H$ is the length of each episode. For each stage $h \in [H]$, $r_h : \mathcal{S} \times \mathcal{A} \to [0, 1]$ is the reward function[1] and $\mathbb{P}_h(s'|s, a)$ is the transition probability function, which denotes the probability for state $s$ to transfer to next state $s'$ with current action $a$. A policy $\pi := \{\pi_h\}_{h=1}^H$ is a collection of mappings $\pi_h$ from a state $s \in \mathcal{S}$ to the simplex of action space $\mathcal{A}$. For simplicity, we denote the state-action pair as $z := (s, a)$. For any policy $\pi$ and stage $h \in [H]$, we define the value function $V_h^\pi(s)$ and the action-value function $Q_h^\pi(s, a)$ as the expected cumulative rewards starting at stage $h$, which can be denoted as follows:

$$Q_h^\pi(s, a) = r_h(s, a) + \mathbb{E}\bigg[ \sum_{h'=h+1}^H r_{h'}\big(s_{h'}, \pi_{h'}(s_{h'})\big) \big| s_h = s, a_h = a \bigg], \; V_h^\pi(s) = Q_h^\pi\big(s, \pi_h(s)\big),$$

where $s_{h'+1} \sim \mathbb{P}_h(\cdot|s_{h'}, a_{h'})$ denotes the observed state at stage $h' + 1$. By this definition, the value function $V_h^\pi(s)$ and action-value function $Q_h^\pi(s, a)$ are bounded in $[0, H]$. In addition, we define the optimal value function $V_h^*$ and the optimal action-value function $Q_h^*$ as $V_h^*(s) = \max_\pi V_h^\pi(s)$ and $Q_h^*(s, a) = \max_\pi Q_h^\pi(s, a)$. We denote the corresponding optimal policy by $\pi^*$. For any function $V : \mathcal{S} \to \mathbb{R}$, we denote $[\mathbb{P}_h V](s, a) = \mathbb{E}_{s' \sim \mathbb{P}_h(\cdot|s,a)} V(s')$ and $[\mathrm{Var}_h V](s, a) = [\mathbb{P}_h V^2](s, a) - \big([\mathbb{P}_h V](s, a)\big)^2$ for simplicity. For any function $f : \mathcal{S} \to \mathbb{R}$, we define the Bellman operator $\mathcal{T}_h$ as $\mathcal{T}_h f(s_h, a_h) = \mathbb{E}_{s_{h+1} \sim \mathbb{P}_h(\cdot|s_h, a_h)} [r_h(s_h, a_h) + f(s_{h+1})]$, where we use the shorthand $f(s) = \max_{a \in \mathcal{A}} f(s, a)$ for simplicity. Based on this definition, for every stage $h \in [H]$ and policy $\pi$, we

---

[1]While we study the deterministic reward functions for simplicity, it is not difficult to generalize our results to stochastic reward functions.

have the following Bellman equation for value functions $Q_h^\pi(s, a)$ and $V_h^\pi(s)$, as well as the Bellman optimality equation for optimal value functions:

$$Q_h^\pi(s_h, a_h) = \mathcal{T}_h V_{h+1}^\pi(s_h, a_h), \ Q_h^*(s_h, a_h) = \mathcal{T}_h V_{h+1}^*(s_h, a_h),$$

where $V_{H+1}^\pi(s) = V_{H+1}^*(s) = 0$. We also define the Bellman operator for second moment as $\mathcal{T}_{2,h}f(s_h, a_h) = \mathbb{E}_{s_{h+1} \sim \mathbb{P}_h(\cdot|s_h, a_h)}\left[\left(r_h(s_h, a_h) + f(s_{h+1})\right)^2\right]$. For simplicity, we omit the subscripts $h$ in the Bellman operator without causing confusion.

**Offline Reinforcement Learning:** In offline RL, the agent only have access to a batch-dataset $D = \{s_h^k, a_h^k, r_h^k : h \in [H], k \in [K]\}$, which is collected by a behavior policy $\mu$, and the agent cannot interact with the environment. Given the batch dataset, the goal of offline RL is finding a near-optimal policy $\pi$ that minimize the sub-optimality $V_1^*(s) - V_1^\pi(s)$. In addition, for each stage $h$ and behavior policy $\mu$, we denote the induced distribution of the state-action pair as $d_h^\mu$.

**General Function Approximation:** In this work, we focus on a special class of episodic MDPs, where the value function satisfies the following completeness assumption.

**Assumption 3.1** ($\epsilon$-completeness under general function approximation, Agarwal et al. 2022). Given a general function class $\{\mathcal{F}_h\}_{h \in [H]}$, where each function class $\mathcal{F}_h$ is composed of functions $f_h : \mathcal{S} \times \mathcal{A} \to [0, L]$. We assume for each stage $h \in [H]$, and any function $V : \mathcal{S} \to [0, H]$, there exists functions $f_h, f_{2,h} \in \mathcal{F}_h$ such that

$$\max_{(s,a) \in \mathcal{S} \times \mathcal{A}} |f_h(s, a) - \mathcal{T}_h V(s, a)| \leq \epsilon, \text{ and } \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} |f_{2,h}(s, a) - \mathcal{T}_{2,h}V(s, a)| \leq \epsilon.$$

In addition, for each stage $h \in [H]$, we assume there exists a function $f_h^* \in \mathcal{F}_h$ closed to the optimal value function such that $\|f_h^* - Q_h^*\|_\infty \leq \epsilon$. For simplicity, we assume $L = O(H)$ throughout the paper and denote $\mathcal{N} = \max_{h \in [H]} |\mathcal{F}_h|$.

To deal with general function class $\mathcal{F}$, Agarwal et al. (2022) introduce the following measure to capture the function class complexity for online learning.

**Definition 3.2** (Generalized Eluder dimension, Agarwal et al. 2022). Given $\lambda > 0$, a sequence of state-action pairs $Z = \{z_i\}_{i \in [K]}$ and a sequence of non-negative weights $\boldsymbol{\sigma} = \{\sigma_i\}_{i \in [K]}$. Let $\mathcal{F}$ be a function class consisting of functions $f : \mathcal{S} \times \mathcal{A} \to [0, L]$. The generalized Eluder dimension of $\mathcal{F}$ is given by $\dim_{\alpha, K}(\mathcal{F}) := \sup_{Z, \boldsymbol{\sigma}: |Z| = K, \boldsymbol{\sigma} \geq \alpha} \dim(\mathcal{F}, Z, \boldsymbol{\sigma})$, where

$$\dim(\mathcal{F}, Z, \boldsymbol{\sigma}) := \sum_{i=1}^{K} \min\left(1, \frac{1}{\sigma_i^2} D_\mathcal{F}^2(z_i; z_{[i-1]}, \sigma_{[i-1]})\right),$$

$$D_\mathcal{F}^2(z; z_{[k-1]}, \sigma_{[k-1]}) := \sup_{f_1, f_2 \in \mathcal{F}} \frac{(f_1(z) - f_2(z))^2}{\sum_{s \in [k-1]} \frac{1}{\sigma_s^2}(f_1(z_s) - f_2(z_s))^2 + \lambda}.$$

Here, the inequality $\boldsymbol{\sigma} \geq \alpha$ represents that $\sigma_i \geq \alpha$ holds for all $i \in [K]$ and we use the notation $z_{[i-1]}, \sigma_{[i-1]}$ to represent the sequences $\{z_s\}_{s=1}^{i-1}, \{\sigma_s\}_{s=1}^{i-1}$.

However, in offline RL, the proposed Generalized Eluder dimension fails to capture the relationship between function class $\mathcal{F}$ and the pre-collected dataset $\mathcal{D}$. To generalize this definition to offline environment, for a batch dataset $\mathcal{D} = \{(s_h^k, a_h^k, r_h^k)\}_{h,k=1}^{H,K}$ and a function class $\mathcal{F}_h$ consisting of functions $f : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$. We denote $\mathcal{D}_h = \{(s_h^k, a_h^k, r_h^k)\}_{k \in [K]}$ as the subset of the dataset $D$ that corresponds to the observations collected up to stage $h$ in the MDP. Then for any weight function $\sigma_h(\cdot, \cdot) : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$, we introduce the following $D^2$-divergence:

$$D_{\mathcal{F}_h}^2(z; \mathcal{D}_h; \sigma_h) = \sup_{f_1, f_2 \in \mathcal{F}_h} \frac{(f_1(z) - f_2(z))^2}{\sum_{k \in [K]} \frac{1}{(\sigma_h(z_h^k))^2}(f_1(z_h^k) - f_2(z_h^k))^2 + \lambda}.$$

**Data Coverage Assumption:** In offline RL, there exists a discrepancy between the state-action distribution generated by the behavior policy and the distribution from the learned policy. Under this situation, the distribution shift problem can cause the learned policy to perform poorly or even fail in offline RL. Therefore, we propose the following data coverage assumption to control the distribution shift.

4

---

**Algorithm 1** Pessimistic Nonlinear Least-Squares Value Iteration (PNLSVI)

---

**Require:** Input confidence parameters $\beta'_{1,h}, \beta'_{2,h}, \beta_h$ and $\epsilon > 0$.

1: **Initialize**: Split the input dataset into $\mathcal{D} = \{s_h^k, a_h^k, r_h^k\}_{k,h=1}^{K,H}, \mathcal{D}' = \{\bar{s}_h^k, \bar{a}_h^k, \bar{r}_h^k\}_{k,h=1}^{K,H}$ ; Set the value function $\widehat{f}_{H+1}(\cdot) = \widehat{f}'_{H+1}(\cdot) = 0$.

2: **for** stage $h = H, \ldots, 1$ **do**

3: $\quad \widetilde{f}'_h = \text{argmin}_{f_h \in \mathcal{F}_h} \sum_{k \in [K]} \left( f_h(\bar{s}_h^k, \bar{a}_h^k) - \bar{r}_h^k - \widehat{f}'_{h+1}(\bar{s}_{h+1}^k) \right)^2$.

4: $\quad \widetilde{g}'_h = \text{argmin}_{g_h \in \mathcal{F}_h} \sum_{k \in [K]} \left( g_h(\bar{s}_h^k, \bar{a}_h^k) - \left( \bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k) \right)^2 \right)^2$.

5: $\quad$ Use the bonus oracle (Definition 4.1) to calculate the bonus function
$\quad b'_h = \mathcal{B}(1, \mathcal{D}'_h, \mathcal{F}_h, \widetilde{f}'_h, \beta'_{1,h} + \beta'_{2,h}, \lambda, \epsilon)$,

6: $\quad \widehat{f}'_h \leftarrow \{\widetilde{f}'_h - b'_h - \epsilon\}_{[0, H-h+1]};$

7: $\quad$ Construct the variance estimator
$\quad \widehat{\sigma}_h^2(s,a) = \max \left\{ 1, \widetilde{g}'_h(s,a) - (\widetilde{f}'_h(s,a))^2 - O\left( \frac{\sqrt{\log \mathcal{N} \mathcal{N}_b H^3}}{\sqrt{K\kappa}} \right) \right\}.$

8: **end for**

9: **for** stage $h = H, \ldots, 1$ **do**

10: $\quad \widetilde{f}_h = \text{argmin}_{f_h \in \mathcal{F}_h} \sum_{k \in [K]} \frac{1}{\widehat{\sigma}_h^2(s_h^k, a_h^k)} \left( f_h(s_h^k, a_h^k) - r_h^k - \widehat{f}_{h+1}(s_{h+1}^k) \right)^2$

11: $\quad$ Use the bonus oracle (Definition 4.1) to calculate the bonus function
$\quad b_h = \mathcal{B}(\widehat{\sigma}_h, \mathcal{D}_h, \mathcal{F}_h, \widetilde{f}_h, \beta_h, \lambda, \epsilon);$

12: $\quad \widehat{f}_h \leftarrow \{\widetilde{f}_h - b_h - \epsilon\}_{[0, H-h+1]};$

13: $\quad \widehat{\pi}_h(\cdot|s) = \text{argmax}_a \widehat{f}_h(s,a).$

14: **end for**

15: **Output:** $\widehat{\pi} = \{\widehat{\pi}_h\}_{h=1}^H.$

---

**Assumption 3.3** (Uniform Data Coverage). there exists a constant $\kappa > 0$, such that for any stage $h$ and functions $f_1, f_2 \in \mathcal{F}_h$, the following inequality holds,

$$\mathbb{E}_{\mu,h} \left[ \left( f_1(s_h, a_h) - f_2(s_h, a_h) \right)^2 \right] \geq \kappa \|f_1 - f_2\|_\infty^2,$$

where the state-action pair (at stage $h$) $(s_h, a_h)$ is stochastic generated from behavior policy $\mu$.

**Remark 3.4.** Data coverage assumption is widely used in offline RL to guarantee that the collected dataset contains enough information of the state-action space to learn an effective policy. In Yin et al. (2022b), they studied the general differentiable function, where the function class is defined as

$$\mathcal{F} := \left\{ f(\boldsymbol{\theta}, \boldsymbol{\phi}(\cdot, \cdot)) : \mathcal{S} \times \mathcal{A} \to \mathbb{R}, \boldsymbol{\theta} \in \Theta \right\}.$$

Under this definition, Yin et al. (2022b) introduce the following coverage assumption (Assumption 2.3) such that for all stage $h \in [H]$, there exists a constant $\kappa$,

$$\mathbb{E}_{\mu,h} \left[ (f(\boldsymbol{\theta}_1, \boldsymbol{\phi}(s,a)) - f(\boldsymbol{\theta}_2, \boldsymbol{\phi}(s,a)))^2 \right] \geq \kappa \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2^2, \forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \Theta; \quad (*)$$

$$\mathbb{E}_{\mu,h} \left[ \nabla f(\boldsymbol{\theta}, \boldsymbol{\phi}(s,a)) \nabla f(\boldsymbol{\theta}, \boldsymbol{\phi}(s,a))^\top \right] \succ \kappa I, \forall \boldsymbol{\theta} \in \Theta. \quad (**)$$

We can prove that our assumption is weaker than the first assumption (*). For the second assumption (**), there is no direct counterpart in the general setting.

In addition, for the linear function class, the coverage assumption in Yin et al. (2022b) will reduce to the following linear function coverage assumption (Wang et al., 2020a; Min et al., 2021; Yin et al., 2022a; Xiong et al., 2022).

$$\lambda_{\min}(\mathbb{E}_{\mu,h}[\boldsymbol{\phi}(s,a)\boldsymbol{\phi}(s,a)^\top]) = \kappa > 0, \ \forall h \in [H].$$

Therefore, our assumption is also weaker than the linear function coverage assumption when dealing with the linear function class. Due to space limitations, we provide the detailed proof in the appendix.

## 4 Algorithm

In this section, we provide a comprehensive and detailed description of our algorithm (PNLSVI), as displayed in Algorithm 1. In the sequel, we introduce the key ideas of the proposed algorithm.

## 4.1 Pessimistic Value Iteration Based Planning

Our algorithm operates in two distinct phases, Variance Estimate Phase and Pessimistic Planning Phase. At the beginning of the algorithm, the data-set is divided into two disjoint subsets $\mathcal{D}, \mathcal{D}'$, and each assigned to a specific phase.

The basic framework of our algorithm follows the pessimistic value iteration, which was initially introduced by Jin et al. (2021b). In details, for each stage $h \in [H]$, we construct the estimator value function $\widetilde{f}_h$ by solving the following variance-weighted ridge regression (Line 11):

$$\widetilde{f}_h = \operatorname*{argmin}_{f_h \in \mathcal{F}_h} \sum_{k \in [K]} \frac{1}{\widehat{\sigma}_h^2(s_h^k, a_h^k)} \left( f_h(s_h^k, a_h^k) - r_h^k - \widehat{f}_{h+1}(s_{h+1}^k) \right)^2,$$

where $\widehat{\sigma}_h^2$ is the estimated variance and will be discussed in section 4.2. In Line 12, we subtract the confidence bonus function $b_h$ from the estimator value function $\widetilde{f}_h$ to construct the pessimistic value function $\widehat{f}_h$. With the help of the confidence bonus function $b_h$, the pessimistic value function $\widehat{f}_h$ is almost a lower bound for the optimal value function $f_h^*$. The details of the bonus function and bonus oracle will be discussed in section 4.3.

Based on the pessimistic value function $\widehat{f}_h$ for stage $h$, we recursively perform the value iteration for the stage $h - 1$. Finally, we use the pessimistic value function $\widehat{f}_h$ to do planning and output the greedy policy with respect to the pessimistic value function $\widehat{f}_h$ (Line 13 - Line 15).

## 4.2 Variance Estimate Phase

In this phase, we provide a estimator for the variance $\widehat{\sigma}_h$ in the weighted ridge regression. According to the definition of Bellman operators $\mathcal{T}$ and $\mathcal{T}_2$, the variance of the function $\widehat{f}'_{h+1}$ for each state-action pair $(s, a)$ can be denoted by

$$[\mathrm{Var}_h \widehat{f}_{h+1}](s, a) = \mathcal{T}_{2,h} \widehat{f}'_{h+1}(s, a) - \left( \mathcal{T}_h \widehat{f}'_{h+1}(s, a) \right)^2.$$

Therefore, we need the evaluate the first-order and second-order moments for $\widehat{f}'_h$. We perform nonlinear least-squares regression separately for each of these moments. Specifically, in Line 3, we conduct regression to estimate the first-order moment.

$$\widetilde{f}'_h = \operatorname*{argmin}_{f_h \in \mathcal{F}_h} \sum_{k \in [K]} \left( f_h(\bar{s}_h^k, \bar{a}_h^k) - \bar{r}_h^k - \widehat{f}'_{h+1}(\bar{s}_{h+1}^k) \right)^2.$$

In Line 4, we perform regression for the second-order moment.

$$\widetilde{g}'_h = \operatorname*{argmin}_{g_h \in \mathcal{F}_h} \sum_{k \in [K]} \left( g_h(\bar{s}_h^k, \bar{a}_h^k) - \left( \bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k) \right)^2 \right)^2.$$

In this phase, we set the variance function to 1 for each state-action pair $(s, a)$ and derive an estimator with confidence radius $\beta'_{1,h}, \beta'_{2,h}$. Combing these two regression results and subtracting a confidence bonus function $b'_h$, we create a pessimistic estimator for the variance function (Lines 6 to 7).

## 4.3 Nonlinear Bonus Oracle

As we discussed in sections 4.1 and 4.2, we introduce a uncertainty bonus function to construct a pessimistic estimate of the value function. Unfortunately, for a general class, the uncertainty bonus may varies greatly across different state-action pair. Under this situation, the addition uncertainty bonus function will highly increase the complexity of the pessimistic function class, which make it difficult to construct a accurate estimation and may significant deteriorate the final performance. To address this issue, we assume there exists a function class $\mathcal{W}$ with cardinally $|\mathcal{W}| = \mathcal{N}_b$ and can approximate the bonus function well. In addition, we assume there exists a nonlinear bonus oracle (Agarwal and Zhang, 2022), which can output the approximate bonus function in the class $\mathcal{W}$ for each dataset $\mathcal{D}_h$.

6

**Definition 4.1** (Oracle for bonus function). For an offline dataset $\mathcal{D} = \{s_h^k, a_h^k, r_h^k\}_{h,k=1}^{H,K}$, given index $h \in [H]$, let $\mathcal{D}_h = \{(s_h^k, a_h^k, r_h^k)\}_{k \in [K]}$ denote the subset of the dataset $D$ that corresponds to the observations collected up to stage h in the MDP. $\widehat{\sigma}_h(\cdot, \cdot) : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ is a variance function. $\mathcal{F}_h$ is a function class such that $\widehat{f}_h \in \mathcal{F}_h$. The parameters $\beta_h, \lambda \geq 0$, error parameter $\epsilon \geq 0$. The bonus oracle $\mathcal{B}(\widehat{\sigma}, \mathcal{D}_h, \mathcal{F}_h, \widehat{f}_h, \beta_h, \lambda, \epsilon)$ outputs a bonus function $b_h(\cdot)$ such that

- $b_h : \mathcal{S} \times \mathcal{A} \to \mathbb{R}_{\geq 0}$ belongs to function class $\mathcal{W}$.

- $b_h(z_h) \geq \max \left\{ |f_h(z_h) - \widehat{f}_h(z_h)|, f_h \in \mathcal{F}_h : \sum_{k \in [K]} \frac{(f_h(z_h^k) - \widehat{f}_h(z_h^k))^2}{(\widehat{\sigma}_h(s_h^k, a_h^k))^2} \leq (\beta_h)^2 \right\}$ for any $z_h \in \mathcal{S} \times \mathcal{A}$.

- $b_h(z_h) \leq C \cdot \left( D_{\mathcal{F}_h}(z_h; \mathcal{D}_h; \widehat{\sigma}_h) \cdot \sqrt{(\beta_h)^2 + \lambda} + \epsilon \beta_h \right)$ for all $z_h \in \mathcal{S} \times \mathcal{A}$ with constant $0 < C \leq \infty$.

**Remark 4.2.** To address the concern of function class complexity, some previous studies (Xie et al., 2021a) have approached the problem differently. Instead of introducing a pointwise bonus in the estimated value function, they solve a complicated optimization problem to guarantee the optimism solely in the intial state. This method can prevent the complexity from bonus function, as a cost, they requires solving an optimization problem over all potential policy and corresponding version space, which includes all functions with lower Bellman error.

# 5 Main Results

In this section we prove an problem-dependent regret bound of Algorithm 1.

**Theorem 5.1.** Under Assumption 3.3, for $K \geq \widetilde{\Omega} \left( \frac{\log(\mathcal{N}\mathcal{N}_b)H^6}{\kappa^2} \right)$, if we set the parameters $\beta'_{1,h}, \beta'_{2,h} = \widetilde{O}(\sqrt{\log \mathcal{N}\mathcal{N}_b}H^2)$ and $\beta_h = \widetilde{O}(\sqrt{\log \mathcal{N}})$ in Algorithm 1, then with probability at least $1 - \delta$, for any state $s \in \mathcal{S}$, we have

$$V_1^*(s) - V_1^{\widehat{\pi}}(s) \leq \widetilde{O}(\sqrt{\log \mathcal{N}}) \sum_{h=1}^{H} \mathbb{E}_{\pi^*} \left[ D_{\mathcal{F}_h}(z_h; \mathcal{D}_h; [\mathbb{V}_h V_{h+1}^*](\cdot, \cdot)) | s_1 = s \right],$$

where $[\mathbb{V}_h V_{h+1}^*](s, a) = \max\{1, [\text{Var}_h V_{h+1}^*](s, a)\}$ is the truncated conditional variance.

**Remark 5.2.** When reduce to the linear MDP environment, the following function classes

$$\mathcal{F}_h^{\text{lin}} = \{\langle \boldsymbol{\phi}_h(\cdot, \cdot), \boldsymbol{\theta}_h \rangle : \boldsymbol{\theta}_h \in \mathbb{R}^d, \|\boldsymbol{\theta}_h\|_2 \leq B_h\} \text{ for any } h \in [H],$$

satisfy the completeness assumption (Assumption 3.1) (Jin et al., 2020). Let $\mathcal{F}_h^{\text{lin}}(\epsilon)$ be a $\epsilon$-net of the linear function class $\mathcal{F}_h^{\text{lin}}$. In this case, the covering number satisfies $\log |\mathcal{F}_h^{\text{lin}}(\epsilon)| = \widetilde{O}(d)$ and the dependency of function class will reduce to $\widetilde{O}(\sqrt{\log \mathcal{N}}) = \widetilde{O}(\sqrt{d})$. For linear function class, Xiong et al. (2022) proposed the following regret guarantee,

$$V_1^*(s) - V_1^{\widehat{\pi}}(s) \leq \widetilde{O}(\sqrt{d}) \cdot \sum_{h=1}^{H} \mathbb{E}_{\pi^*} \left[ \|\boldsymbol{\phi}(s_h, a_h)\|_{\boldsymbol{\Sigma}_h^{*-1}} | s_1 = s \right],$$

where $\boldsymbol{\Sigma}_h^* = \sum_{k \in [K]} \boldsymbol{\phi}(s_h^k, a_h^k) \boldsymbol{\phi}(s_h^k, a_h^k)^\top / [\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k) + \lambda \mathbf{I}$. In comparison, we can prove the following inequality:

$$D_{\mathcal{F}_h^{\text{lin}}(\epsilon)}(z; \mathcal{D}_h; [\mathbb{V}_h V_{h+1}^*](\cdot, \cdot)) \leq \|\boldsymbol{\phi}_h(z)\|_{\boldsymbol{\Sigma}_h^{*-1}}.$$

This shows that Theorem 5.1 matches the optimal result in Xiong et al. (2022) for linear function class.

# 6 Key Techniques

In this section, we provide an overview of the key techniques in our algorithm design and analysis.

## 6.1 Variance Estimator with Nonlinear Function Class

The technique of variance-weighted ridge regression, first introduced in Zhou et al. (2021), has demonstrated its effectiveness in the online RL setting with linear function approximation. For offline

setting, Xiong et al. (2022) modified the variance-weighted ridge regression technique, and showed that using an accurate and independent variance estimator can improves the performance of the pessimistic value iteration (PEVI) algorithm (Jin et al., 2021b).

In our work, we extend this technique to general nonlinear function class $\mathcal{F}$, and use the following nonlinear least-squares regression to estimate the underlying value function:

$$\widetilde{f}_h = \underset{f_h \in \mathcal{F}_h}{\operatorname{argmin}} \sum_{k \in [K]} \frac{1}{\widehat{\sigma}_h^2(s_h^k, a_h^k)} \left( f_h(s_h^k, a_h^k) - r_h^k - \widehat{f}_{h+1}(s_{h+1}^k) \right)^2.$$

For this regression, it is crucial to obtain a reliable evaluation for the variance of the estimated cumulative reward $r_h^k + \widehat{f}_{h+1}(s_{h+1}^k)$. According to the definition of Bellman operators $\mathcal{T}$ and $\mathcal{T}_2$, the variance of the function $\widehat{f}'_{h+1}$ for each state-action pair $(s, a)$ can be denoted by

$$[\operatorname{Var}_h \widehat{f}'_{h+1}](s, a) = \mathcal{T}_{2,h} \widehat{f}'_{h+1}(s, a) - \left( \mathcal{T}_h \widehat{f}'_{h+1}(s, a) \right)^2.$$

To evaluate the first and second moment for the Bellman operator, we perform nonlinear least-squares regression on a separate dataset $\mathcal{D}'$ with uniform weight ($\widehat{\sigma}_h(s, a) = 1$ for all state-action pair $(s, a)$).

For simplicity, we denote the empirical variance as $\mathbb{B}_h(s, a) = \widetilde{g}'_h(s, a) - \left( \widetilde{f}'_h(s, a) \right)^2$, and the difference between empirical variance $\mathbb{B}_h(s, a)$ with actually variance $[\operatorname{Var}_h \widehat{f}'_{h+1}](s, a)$ is upper bound by

$$\left| \mathbb{B}_h(s, a) - [\operatorname{Var}_h \widehat{f}'_{h+1}](s, a) \right| \leq \left| \widetilde{g}'_h(s, a) - \mathcal{T}_{2,h} \widehat{f}'_{h+1}(s, a) \right| + \left| \left( \widetilde{f}'_h(s, a) \right)^2 - \left( \mathcal{T}_h \widehat{f}'_{h+1}(s, a) \right)^2 \right|.$$

For these nonlinear function estimator, the following Lemmas provide coarse concentration properties for the first and second order Bellman operators.

**Lemma 6.1.** Given a stage $h \in [H]$, let $\widehat{f}'_{h+1}(\cdot, \cdot) \leq H$ be the estimated value function constructed in Algorithm 1 Line 6. By utilizing Assumption 3.1, there exists a function $\bar{f}'_h \in \mathcal{F}_h$, such that $|\bar{f}'_h(z_h) - \mathcal{T}_h \widehat{f}'_{h+1}(z_h)| \leq \epsilon$ holds for all state-action pair $z_h = (s_h, a_h)$. Then with probability at least $1 - \delta/(4H)$, it holds that $\sum_{k \in [K]} \left( \bar{f}'_h(\bar{z}_h^k) - \widetilde{f}'_h(\bar{z}_h^k) \right)^2 \leq (\beta'_{1,h})^2$, where $\beta'_{1,h} = \widetilde{O}\left( \sqrt{\log \mathcal{N} \mathcal{N}_b} H^2 \right)$, and $\widetilde{f}'_h$ is the estimated function for first-moment Bellman operator (Line 3 in Algorithm 1).

**Lemma 6.2.** Given a stage $h \in [H]$, let $\widehat{f}'_{h+1}(\cdot, \cdot) \leq H$ be the estimated value function constructed in Algorithm 1 Line 6. By utilizing Assumption 3.1, there exists a function $\bar{g}'_h \in \mathcal{F}_h$, such that $|\bar{g}'_h(z_h) - \mathcal{T}_{2,h} \widehat{f}'_{h+1}(z_h)| \leq \epsilon$ holds for all state-action pair $z_h = (s_h, a_h)$. Then with probability at least $1 - \delta/(4H)$, it holds that $\sum_{k \in [K]} \left( \bar{g}'_h(\bar{z}_h^k) - \widetilde{g}'_h(\bar{z}_h^k) \right)^2 \leq (\beta'_{2,h})^2$, where $\beta'_{2,h} = \widetilde{O}\left( \sqrt{\log \mathcal{N} \mathcal{N}_b} H^2 \right)$, and $\widetilde{g}'_h$ is the estimated function for second-moment Bellman operator (Line 4 in Algorithm 1).

Notice that all of the previous analysis focuses on the estimated function $\widehat{f}'_{h+1}$. By leveraging an induction procedure similar to existing works in the linear case (Jin et al., 2021b; Xiong et al., 2022), we can control the distance between the estimated function $\widehat{f}'_{h+1}$ and the optimal value function $f_h^*$. In details, with high probability, for all stage $h \in [H]$, the distance is upper bounded by $O\left( \sqrt{\log \mathcal{N} \mathcal{N}_b} H^3 / \sqrt{K\kappa} \right)$. This result allows us to further bound $[\operatorname{Var}_h \widehat{f}'_{h+1}](s, a)$ and $[\operatorname{Var}_h f_{h+1}^*](s, a)$.

Therefore, the concentration properties in Lemmas 6.1 and 6.2 enable us to construct the pessimistic variance estimator, which satisfies the following property:

$$[\mathbb{V}_h V_{h+1}^*](s, a) - \widetilde{O}\left( \frac{\sqrt{\log \mathcal{N} \mathcal{N}_b} H^3}{\sqrt{K\kappa}} \right) \leq \widehat{\sigma}_h^2(s, a) \leq [\mathbb{V}_h V_{h+1}^*](s, a). \tag{6.1}$$

where $[\mathbb{V}_h V_{h+1}^*](s, a) = \max\{1, [\operatorname{Var}_h V_{h+1}^*](s, a)\}$ is the truncated conditional variance. Compared with the results in the linear function class, we utilize the logarithm of the covering number of the function class as a substitute for the linear dimension $d$, which is a common technique in nonlinear function approximation.

8

## 6.2 Reference-Advantage Decomposition

The reference-advantage decomposition is a powerful technique to tackle the challenge of additional error from uniform concentration over whole function class $\mathcal{F}_h$. Such an analysis approach has been first studied in the online RL setting Azar et al. (2017); Zhang et al. (2021); Hu et al. (2022); He et al. (2022); Agarwal et al. (2022) and later in the offline environment by Xiong et al. (2022).

For offline RL, in the context of nonlinear function classes, without a explicit linear expression, the increased complexity of the function class structure poses a significant obstacle to effectively utilizing this technique. Previous works, such as Yin et al. (2022b), have struggled to adapt the reference-advantage decomposition to their nonlinear function class, resulting in a parameter space dependence that scales with $d$, instead of the optimal $\sqrt{d}$. We provide detailed insights into this approach as follows:

$$
r_h(s,a) + \widehat{f}_{h+1}(s,a) - \mathcal{T}_h \widehat{f}_{h+1}(s,a) = \underbrace{r_h(s,a) + f^*_{h+1}(s,a) - \mathcal{T}_h f^*_{h+1}(s,a)}_{\text{Reference uncertainty}}
$$

$$
+ \underbrace{\widehat{f}_{h+1}(s,a) - f^*_{h+1}(s,a) - ([\mathbb{P}_h \widehat{f}_{h+1}](s,a) - [\mathbb{P}_h f^*_{h+1}](s,a))}_{\text{Advantage uncertainty}}.
$$

We decompose the Bellman error into two parts: the Reference uncertainty and the Advantage uncertainty. For the first term, the optimal value function $f^*_{h+1}$ is fixed and not related to the pre-collected dataset, which circumvents additional uniform concentration over the whole function class and avoid any dependence on the function class size. For the second term, it is worth to notice that the distance between the estimated function $\widehat{f'}_{h+1}$ and the optimal value function $f^*_h$ is decreased as $O(1/\sqrt{K\kappa})$. Though, we still need to maintain the uniform convergence guarantee, the Advantage uncertainty is dominated by the Reference uncertainty when the number of episode $K$ is large enough. By integrating these results, we can prove a variance-weighted concentration inequality for Bellman operators.

**Lemma 6.3.** For each stage $h \in [H]$, assuming the variance estimator $\widehat{\sigma}_h$ satisfies (6.1), let $\widehat{f}_{h+1}(\cdot, \cdot) \leq H$ be the estimated value function constructed in Algorithm 1 Line 12. By utilizing Assumption 3.1, there exists a function $\bar{f}_h \in \mathcal{F}_h$, such that $|\bar{f}_h(z_h) - \mathcal{T}_h \widehat{f}_{h+1}(z_h)| \leq \epsilon$ holds for all state-action pair $z_h = (s_h, a_h)$. Then with probability at least $1 - \delta/(4H)$, it holds that $\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( \bar{f}_h(z_h^k) - \widetilde{f}_h(z_h^k) \right)^2 \leq (\beta_h)^2$, where $\beta_h = \widetilde{O}(\sqrt{\log \mathcal{N}})$ and $\widetilde{f}_h$ is the estimated function from the weighted ridge regression (Line 10 in Algorithm 1).

After controlling the Bellman error, with a similar argument to Jin et al. (2021b); Xiong et al. (2022), we obtain the following lemma, which provide an upper bound for the regret.

**Lemma 6.4** (Regret Decomposition Property). If $|\mathcal{T}_h \widehat{f}_{h+1}(z) - \widetilde{f}_h(z)| \leq b_h(z)$ holds for all stage $h \in [H]$ and state-action pair $z = (s,a) \in \mathcal{S} \times \mathcal{A}$, then the regret of Algorithm 1 can be bounded as

$$
V^*_1(s) - V^{\widehat{\pi}}_1(s) \leq 2 \sum_{h=1}^{H} \mathbb{E}_{\pi^*} \left[ b_h(s_h, a_h) \mid s_1 = s \right].
$$

Here, the expectation $\mathbb{E}_{\pi^*}$ is with respect to the trajectory induced by $\pi^*$ in the underlying MDP.

Combing the results in Lemmas 6.3 and 6.4, we have proved Theorem 5.1.

# 7 Conclusion and Future Work

In this paper, we present PNLSVI, an oracle-efficient algorithm for offline RL with non-linear function approximation. It achieves minimax optimal problem-dependent regret when specialized to linear function approximation.

Regarding future work, we observe that instead of using the uniform coverage assumption, a series of works, such as (Liu et al., 2020; Xie et al., 2021a; Uehara and Sun, 2021; Zhan et al., 2022), only relies on partial coverage assumption. In these works, the offline data distribution only encompasses the state-action distribution of a select high-quality comparator policy $\pi^*$. It would be of significant interest to investigate whether it's possible to design practical algorithms for nonlinear function classes under this weaker partial coverage assumption, while still preserving the inherent efficiency found in linear function approximation.

## References

AGARWAL, A., JIN, Y. and ZHANG, T. (2022). Vo $q$ l: Towards optimal regret in model-free rl with nonlinear function approximation. *arXiv preprint arXiv:2212.06069* .

AGARWAL, A. and ZHANG, T. (2022). Model-based rl with optimistic posterior sampling: Structural conditions and sample complexity. *arXiv preprint arXiv:2206.07659* .

AYOUB, A., JIA, Z., SZEPESVARI, C., WANG, M. and YANG, L. (2020). Model-based reinforcement learning with value-targeted regression. In *International Conference on Machine Learning*. PMLR.

AZAR, M. G., OSBAND, I. and MUNOS, R. (2017). Minimax regret bounds for reinforcement learning. In *International Conference on Machine Learning*. PMLR.

CESA-BIANCHI, N. and LUGOSI, G. (2006). *Prediction, learning, and games*. Cambridge university press.

CHEN, J. and JIANG, N. (2019). Information-theoretic considerations in batch reinforcement learning. In *International Conference on Machine Learning*. PMLR.

DANN, C., JIANG, N., KRISHNAMURTHY, A., AGARWAL, A., LANGFORD, J. and SCHAPIRE, R. E. (2018). On oracle-efficient pac rl with rich observations. *Advances in neural information processing systems* **31**.

DEGRAVE, J., FELICI, F., BUCHLI, J., NEUNERT, M., TRACEY, B., CARPANESE, F., EWALDS, T., HAFNER, R., ABDOLMALEKI, A., DE LAS CASAS, D. ET AL. (2022). Magnetic control of tokamak plasmas through deep reinforcement learning. *Nature* **602** 414–419.

DU, S., KAKADE, S., LEE, J., LOVETT, S., MAHAJAN, G., SUN, W. and WANG, R. (2021). Bilinear classes: A structural framework for provable generalization in rl. In *International Conference on Machine Learning*. PMLR.

DU, S. S., KAKADE, S. M., WANG, R. and YANG, L. F. (2019). Is a good representation sufficient for sample efficient reinforcement learning? *arXiv preprint arXiv:1910.03016* .

FOSTER, D. J., KAKADE, S. M., QIAN, J. and RAKHLIN, A. (2021). The statistical complexity of interactive decision making. *arXiv preprint arXiv:2112.13487* .

GENTILE, C., WANG, Z. and ZHANG, T. (2022). Achieving minimax rates in pool-based batch active learning. In *International Conference on Machine Learning*. PMLR.

GU, S., HOLLY, E., LILLICRAP, T. and LEVINE, S. (2017). Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In *2017 IEEE international conference on robotics and automation (ICRA)*. IEEE.

HE, J., ZHAO, H., ZHOU, D. and GU, Q. (2022). Nearly minimax optimal reinforcement learning for linear markov decision processes. *arXiv preprint arXiv:2212.06132* .

HE, J., ZHOU, D. and GU, Q. (2021). Logarithmic regret for reinforcement learning with linear function approximation. In *International Conference on Machine Learning*. PMLR.

HU, P., CHEN, Y. and HUANG, L. (2022). Nearly minimax optimal reinforcement learning with linear function approximation. In *International Conference on Machine Learning*. PMLR.

JIANG, N., KRISHNAMURTHY, A., AGARWAL, A., LANGFORD, J. and SCHAPIRE, R. E. (2017). Contextual decision processes with low bellman rank are pac-learnable. In *International Conference on Machine Learning*. PMLR.

JIN, C., LIU, Q. and MIRYOOSEFI, S. (2021a). Bellman eluder dimension: New rich classes of rl problems, and sample-efficient algorithms. *Advances in neural information processing systems* **34** 13406–13418.

JIN, C., YANG, Z., WANG, Z. and JORDAN, M. I. (2020). Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*. PMLR.

JIN, Y., YANG, Z. and WANG, Z. (2021b). Is pessimism provably efficient for offline rl? In *International Conference on Machine Learning*. PMLR.

LEVINE, S., PASTOR, P., KRIZHEVSKY, A., IBARZ, J. and QUILLEN, D. (2018). Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *The International journal of robotics research* **37** 421–436.

LI, G., SHI, L., CHEN, Y., CHI, Y. and WEI, Y. (2022). Settling the sample complexity of model-based offline reinforcement learning. *arXiv preprint arXiv:2204.05275* .

LIU, Y., SWAMINATHAN, A., AGARWAL, A. and BRUNSKILL, E. (2020). Provably good batch off-policy reinforcement learning without great exploration. *Advances in neural information processing systems* **33** 1264–1274.

MIN, Y., WANG, T., ZHOU, D. and GU, Q. (2021). Variance-aware off-policy evaluation with linear function approximation. *Advances in neural information processing systems* **34** 7598–7610.

MODI, A., JIANG, N., TEWARI, A. and SINGH, S. (2020). Sample complexity of reinforcement learning using linearly combined model ensembles. In *International Conference on Artificial Intelligence and Statistics*. PMLR.

NGUYEN-TANG, T. and ARORA, R. (2023). Viper: Provably efficient algorithm for offline rl with neural function approximation. In *The Eleventh International Conference on Learning Representations*.

RASHIDINEJAD, P., ZHU, B., MA, C., JIAO, J. and RUSSELL, S. (2021). Bridging offline reinforcement learning and imitation learning: A tale of pessimism. *Advances in Neural Information Processing Systems* **34** 11702–11716.

SCHRITTWIESER, J., ANTONOGLOU, I., HUBERT, T., SIMONYAN, K., SIFRE, L., SCHMITT, S., GUEZ, A., LOCKHART, E., HASSABIS, D., GRAEPEL, T. ET AL. (2020). Mastering atari, go, chess and shogi by planning with a learned model. *Nature* **588** 604–609.

SHI, L., LI, G., WEI, Y., CHEN, Y. and CHI, Y. (2022). Pessimistic q-learning for offline reinforcement learning: Towards optimal sample complexity. In *International Conference on Machine Learning*. PMLR.

SILVER, D., SCHRITTWIESER, J., SIMONYAN, K., ANTONOGLOU, I., HUANG, A., GUEZ, A., HUBERT, T., BAKER, L., LAI, M., BOLTON, A. ET AL. (2017). Mastering the game of go without human knowledge. *nature* **550** 354–359.

SUN, W., JIANG, N., KRISHNAMURTHY, A., AGARWAL, A. and LANGFORD, J. (2019). Model-based rl in contextual decision processes: Pac bounds and exponential improvements over model-free approaches. In *Conference on learning theory*. PMLR.

UEHARA, M. and SUN, W. (2021). Pessimistic model-based offline reinforcement learning under partial coverage. *arXiv preprint arXiv:2107.06226* .

WANG, R., FOSTER, D. P. and KAKADE, S. M. (2020a). What are the statistical limits of offline rl with linear function approximation? *arXiv preprint arXiv:2010.11895* .

WANG, R., SALAKHUTDINOV, R. R. and YANG, L. (2020b). Reinforcement learning with general value function approximation: Provably efficient approach via bounded eluder dimension. *Advances in Neural Information Processing Systems* **33** 6123–6135.

WANG, Y., WANG, R., DU, S. S. and KRISHNAMURTHY, A. (2020c). Optimism in reinforcement learning with generalized linear function approximation. In *International Conference on Learning Representations*.

WEISZ, G., AMORTILA, P. and SZEPESVÁRI, C. (2021). Exponential lower bounds for planning in mdps with linearly-realizable optimal action-value functions. In *Algorithmic Learning Theory*. PMLR.

XIE, T., CHENG, C.-A., JIANG, N., MINEIRO, P. and AGARWAL, A. (2021a). Bellman-consistent pessimism for offline reinforcement learning. *Advances in neural information processing systems* **34** 6683–6694.

XIE, T., JIANG, N., WANG, H., XIONG, C. and BAI, Y. (2021b). Policy finetuning: Bridging sample-efficient offline and online reinforcement learning. *Advances in neural information processing systems* **34** 27395–27407.

XIONG, W., ZHONG, H., SHI, C., SHEN, C., WANG, L. and ZHANG, T. (2022). Nearly minimax optimal offline reinforcement learning with linear function approximation: Single-agent mdp and markov game. *arXiv preprint arXiv:2205.15512* .

YANG, L. and WANG, M. (2019). Sample-optimal parametric q-learning using linearly additive features. In *International Conference on Machine Learning*.

YANG, L. and WANG, M. (2020). Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. In *International Conference on Machine Learning*. PMLR.

YIN, M., DUAN, Y., WANG, M. and WANG, Y.-X. (2022a). Near-optimal offline reinforcement learning with linear representation: Leveraging variance information with pessimism. *arXiv preprint arXiv:2203.05804* .

YIN, M., WANG, M. and WANG, Y.-X. (2022b). Offline reinforcement learning with differentiable function approximation is provably efficient. *arXiv preprint arXiv:2210.00750* .

YIN, M. and WANG, Y.-X. (2021). Towards instance-optimal offline reinforcement learning with pessimism. *Advances in neural information processing systems* **34** 4065–4078.

ZANETTE, A., BRANDFONBRENER, D., BRUNSKILL, E., PIROTTA, M. and LAZARIC, A. (2020a). Frequentist regret bounds for randomized least-squares value iteration. In *International Conference on Artificial Intelligence and Statistics*. PMLR.

ZANETTE, A., LAZARIC, A., KOCHENDERFER, M. and BRUNSKILL, E. (2020b). Learning near optimal policies with low inherent bellman error. In *International Conference on Machine Learning*. PMLR.

ZANETTE, A., WAINWRIGHT, M. J. and BRUNSKILL, E. (2021). Provable benefits of actor-critic methods for offline reinforcement learning. *Advances in neural information processing systems* **34** 13626–13640.

ZHAN, W., HUANG, B., HUANG, A., JIANG, N. and LEE, J. (2022). Offline reinforcement learning with realizability and single-policy concentrability. In *Conference on Learning Theory*. PMLR.

ZHANG, Z., JI, X. and DU, S. (2021). Is reinforcement learning more difficult than bandits? a near-optimal algorithm escaping the curse of horizon. In *Conference on Learning Theory*. PMLR.

ZHOU, D., GU, Q. and SZEPESVARI, C. (2021). Nearly minimax optimal reinforcement learning for linear mixture markov decision processes. In *Conference on Learning Theory*. PMLR.

# A   Comparison of data coverage assumptions

In Yin et al. (2022b), they studied the general differentiable function class, where the function class can be denoted by

$$\mathcal{F} := \Big\{ f\big(\boldsymbol{\theta}, \phi(\cdot, \cdot)\big) : \mathcal{X} \times \mathcal{A} \to \mathbb{R}, \boldsymbol{\theta} \in \Theta \Big\}.$$

In this definition, $\Psi$ is a compact subset and $\phi(\cdot, \cdot) : \mathcal{X} \times \mathcal{A} \to \Psi \subseteq \mathbb{R}^m$ is a feature map. The parameter space $\Theta$ is a compact subset $\Theta \subseteq \mathbb{R}^d$. The function $f : \mathbb{R}^d \times \mathbb{R}^m \to \mathbb{R}$ satisfies the following smoothness conditions:

- For any vector $\phi \in \mathbb{R}^m$, $f(\boldsymbol{\theta}, \phi)$ is third-time differentiable with respect to the parameter $\boldsymbol{\theta}$.
- Functions $f, \partial_{\boldsymbol{\theta}} f, \partial^2_{\boldsymbol{\theta},\boldsymbol{\theta}} f, \partial^3_{\boldsymbol{\theta},\boldsymbol{\theta},\boldsymbol{\theta}} f$ are jointly continuous for $(\boldsymbol{\theta}, \phi)$.

Under this definition, Yin et al. (2022b) introduce the following coverage assumption (Assumption 2.3) such that for all stage $h \in [H]$, there exists a constant $\kappa$,

$$\mathbb{E}_{\mu,h} \left[ \big( f\left( \boldsymbol{\theta}_1, \phi(x, a) \right) - f(\boldsymbol{\theta}_2, \phi(x, a)) \big)^2 \right] \geq \kappa \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2^2, \forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \Theta; (*)$$

$$\mathbb{E}_{\mu,h} \left[ \nabla f(\boldsymbol{\theta}, \phi(x, a)) \nabla f(\boldsymbol{\theta}, \phi(x, a))^\top \right] \succ \kappa I, \forall \boldsymbol{\theta} \in \Theta. (**)$$

It is worth noting that our assumption 3.3 is weaker than this assumption. For any compact sets $\Theta$, $\Psi$ and continuous function $f$, there always exist a constant $\kappa_0 > 0$ such that $f$ is $\kappa_0$-Lipschitz with respect to the parameter $\boldsymbol{\theta}$,i.e:

$$|f(\boldsymbol{\theta}_1, \phi) - f(\boldsymbol{\theta}_2, \phi)| \leq \kappa_0 \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2, \forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \Theta, \phi \in \Psi.$$

Therefore, the coverage assumption in Yin et al. (2022b) implies that

$$\mathbb{E}_{\mu,h} \left[ \big( f(\boldsymbol{\theta}_1, \phi(\cdot, \cdot)) - f(\boldsymbol{\theta}_2, \phi(\cdot, \cdot)) \big)^2 \right] \geq \kappa \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2^2$$

$$\geq \frac{\kappa}{\kappa_0^2} \sup_{(x,a) \in \mathcal{X} \times \mathcal{A}} \big( f(\boldsymbol{\theta}_1, \phi(x, a)) - f(\boldsymbol{\theta}_2, \phi(x, a)) \big)^2.$$

Thus, our assumption is weaker than the first assumption (*). For the second assumption (**), there is no direct counterpart in the general setting.

In addition, for the linear function class, the coverage assumption in Yin et al. (2022b) will reduce to the following linear function coverage assumption(Wang et al., 2020a; Min et al., 2021; Yin et al., 2022a; Xiong et al., 2022).

$$\lambda_{\min}(\mathbb{E}_{\mu,h}[\phi(x, a)\phi(x, a)^\top]) = \kappa > 0, \ \forall h \in [H].$$

Therefore, our assumption is also weaker than the linear function coverage assumption when dealing with the linear function class.

# B   Proof of Theorem 5.1

We need the following lemmas to prove Theorem 5.1. To start with, we prove the result that our data coverage assumption (Assumption 3.3) can lead to an upper bound of the $D^2$-divergence for large dataset.

**Lemma B.1.** Let $\mathcal{D}_h$ be the dataset satisfying Assumption 3.3. When the size of data set satisfies $K \geq \widetilde{\Omega}\left(\frac{\log \mathcal{N}}{\kappa^2}\right)$, with probability at least $1 - \delta$, for each state-action pair $z$, we have

$$D_{\mathcal{F}_h}(z, \mathcal{D}_h, 1) = \widetilde{O}\left(\frac{1}{\sqrt{K\kappa}}\right).$$

**Lemma B.2.** Let $\mathcal{D}_h$ be a dataset satisfying Assumption 3.3. When the size of data set satisfies $K \geq \widetilde{\Omega}\left(\frac{\log \mathcal{N}}{\kappa^2}\right), \widehat{\sigma}_h \leq H$, with probability at least $1 - \delta$, for each state-action pair $z$, we have

$$D_{\mathcal{F}_h}(z, \mathcal{D}_h, \widehat{\sigma}_h) = \widetilde{O}\left(\frac{H}{\sqrt{K\kappa}}\right).$$

The following lemmas show the confidence radius for the first and second-order Bellman error.

**Lemma B.3** (Restatement of Lemma 6.1). At stage $h \in [H]$, the estimated value function $\widehat{f}'_{h+1}$ in Algorithm 1 is bounded by $H$. According to Assumption 3.1, there exists some function $\bar{f}'_h \in \mathcal{F}_h$, such that $|\bar{f}'_h(z_h) - \mathcal{T}_h \widehat{f}'_{h+1}(z_h)| \leq \epsilon$ for all $z_h = (s_h, a_h)$. Then with probability at least $1 - \delta/(4H^2)$, the following inequality holds:

$$\sum_{k \in [K]} \left( \bar{f}'_h(\bar{z}^k_h) - \widetilde{f}'_h(\bar{z}^k_h) \right)^2 \leq (\beta'_{1,h})^2,$$

where $\beta'_{1,h} = \widetilde{O}(\sqrt{\log \mathcal{N} \mathcal{N}_b} H)$.

The following lemma for second-order function approximation parallels the lemma we have proved.

**Lemma B.4** (Restatement of Lemma 6.2). At stage $h \in [H]$, the estimated value function $\widehat{f}'_{h+1}$ in Algorithm 1 is bounded by $H$. According to Assumption 3.1, there exists some functions $\bar{g}'_h \in \mathcal{F}_h$, such that $|\bar{g}_h(z_h) - \mathcal{T}_{2,h} \widehat{f}_{h+1}(z_h)| \leq \epsilon$ for all $z_h = (x_h, a_h)$. Then with probability at least $1 - \delta/(4H^2)$, the following inequality holds:

$$\sum_{k \in [K]} \left( \bar{g}'_h(\bar{z}^k_h) - \widetilde{g}'_h(\bar{z}^k_h) \right)^2 \leq (\beta'_{2,h})^2,$$

where $\beta'_{2,h} = \widetilde{O}(\sqrt{\log \mathcal{N} \mathcal{N}_b} H^2)$.

Using Lemma B.3 and Lemma B.4, we can prove a high probability bound of the variance estimator. We first recall the definition of the variance estimator.

$$\widetilde{f}'_h = \underset{f_h \in \mathcal{F}_h}{\arg\min} \sum_{k \in [K]} \left( f_h(\bar{s}^k_h, \bar{a}^k_h) - \bar{r}^k_h - \widehat{f}'_{h+1}(\bar{s}^k_{h+1}) \right)^2$$

$$\widetilde{g}'_h = \underset{g_h \in \mathcal{F}_h}{\arg\min} \sum_{k \in [K]} \left( g_h(\bar{s}^k_h, \bar{a}^k_h) - \left( \bar{r}^k_h + \widehat{f}'_{h+1}(\bar{s}^k_{h+1}) \right)^2 \right)^2.$$

We then employ the following variance estimator:

$$\widehat{\sigma}^2_h(s, a) := \max \left\{ 1, \widetilde{g}'_h(s, a) - \left( \widetilde{f}'_h(s, a) \right)^2 - \widetilde{O}\left( \frac{\sqrt{\log \mathcal{N} \mathcal{N}_b} H^3}{\sqrt{K\kappa}} \right) \right\}.$$

The following lemma shows our constructed estimator is closed to the actual variance of the optimal value function $[\mathbb{V}_h V^*_{h+1}](s, a)$.

**Lemma B.5.** with probability at least $1 - \delta/2$, for any $h \in [H]$, the variance estimator designed above satisfies:

$$[\mathbb{V}_h V^*_{h+1}](s, a) - \widetilde{O}\left( \frac{\sqrt{\log \mathcal{N} \mathcal{N}_b} H^3}{\sqrt{K\kappa}} \right) \leq \widehat{\sigma}^2_h(s, a) \leq [\mathbb{V}_h V^*_{h+1}](s, a).$$

With the property of $\widehat{\sigma}_h$ in Lemma B.5, we can prove a variance weighted version of concentration inequality.

**Lemma B.6** (Restatement of Lemma 6.3). Suppose the variance function $\widehat{\sigma}_h$ satisfies the inequality in Lemma B.5. at stage $h \in [H]$, the estimated value function $\widehat{f}_{h+1}$ in Algorithm 1 is bounded by $H$. According to Assumption 3.1, there exists some function $\bar{f}_h \in \mathcal{F}_h$, such that $|\bar{f}_h(z_h) - \mathcal{T}_h \widehat{f}_{h+1}(z)| \leq \epsilon$ for all $z_h = (s_h, a_h)$. Then with probability at least $1 - \delta/2$, the following inequality holds for all stage $h \in [H]$ simultaneously,

$$\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z^k_h))^2} \left( \bar{f}_h(z^k_h) - \widetilde{f}_h(z^k_h) \right)^2 \leq (\beta_h)^2.$$

Finally, we start the proof of Theorem 5.1.

14

*Proof of Theorem 5.1.* For any state-action pair $z = (s, a) \in \mathcal{S} \times \mathcal{A}$, we have

$$\left| \mathcal{T}_h \widehat{f}_{h+1}(z) - \widetilde{f}_h(z) \right| \leq \left| \mathcal{T}_h \widehat{f}_{h+1}(z) - \bar{f}_h(z) \right| + \left| \bar{f}_h(z) - \widetilde{f}_h(z) \right|$$
$$\leq \epsilon + b_h(z),$$

where we bound the first term with the Bellman completeness assumption (Assumption 3.1). For the second term, we use the bonus oracle (Definition 4.1) and Lemma B.6. Therefore, using Lemma D.2 we have

$$V_1^*(s) - \widehat{V}_1(s) \leq 2 \sum_{h=1}^{H} \mathbb{E}_{\pi^*} \left[ b_h \left( s_h, a_h \right) \mid s_1 = s \right] + 2\epsilon H$$

$$\leq \sum_{h=1}^{H} \mathbb{E}_{\pi^*} \left[ D_{\mathcal{F}_h}(z_h; \mathcal{D}_h; \widehat{\sigma}_h) \cdot \sqrt{(\beta_h)^2 + \lambda} \mid s_1 = s \right] + 2\epsilon H$$

$$\leq \widetilde{O} \left( \sqrt{\log \mathcal{N}} \right) \sum_{h=1}^{H} \mathbb{E}_{\pi^*} \left[ D_{\mathcal{F}_h}(z_h; \mathcal{D}_h; \widehat{\sigma}_h) \mid s_1 = s \right]$$

$$\leq \widetilde{O} \left( \sqrt{\log \mathcal{N}} \right) \sum_{h=1}^{H} \mathbb{E}_{\pi^*} \left[ D_{\mathcal{F}_h} \left( z_h; \mathcal{D}_h; [\mathbb{V}_h V_{h+1}^*](\cdot, \cdot) \right) \mid s_1 = s \right].$$

Here the second inequality holds because of our choice of bonus function (Definition 4.1). We use the definition of $\beta_h = \widetilde{O} \left( \sqrt{\log \mathcal{N}} \right)$ in the third inequality. Finally, due to Lemma B.5, with probability at least $1 - \delta$, for any $z \in \mathcal{S} \times \mathcal{A}$, we have $\widehat{\sigma}_h(z) \leq [\mathbb{V}_h V_{h+1}^*](z)$. Therefore, using the fact that $D^2$-divergence is increasing with respect to the variance function (Definition 3.2), we have proved the last inequality. $\qquad\square$

## C  Proof of the Lemmas in Section B

### C.1  Proof of Lemma B.1 and Lemma B.2

*Proof of Lemma B.1.* From the definition of $D^2$ divergence (Definition 3.2), we have

$$D^2_{\mathcal{F}_h}(z; \mathcal{D}_h; 1) = \sup_{f_1, f_2 \in \mathcal{F}_h} \frac{(f_1(z) - f_2(z))^2}{\sum_{k \in [K]} \left( f_1(z_h^k) - f_2(z_h^k) \right)^2 + \lambda} \tag{C.1}$$

By Hoeffding inequality (Lemma D.3), with probability at least $1 - \delta/(\mathcal{N}^2)$, we have

$$\sum_{k \in [K]} \left( f_1(z_h^k) - f_2(z_h^k) \right)^2 - K \mathbb{E}_{\mu, h} \left[ (f_1(z_h) - f_2(z_h))^2 \right] \geq -2\sqrt{2K \log(\mathcal{N}^2/\delta)} \cdot \|f_1 - f_2\|_\infty^2.$$

Hence, after taking a union bound, we have with probability at least $1 - \delta$, for all $f_1, f_2 \in \mathcal{F}_h$,

$$\sum_{k \in [K]} \left( f_1(z_h^k) - f_2(z_h^k) \right)^2 \geq K \mathbb{E}_{\mu, h} \left[ (f_1(z_h) - f_2(z_h))^2 \right] - 2\sqrt{2K \log(\mathcal{N}^2/\delta)} \cdot \|f_1 - f_2\|_\infty^2$$

$$\geq K \cdot \kappa \|f_1 - f_2\|_\infty^2 - 2\sqrt{2K \log(\mathcal{N}^2/\delta)} \cdot \|f_1 - f_2\|_\infty^2, \tag{C.2}$$

where the second inequality holds due to Assumption 3.3. Substituting (C.2) into (C.1), when the size of dataset $K \geq \widetilde{\Omega} \left( \frac{\log \mathcal{N}}{\kappa^2} \right)$, we have

$$D^2_{\mathcal{F}_h}(z; \mathcal{D}_h; 1) \leq \sup_{f_1, f_2 \in \mathcal{F}_h} \frac{(f_1(z) - f_2(z))^2}{\frac{1}{2} K \cdot \kappa \|f_1 - f_2\|_\infty^2 + \lambda} = \widetilde{O} \left( \frac{1}{K\kappa} \right).$$

$\qquad\square$

545 *Proof of Lemma B.2.* From the definition of $D^2$ divergence, we have

$$D^2_{\mathcal{F}_h}(z; \mathcal{D}_h; \widehat{\sigma}_h) = \sup_{f_1, f_2 \in \mathcal{F}_h} \frac{(f_1(z) - f_2(z))^2}{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left(f_1(z_h^k) - f_2(z_h^k)\right)^2 + \lambda} \tag{C.3}$$

546 By Hoeffding inequality (Lemma D.3), with probability at least $1 - \delta/(\mathcal{N}^2)$,

$$\sum_{k \in [K]} \left(f_1(z_h^k) - f_2(z_h^k)\right)^2 - K \mathbb{E}_{\mu, h} \left[(f_1(z_h) - f_2(z_h))^2\right] \geq -2\sqrt{2K \log(\mathcal{N}^2/\delta)} \cdot \|f_1 - f_2\|_\infty^2.$$

547 Hence, after taking a union bound, we have with probability at least $1 - \delta$, for all $f_1, f_2 \in \mathcal{F}_h$,

$$\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left(f_1(z_h^k) - f_2(z_h^k)\right)^2$$

$$\geq \frac{1}{H^2} \left(K \mathbb{E}_{\mu, h} \left[(f_1(z_h) - f_2(z_h))^2\right] - 2\sqrt{2K \log(\mathcal{N}^2/\delta)} \cdot \|f_1 - f_2\|_\infty^2\right)$$

$$\geq \frac{1}{H^2} \left(K \cdot \kappa \|f_1 - f_2\|_\infty^2 - 2\sqrt{2K \log(\mathcal{N}^2/\delta)} \cdot \|f_1 - f_2\|_\infty^2\right), \tag{C.4}$$

548 where we use Assumption 3.3 and substituting (C.4) into (C.3), when $K \geq \widetilde{\Omega}\left(\frac{\log \mathcal{N}}{\kappa}\right)$, we have

$$D^2_{\mathcal{F}_h}(z; \mathcal{D}_h; \widehat{\sigma}_h) \leq \sup_{f_1, f_2 \in \mathcal{F}_h} \frac{H^2(f_1(z) - f_2(z))^2}{\frac{1}{2}K \cdot \kappa \|f_1 - f_2\|_\infty^2 + \lambda} = \widetilde{O}\left(\frac{H^2}{K\kappa^2}\right).$$

549 $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

## 550 C.2 Proof of Lemma B.3

551 We need to prove the following concentration inequality first.

552 **Lemma C.1.** Based on the dataset $\mathcal{D}' = \{\bar{s}_h^k, \bar{a}_h^k, \bar{r}_h^k\}_{k,h=1}^{K,H}$, we define the filtration

$$\bar{\mathcal{H}}_h^k = \sigma\left(\bar{s}_1^1, \bar{a}_1^1, \bar{r}_1^1, \bar{s}_2^1, \ldots, \bar{r}_H^1, \bar{s}_{H+1}^1; \bar{s}_1^2, \bar{a}_1^2, \bar{r}_1^2, \bar{s}_2^2, \ldots, \bar{r}_H^2, \bar{s}_{H+1}^2; \ldots; \bar{s}_1^k, \bar{a}_1^k, \bar{r}_1^k, \bar{s}_2^k, \ldots, \bar{r}_h^k, \bar{s}_{h+1}^k\right).$$

553 For any fixed functions $f, f' : \mathcal{S} \to [0, L]$, we make the following definitions:

$$\bar{\eta}_h^k[f'] := f'(\bar{s}_{h+1}^k) - [\mathbb{P}_h f'](\bar{s}_h^k, \bar{a}_h^k)$$

$$\bar{D}_h^k[f, f'] := 2\bar{\eta}_h^k[f'] \left(f(\bar{z}_h^k) - \mathcal{T}_h f'(\bar{z}_h^k)\right).$$

554 Then with probability at least $1 - \delta/(4H^2\mathcal{N}^2\mathcal{N}_b^2)$, the following inequality holds,

$$\sum_{k \in [K]} \bar{D}_h^k[f, f'] \leq (24L + 5)i^2(\delta) + \frac{\sum_{k \in [K]} \left(f(\bar{z}_h^k) - \mathcal{T}_h f'(\bar{z}_h^k)\right)^2}{2},$$

555 where $i(\delta) = \sqrt{2 \log \frac{\mathcal{N}\mathcal{N}_b H(2\log(4LK)+2)(\log(2L)+2)}{\delta}}$.

556 *Proof.* We use Lemma D.1, with the following conditions:

$\bar{D}_h^k[f, f']$ is adapted to the filtration $\bar{\mathcal{H}}_h^k$ and $\mathbb{E}\left[\bar{D}_h^k[f, f'] \mid \bar{\mathcal{H}}_h^{k-1}\right] = 0$.

$\left|\bar{D}_h^k[f, f']\right| \leq 2\left|\bar{\eta}_h^k\right| \max_z |f(z) - \mathcal{T}_h f'(z)| \leq 4L^2 = M.$

$$\sum_{k \in [K]} \mathbb{E}\left[\left(\bar{D}_h^k[f, f']\right)^2 \Big| \bar{z}_h^k\right] = \sum_{k \in [K]} \mathbb{E}\left[4\left(\bar{\eta}_h^k[f']\right)^2 \Big| \bar{z}_h^k\right]\left(f(\bar{z}_h^k) - \mathcal{T}_h f'(\bar{z}_h^k)\right)^2 \leq (4LK)^2 = V^2.$$

557 On the other hand,

$$\sum_{k \in [K]} \mathbb{E}\left[\left(\bar{D}_h^k[f, f']\right)^2 \Big| \bar{z}_h^k\right] = \sum_{k \in [K]} \mathbb{E}\left[4\left(\bar{\eta}_h^k[f']\right)^2 \Big| \bar{z}_h^k\right]\left(f(\bar{z}_h^k) - \mathcal{T}_h f'(\bar{z}_h^k)\right)^2$$

$$\leq 8L^2 \sum_{k \in [K]} \left(f(\bar{z}_h^k) - \mathcal{T}_h f'(\bar{z}_h^k)\right)^2.$$

16

Then using Lemma D.1 with $v = 1$, $m = 1$, with high probability, we have:

$$\sum_{k \in [K]} 2\bar{\eta}_h^k[f'] \left( f(\bar{z}_h^k) - \mathcal{T}_h f'(\bar{z}_h^k) \right) \le i(\delta) \sqrt{2(2 \cdot 8L^2) \sum_{k \in [K]} \left( f(\bar{z}_h^k) - \mathcal{T}_h f'(\bar{z}_h^k) \right)^2}$$

$$+ \frac{2}{3} i^2(\delta) + \frac{4}{3} i^2(\delta) \cdot 4L^2$$

$$\le (24L^2 + 5) i^2(\delta) + \frac{\sum_{k \in [K]} \left( f(\bar{z}_h^k) - \mathcal{T}_h f'(\bar{z}_h^k) \right)^2}{2}.$$

$\square$

*Proof of Lemma B.3.* Let $(\beta'_{1,h})^2 = (24L^2 + 5) i^2(\delta) + 8KL\epsilon$. We define the event $\mathcal{E}'_{1,h} :=$ $\left\{ \sum_{k \in [K]} \left( \bar{f}'_h(\bar{z}_h^k) - \widetilde{f}'_h(\bar{z}_h^k) \right)^2 > (\beta'_{1,h})^2 \right\}$. The following inequality will be useful in our proof.

$$\sum_{k \in [K]} \left( \bar{f}'_h(\bar{z}_h^k) - \widetilde{f}'_h(\bar{z}_h^k) \right)^2 = \sum_{k \in [K]} \left[ \left( \bar{r}_h^k + \widehat{f'}_{h+1}(\bar{s}_{h+1}^k) - \bar{f}'_h(\bar{z}_h^k) \right) + \left( \widetilde{f}'_h(\bar{z}_h^k) - \bar{r}_h^k - \widehat{f'}_{h+1}(\bar{s}_{h+1}^k) \right) \right]^2$$

$$= \sum_{k \in [K]} \left( \bar{r}_h^k + \widehat{f'}_{h+1}(\bar{s}_{h+1}^k) - \bar{f}'_h(\bar{z}_h^k) \right)^2 + \sum_{k \in [K]} \left( \widetilde{f}'_h(\bar{z}_h^k) - \bar{r}_h^k - \widehat{f'}_{h+1}(\bar{s}_{h+1}^k) \right)^2$$

$$+ 2 \sum_{k \in [K]} \left( \bar{r}_h^k + \widehat{f'}_{h+1}(\bar{s}_{h+1}^k) - \bar{f}'_h(\bar{z}_h^k) \right) \left( \widetilde{f}'_h(\bar{z}_h^k) - \bar{r}_h^k - \widehat{f'}_{h+1}(\bar{s}_{h+1}^k) \right)$$

$$\le 2 \sum_{k \in [K]} \left( \bar{r}_h^k + \widehat{f'}_{h+1}(\bar{s}_{h+1}^k) - \bar{f}'_h(\bar{z}_h^k) \right)^2$$

$$+ 2 \sum_{k \in [K]} \left( \bar{r}_h^k + \widehat{f'}_{h+1}(\bar{s}_{h+1}^k) - \bar{f}'_h(\bar{z}_h^k) \right) \left( \widetilde{f}'_h(\bar{z}_h^k) - \bar{r}_h^k - \widehat{f'}_{h+1}(\bar{s}_{h+1}^k) \right)$$

$$\le 2 \sum_{k \in [K]} \left( \bar{r}_h^k + \widehat{f'}_{h+1}(\bar{s}_{h+1}^k) - \bar{f}'_h(\bar{z}_h^k) \right) \left( \widetilde{f}'_h(\bar{z}_h^k) - \bar{f}'_h(\bar{z}_h^k) \right).$$

Here we use our choice of $\widetilde{f}'_h$, i.e. $\widetilde{f}'_h = \operatorname{argmin}_{f_h \in \mathcal{F}_h} \sum_{k \in [K]} \left( f_h(\bar{s}_h^k, \bar{a}_h^k) - \bar{r}_h^k - \widehat{f'}_{h+1}(\bar{s}_{h+1}^k) \right)^2$.

Next, we will use Lemma C.1. For any fixed $h$, let $f = \widetilde{f}'_h \in \mathcal{F}_h$, $f' = \widehat{f'}_{h+1} = \{\widetilde{f} - \epsilon\}_{[0, H-h+1]}$, where $\widetilde{f} = \widetilde{f}'_h - b'_h \in \mathcal{F}_h - \mathcal{W}$. Following the construction in Lemma C.1, we define

$$\bar{\eta}_h^k[f'] = \bar{r}_h^s + f'(\bar{s}_{h+1}^k) - \mathbb{E}\left[ \bar{r}_h^k + f'(\bar{s}_{h+1}^k) | \bar{z}_h^k \right],$$
$$\text{and } \bar{D}_h^k[f, f'] = 2\bar{\eta}_h^k[f'] \left( f(\bar{z}_h^k) - \mathcal{T}_h \bar{f}'(\bar{z}_h^k) \right).$$

Due to the result of Lemma C.1, taking a union bound, we have with probability at least $1 - \delta/(4H^2)$, the following inequality holds,

$$\sum_{k \in [K]} \bar{D}_h^k[f, f'] \le (24L^2 + 5) i^2(\delta) + \frac{\sum_{k \in [K]} \left( f(\bar{z}_h^k) - \mathcal{T}_h f'(\bar{z}_h^k) \right)^2}{2}. \tag{C.5}$$

17

567    Therefore, with probability at least $1 - \delta/(4H^2)$, we have

$$2 \sum_{k \in [K]} \left( \bar{r}_h^k + \widehat{f}_{h+1}'(\bar{s}_{h+1}^k) - \bar{f}_h'(\bar{z}_h^k) \right) \left( \widetilde{f}_h'(\bar{z}_h^k) - \bar{f}_h'(\bar{z}_h^s) \right)$$

$$= 2 \sum_{k \in [K]} \left( \bar{r}_h^k + \widehat{f}_{h+1}'(\bar{s}_{h+1}^k) - \mathcal{T}_h \widehat{f}_{h+1}'(\bar{z}_h^k) \right) \left( \widetilde{f}_h'(\bar{z}_h^k) - \bar{f}_h'(\bar{z}_h^k) \right)$$

$$+ 2 \sum_{k \in [K]} \left( \mathcal{T}_h \widehat{f}_{h+1}'(\bar{z}_h^k) - \bar{f}_h'(\bar{z}_h^k) \right) \left( \widetilde{f}_h'(\bar{z}_h^k) - \bar{f}_h'(\bar{z}_h^k) \right)$$

$$\leq 2 \sum_{k \in [K]} \left( \bar{r}_h^k + \widehat{f}_{h+1}'(\bar{s}_{h+1}^k) - \mathcal{T}_h \widehat{f}_{h+1}'(\bar{z}_h^k) \right) \left( \widetilde{f}_h'(\bar{z}_h^k) - \bar{f}_h'(\bar{z}_h^k) \right) + 4KL\epsilon$$

$$\leq (24L^2 + 5)i^2(\delta) + 4KL\epsilon + \frac{\sum_{k \in [K]} \left( \widetilde{f}_h'(\bar{z}_h^k) - \mathcal{T}_h \widehat{f}_{h+1}'(\bar{z}_h^k) \right)^2}{2}$$

$$\leq (24L^2 + 5)i^2(\delta) + 8KL\epsilon + \frac{\sum_{k \in [K]} \left( \bar{f}_h'(\bar{z}_h^k) - \widetilde{f}_h'(\bar{z}_h^k) \right)^2}{2}$$

$$\leq \frac{(\beta_{1,h}')^2}{2} + \frac{\sum_{k \in [K]} \left( \bar{f}_h'(\bar{z}_h^k) - \widetilde{f}_h'(\bar{z}_h^k) \right)^2}{2}.$$

568    Here the second inequality holds because of the Bellman completeness assumption (Assumption 3.1).
569    The third inequality arises from (C.5). The last inequality holds due to the choice of

$$\beta_{1,h}' = \sqrt{2(24L^2 + 5)i^2(\delta) + 16KL\epsilon} = \widetilde{O}\left( \sqrt{\log \mathcal{N}\mathcal{N}_b} H \right).$$

570    But conditioned on the event $\mathcal{E}_{1,h}'$, we have

$$\sum_{k \in [K]} \left( \bar{r}_h^k + \widehat{f}_{h+1}'(\bar{s}_{h+1}^k) - \bar{f}_h'(\bar{z}_h^k) \right) \left( \widetilde{f}_h'(\bar{z}_h^k) - \bar{f}_h'(\bar{z}_h^k) \right)$$

$$\geq \sum_{k \in [K]} \left( \bar{f}_h'(\bar{z}_h^k) - \widetilde{f}_h'(\bar{z}_h^k) \right)^2$$

$$> \frac{(\beta_{1,h}')^2}{2} + \frac{\sum_{k \in [K]} \left( \bar{f}_h'(\bar{z}_h^k) - \widetilde{f}_h'(\bar{z}_h^k) \right)^2}{2}.$$

571    Thus, we have $\mathbb{P}[\mathcal{E}_{1,h}'] \leq \delta/(4H^2)$.      $\square$

572    ## C.3   Proof of Lemma B.4

573    To prove this lemma, we need a lemma similar to Lemma C.1

574    **Lemma C.2.** On dataset $\mathcal{D}' = \{ \bar{s}_h^k, \bar{a}_h^k, \bar{r}_h^k \}_{k,h=1}^{K,H}$, we define the filtration

$$\bar{\mathcal{H}}_h^k = \sigma(\bar{s}_1^1, \bar{a}_1^1, \bar{r}_1^1, \bar{s}_2^1, \ldots, \bar{r}_H^1, \bar{s}_{H+1}^1; \bar{s}_1^2, \bar{a}_1^2, \bar{r}_1^2, \bar{s}_2^2, \ldots, \bar{r}_H^2, \bar{s}_{H+1}^2; \ldots; \bar{s}_1^k, \bar{a}_1^k, \bar{r}_1^k, \bar{s}_2^k, \ldots, \bar{r}_h^k, \bar{s}_{h+1}^k).$$

575    For any fixed function $f, f' : \mathcal{S} \to [0, L]$, we make the following definitions:

$$\bar{\eta}_h^k[f'] := \left( \bar{r}_h^k + f'(\bar{s}_{h+1}^k) \right)^2 - \left[ \mathbb{P}_h(\bar{r}_h + f')^2 \right] (\bar{s}_h^k, \bar{a}_h^k)$$

$$\bar{D}_h^k[f, f'] := 2\bar{\eta}_h^k[f'] \left( f(\bar{z}_h^k) - \mathcal{T}_{2,h} f'(\bar{z}_h^k) \right).$$

576    Then with probability at least $1 - \delta/(4H^2\mathcal{N}^2\mathcal{N}_b^2)$, the following inequality holds,

$$\sum_{k \in [K]} \bar{D}_h^k[f, f'] \leq (24L + 5)i'^2(\delta) + \frac{\sum_{k \in [K]} \left( f(\bar{z}_h^k) - \mathcal{T}_{2,h} f'(\bar{z}_h^k) \right)^2}{2},$$

577    where $i'(\delta) = \sqrt{4 \log \frac{\mathcal{N}\mathcal{N}_b H(2\log(4LK)+2)(\log(4L)+2)}{\delta}}$.

18

*Proof.* We use Lemma D.1, with the following conditions:

$\bar{D}_h^k[f, f']$ is adapted to the filtration $\bar{\mathcal{H}}_h^k$ and $\mathbb{E}\left[\bar{D}_h^k[f, f'] \mid \bar{\mathcal{H}}_h^{k-1}\right] = 0$.

$$\left|\bar{D}_h^k[f, f']\right| \le 2|\bar{\eta}_h^k| \max_z |f(z) - \mathcal{T}_{2,h} f'(z)| \le 4L^2 = M.$$

$$\sum_{k \in [K]} \mathbb{E}\left[\left(\bar{D}_h^k[f, f']\right)^2 \Big| \bar{z}_h^k\right] = \sum_{k \in [K]} \mathbb{E}\left[4(\bar{\eta}_h^k[f'])^2|\bar{z}_h^k\right]\left(f(\bar{z}_h^k) - \mathcal{T}_{2,h} f'(\bar{z}_h^k)\right)^2 \le (4L^2 K)^2 = V^2.$$

On the other hand,

$$\sum_{k \in [K]} \mathbb{E}\left[\left(\bar{D}_h^k[f, f']\right)^2 \Big| \bar{z}_h^k\right] = \sum_{k \in [K]} \mathbb{E}\left[4\left(\bar{\eta}_h^k[f']\right)^2 \Big| \bar{z}_h^k\right]\left(f(\bar{z}_h^k) - \mathcal{T}_h f'(\bar{z}_h^k)\right)^2$$

$$\le 8L^4 \sum_{k \in [K]} \left(f(\bar{z}_h^k) - \mathcal{T}_{2,h} f'(\bar{z}_h^k)\right)^2.$$

Then using Lemma D.1 with $v = 1$, $m = 1$, we have:

$$\sum_{k \in [K]} 2\bar{\eta}_h^k[f']\left(f(\bar{z}_h^k) - \mathcal{T}_{2,h} f'(\bar{z}_h^k)\right) \le i'(\delta)\sqrt{2(2 \cdot 8L^4) \sum_{k \in [K]}\left(f(\bar{z}_h^k) - \mathcal{T}_{2,h} f'(\bar{z}_h^k)\right)^2}$$

$$+ \frac{2}{3}i'^2(\delta) + \frac{4}{3}i'^2(\delta) \cdot 4L^2$$

$$\le (20L^4 + 5)i'^2(\delta) + \frac{\sum_{k \in [K]}\left(f(\bar{z}_h^k) - \mathcal{T}_h f'(\bar{z}_h^k)\right)^2}{2}$$

$\square$

*Proof of Lemma B.4.* Let $(\beta'_{2,h})^2 = 2(20L^4 + 5)i'^2(\delta) + 16KL\epsilon$. We define the event $\mathcal{E}'_{2,h} := \left\{\sum_{k \in [K]}\left(\bar{g}'_h(\bar{z}_h^k) - \widetilde{g}'_h(\bar{z}_h^k)\right)^2 > (\beta'_{2,h})^2\right\}$. The following inequality will be useful in our proof.

$$\sum_{k \in [K]}\left(\bar{g}'_h(\bar{z}_h^k) - \widetilde{g}'_h(\bar{z}_h^k)\right)^2 = \sum_{k \in [K]}\left[\left((\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2 - \bar{g}'_h(\bar{z}_h^k)\right) + \left(\widetilde{g}'_h(\bar{z}_h^k) - (\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2\right)\right]^2$$

$$= \sum_{k \in [K]}\left((\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2 - \bar{g}'_h(\bar{z}_h^k)\right)^2 + \sum_{k \in [K]}\left(\widetilde{g}'_h(\bar{z}_h^k) - (\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2\right)^2$$

$$+ 2\sum_{k \in [K]}\left((\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2 - \bar{g}'_h(\bar{z}_h^k)\right)\left(\widetilde{g}'_h(\bar{z}_h^k) - (\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2\right)$$

$$\le 2\sum_{k \in [K]}\left((\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2 - \bar{g}'_h(\bar{z}_h^k)\right)^2$$

$$+ 2\sum_{k \in [K]}\left((\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2 - \bar{g}'_h(\bar{z}_h^k)\right)\left(\widetilde{g}'_h(\bar{z}_h^k) - (\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2\right)$$

$$\le 2\sum_{k \in [K]}\left((\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2 - \bar{g}'_h(\bar{z}_h^k)\right)\left(\widetilde{g}'_h(\bar{z}_h^k) - \bar{g}'_h(\bar{z}_h^k)\right). \quad \text{(C.6)}$$

Here we use our choice of $\widetilde{g}'_h$, i.e. $\widetilde{g}'_h = \operatorname{argmin}_{g_h \in \mathcal{F}_h} \sum_{k \in [K]}\left(g_h(\bar{s}_h^k, \bar{a}_h^k) - (\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2\right)^2$.

Next, we will use Lemma C.2. For any fixed $h$, let $f = \widetilde{g}'_h \in \mathcal{F}_h$, $f' = \widehat{f}'_{h+1} = \{\widetilde{f} - \epsilon\}_{[0, H-h+1]}$, where $\widetilde{f} = \widehat{f}'_h - b'_h \in \mathcal{F}_h - \mathcal{W}$. Following the construction in Lemma C.1, we define

$$\bar{\eta}_h^k[f'] := \left(\bar{r}_h^k + f'(\bar{s}_{h+1}^k)\right)^2 - \left[\mathbb{P}_h(\bar{r}_h + f')^2\right](\bar{s}_h^k, \bar{a}_h^k)$$
$$\text{and } \bar{D}_h^k[f, f'] := 2\bar{\eta}_h^k[f']\left(f(\bar{z}_h^k) - \mathcal{T}_{2,h} f'(\bar{z}_h^k)\right).$$

Due to the result of Lemma C.2, taking a union bound, we have with probability at least $1 - \delta/(4H^2)$, the following inequality holds,

$$\sum_{k \in [K]} \bar{D}_h^k[f, f'] \le (20L^4 + 5)i^2(\delta) + \frac{\sum_{k \in [K]}\left(f(\bar{z}_h^k) - \mathcal{T}_{2,h} f'(\bar{z}_h^k)\right)^2}{2}. \quad \text{(C.7)}$$

19

589  Therefore, with probability at least $1 - \delta/(4H^2)$, we have

$$2 \sum_{k \in [K]} \left( (\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2 - \bar{g}'_h(\bar{z}_h^k) \right) \left( \widetilde{g}'_h(\bar{z}_h^k) - \bar{g}'_h(\bar{z}_h^k) \right)$$

$$= 2 \sum_{k \in [K]} \left( (\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2 - \mathcal{T}_{2,h} \widehat{f}'_{h+1}(\bar{z}_h^k) \right) \left( \widetilde{g}'_h(\bar{z}_h^k) - \bar{g}'_h(\bar{z}_h^k) \right)$$

$$\qquad + 2 \sum_{k \in [K]} \left( \mathcal{T}_{2,h} \widehat{f}'_{h+1}(\bar{z}_h^k) - \bar{g}'_h(\bar{z}_h^k) \right) \left( \widetilde{g}'_h(\bar{z}_h^k) - \bar{g}'_h(\bar{z}_h^k) \right)$$

$$\leq 2 \sum_{k \in [K]} \left( (\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2 - \mathcal{T}_{2,h} \widehat{f}'_{h+1}(\bar{z}_h^k) \right) \left( \widetilde{g}'_h(\bar{z}_h^k) - \bar{g}'_h(\bar{z}_h^k) \right) + 4KL\epsilon$$

$$\leq (20L^4 + 5)i'^2(\delta) + 4KL\epsilon + \frac{\sum_{k \in [K]} \left( \bar{g}'_h(\bar{z}_h^k) - \mathcal{T}_{2,h} \widehat{f}'_{h+1}(\bar{z}_h^k) \right)^2}{2}$$

$$\leq (20L^4 + 5)i'^2(\delta) + 8KL\epsilon + \frac{\sum_{k \in [K]} \left( \bar{g}'_h(\bar{z}_h^k) - \widetilde{g}'_h(\bar{z}_h^k) \right)^2}{2}$$

$$\leq \frac{(\beta'_{2,h})^2}{2} + \frac{\sum_{k \in [K]} \left( \bar{g}'_h(\bar{z}_h^k) - \widetilde{g}'_h(\bar{z}_h^k) \right)^2}{2}.$$

590  Here the second inequality holds because of the Bellman completeness assumption (Assumption 3.1).
591  The third inequality arises from (C.7). The last inequality holds due to the choice of

$$\beta'_{2,h} = \sqrt{2(20L^4 + 5)i'^2(\delta) + 16KL\epsilon} = \widetilde{O}(\sqrt{\log \mathcal{N}\mathcal{N}_b} H^2).$$

592  But conditioned on the event $\mathcal{E}'_{2,h}$, we have

$$\sum_{k \in [K]} \left( (\bar{r}_h^k + \widehat{f}'_{h+1}(\bar{s}_{h+1}^k))^2 - \bar{g}'_h(\bar{z}_h^k) \right) \left( \widetilde{g}'_h(\bar{z}_h^k) - \bar{g}'_h(\bar{z}_h^k) \right)$$

$$\geq \sum_{k \in [K]} \left( \bar{g}'_h(\bar{z}_h^k) - \widetilde{g}'_h(\bar{z}_h^k) \right)^2$$

$$> \frac{(\beta'_{2,h})^2}{2} + \frac{\sum_{k \in [K]} \left( \bar{g}'_h(\bar{z}_h^k) - \widetilde{g}'_h(\bar{z}_h^k) \right)^2}{2}.$$

593  Here we use (C.6). Thus, we have $\mathbb{P}[\mathcal{E}'_{2,h}] \leq \delta/(4H^2)$. $\qquad\qquad\square$

594  ## C.4  Proof of Lemma B.5

595  *Proof of Lemma B.5.* We write $\mathbb{B}_h(s,a) = \widetilde{g}'_h(s,a) - \left( \widetilde{f}'_h(s,a) \right)^2$. We first bound the difference
596  between $\mathbb{B}_h(s,a)$ and $[\mathrm{Var}_h \widehat{f}'_{h+1}](s,a)$. By the definition of conditional variance, we have

$$\left| \mathbb{B}_h(s,a) - [\mathrm{Var}_h \widehat{f}'_{h+1}](s,a) \right| \leq \left| \widetilde{g}'_h(s,a) - \mathcal{T}_{2,h} \widehat{f}'_{h+1}(s,a) \right| + \left| \left( \widetilde{f}'_h(s,a) \right)^2 - \left( \mathcal{T}_h \widehat{f}'_{h+1}(s,a) \right)^2 \right|,$$

597  where we use our definition of Bellman operators. By the Bellman completeness assumption, there
598  exists $\bar{f}'_h \in \mathcal{F}_h$, $\bar{g}'_h \in \mathcal{F}_h$, such that $\left| \bar{f}'_h(s,a) - \mathcal{T}_h \widehat{f}'_{h+1}(s,a) \right| \leq \epsilon$, $\left| \bar{g}'_h(s,a) - \mathcal{T}_{2,h} \widehat{f}'_{h+1}(s,a) \right| \leq \epsilon$
599  for all $(s,a)$. Then by Lemma B.3 we can see that with probability at least $1 - \delta/(4H^2)$, the following
600  inequality holds

$$\sum_{k \in [K]} \left( \bar{f}'_h(\bar{z}_h^k) - \widetilde{f}'_h(\bar{z}_h^k) \right)^2 \leq (\beta'_{1,h})^2. \tag{C.8}$$

601  Similarly, for the second order term, using Lemma B.4, we can see that with probability at least
602  $1 - \delta/(4H^2)$, the following inequality holds

$$\sum_{k \in [K]} \left( \bar{g}'_h(\bar{z}_h^k) - \widetilde{g}'_h(\bar{z}_h^k) \right)^2 \leq (\beta'_{2,h})^2. \tag{C.9}$$

20

After taking a union bound, we have that with probability at least $1 - \delta/(2H)$, (C.8) and (C.9) hold for all $h \in [H]$ simultaneously. Under this high-probability event, we have

$$
\left| \widetilde{g}'_h(s,a) - \mathcal{T}_{2,h}\widehat{f}'_{h+1}(s,a) \right| + \left| \left( \widetilde{f}_h(s,a) \right)^2 - \left( \mathcal{T}_h\widehat{f}'_{h+1}(s,a) \right)^2 \right|
$$

$$
\leq \epsilon + |\widetilde{g}'_h(s,a) - \bar{g}'_h(s,a)| + O(H) \cdot \left| \bar{f}'_h(s,a) - \bar{f}'_h(s,a) + \epsilon \right|
$$

$$
\leq O(H) \cdot \epsilon + \frac{|\widetilde{g}'_h(s,a) - \bar{g}'_h(s,a)|}{\sqrt{\sum_{k\in[K]} \left( \bar{g}'_h(\bar{z}_h^k) - \widetilde{g}'_h(\bar{z}_h^k) \right)^2 + \lambda}} \cdot \sqrt{\sum_{k\in[K]} \left( \bar{g}'_h(\bar{z}_h^k) - \widetilde{g}'_h(\bar{z}_h^k) \right)^2 + \lambda}
$$

$$
+ O(H) \cdot \frac{|\widetilde{f}'_h(s,a) - \bar{f}'_h(s,a)|}{\sqrt{\sum_{k\in[K]} \left( \bar{f}'_h(\bar{z}_h^k) - \widetilde{f}'_h(\bar{z}_h^k) \right)^2 + \lambda}} \cdot \sqrt{\sum_{k\in[K]} \left( \bar{f}'_h(\bar{z}_h^k) - \widetilde{f}'_h(\bar{z}_h^k) \right)^2 + \lambda}
$$

$$
\leq O(H) \cdot \epsilon + \frac{|\widetilde{g}'_h(s,a) - \bar{g}'_h(s,a)|}{\sqrt{\sum_{k\in[K]} \left( \bar{g}'_h(\bar{z}_h^k) - \widetilde{g}'_h(\bar{z}_h^k) \right)^2 + \lambda}} \cdot \sqrt{(\beta'_{2,h})^2 + \lambda}
$$

$$
+ O(H) \cdot \frac{|\widetilde{f}'_h(s,a) - \bar{f}'_h(s,a)|}{\sqrt{\sum_{k\in[K]} \left( \bar{f}'_h(\bar{z}_h^k) - \widetilde{f}'_h(\bar{z}_h^k) \right)^2 + \lambda}} \cdot \sqrt{(\beta'_{1,h})^2 + \lambda}
$$

$$
\leq \widetilde{O}(\sqrt{\log \mathcal{N}\mathcal{N}_b}H^2) \cdot D_{\mathcal{F}_h}(z, \mathcal{D}'_h, 1)
$$

$$
\leq \widetilde{O}\left( \frac{\sqrt{\log \mathcal{N}\mathcal{N}_b}H^2}{\sqrt{K\kappa}} \right),
$$

where the first inequality holds due to the completeness assumption. The third inequality holds because of (C.8) and (C.9). The fourth inequality holds due to the definition of $D^2$-divergence. Finally we use Lemma B.1.

To further bound the difference between $\left[ \mathrm{Var}_h\widehat{f}'_{h+1} \right](s,a)$ and $[\mathrm{Var}_h V^*_{h+1}](s,a)$ under the event when (C.8) and (C.9) hold for all $h \in [H]$ simultaneously, we first prove $\|\widehat{f}'_{h+1} - V^*_{h+1}\|_\infty \leq \widetilde{O}\left( \frac{\sqrt{\log \mathcal{N}\mathcal{N}_b}H^3}{\sqrt{K\kappa}} \right)$ by induction.

At stage $H + 1$, $\widehat{f}'_{H+1} = V^*_{H+1} = 0$, the inequality holds naturally. At stage $H$, we have

$$
\begin{aligned}
Q^*_H(s,a) &= \mathcal{T}_H V^*_{H+1}(s,a) \\
&= \mathcal{T}_H \widehat{f}'_{H+1}(s,a) \\
&\geq \widetilde{f}'_H(s,a) - |\mathcal{T}_H \widehat{f}'_{H+1}(s,a) - \widetilde{f}'_H(s,a)| \\
&\geq \widetilde{f}'_H(s,a) - (\epsilon + |\bar{f}'_H(s,a) - \widetilde{f}'_H(s,a)|) \\
&\geq \widetilde{f}'_H(s,a) - b'_H(s,a) - \epsilon \\
&= \widehat{f}'_H(s,a).
\end{aligned}
$$

Here we use the definition of $\widehat{f}'_H$ in Algorithm 1 Line 6. Lemma B.3 shows

$$
\sum_{k\in[K]} \left( \bar{f}'_h(\bar{z}_h^k) - \widetilde{f}'_h(\bar{z}_h^k) \right)^2 \leq (\beta'_{1,h})^2.
$$

Then the fifth inequality follows the bonus oracle assumption (Definition 4.1). Therefore, $V^*_H(s) \geq \widehat{f}'_H(s)$ for all $s \in \mathcal{S}$.

21

We also have

$$V_H^*(s) - \widehat{f}_H'(s) = \langle Q_H^*(s,\cdot) - \widehat{f}_H'(s,\cdot), \pi^*(\cdot|s)\rangle_{\mathcal{A}} + \langle \widehat{f}_H'(s,\cdot), \pi_H^*(\cdot|s) - \widehat{\pi}_H(\cdot|s)\rangle_{\mathcal{A}}$$

$$\leq \langle Q_H^*(s,\cdot) - \widehat{f}_H'(s,\cdot), \pi^*(\cdot|s)\rangle_{\mathcal{A}}$$

$$= \langle \mathcal{T}_H V_H^*(s,\cdot) - \widetilde{f}_H'(s,\cdot) + b_H'(s,\cdot), \pi^*(\cdot|s)\rangle_{\mathcal{A}}$$

$$= \langle \mathcal{T}_H \widehat{f}_{H+1}'(s,\cdot) - \widetilde{f}_H'(s,\cdot) + b_H'(s,\cdot), \pi^*(\cdot|s)\rangle_{\mathcal{A}}$$

$$\quad + \langle \mathcal{T}_H V_{H+1}^*(s,\cdot) - \mathcal{T}_H \widehat{f}_{H+1}'(s,\cdot), \pi_H^*(\cdot|s)\rangle_{\mathcal{A}}$$

$$\leq 2\langle b_H'(s,\cdot), \pi_H^*(\cdot,s)\rangle_{\mathcal{A}} + \epsilon$$

$$\leq \widetilde{O}\left(\frac{\sqrt{\log \mathcal{N}\mathcal{N}_b}H}{\sqrt{K\kappa}}\right),$$

where the second inequality holds due to the selection of policy $\widehat{\pi}_H$. In the fifth inequality, we use the Bellman completeness assumption:

$$\left|\bar{f}_H'(z) - \mathcal{T}_H \widehat{f}_{H+1}'(z)\right| \leq \epsilon, \forall z \in \mathcal{S} \times \mathcal{A}.$$

and Lemma B.3. The last inequality holds because of Definition 4.1 and Lemma B.1.

We define $R_h = \widetilde{O}\left(\frac{\sqrt{\log \mathcal{N}\mathcal{N}_b}H}{\sqrt{K\kappa}}\right) \cdot (H - h + 1)$. To use the method of induction, we define the induction assumption as follows: Suppose with the probability of $1 - \delta_{h+1}$, the event $\mathcal{E}_{h+1} = \{0 \leq V_{h+1}^*(s) - \widehat{V}_{h+1}'(s) \leq R_{h+1}\}$ holds. Then we want to prove that with the probability of $1 - \delta_h$, the event $\mathcal{E}_h = \{0 \leq V_h^*(s) - \widehat{V}_h'(s) \leq R_h\}$ holds.

Conditioned on the event $\mathcal{E}_{h+1}$, using similar argument to stage $H$, we have

$$Q_h^*(s,a) = \mathcal{T}_h V_{h+1}^*(s,a)$$

$$\geq \mathcal{T}_h \widehat{f}_{h+1}'(s,a)$$

$$\geq \widetilde{f}_h'(s,a) - |\mathcal{T}_h \widehat{f}_{h+1}'(s,a) - \widetilde{f}_h'(s,a)|$$

$$\geq \widetilde{f}_h'(s,a) - \left(\epsilon + |\bar{f}_h'(s,a) - \widetilde{f}_h'(s,a)|\right)$$

$$\geq \widetilde{f}_h'(s,a) - b_h'(s,a)$$

$$= \widehat{f}_h'(s,a).$$

Therefore, $V_h^*(\cdot) \geq \widehat{f}_h'(\cdot)$.

On the other hand, similar to the case at stage $H$, we have with probability at least $1 - \delta_h - \delta/(2H^2)$,

$$V_h^*(s) - \widehat{f}_h'(s) = \langle Q_h^*(s,\cdot) - \widehat{f}_h'(s,\cdot), \pi_h^*(\cdot|s)\rangle_{\mathcal{A}} + \langle \widehat{f}_h'(s,\cdot), \pi_h^*(\cdot|s) - \widehat{\pi}_h(\cdot|s)\rangle_{\mathcal{A}}$$

$$\leq \langle Q_h^*(s,\cdot) - \widehat{f}_h'(s,\cdot), \pi_h^*(\cdot|s)\rangle_{\mathcal{A}}$$

$$= \langle \mathcal{T}_h V_{h+1}^*(s,\cdot) - \widetilde{f}_h'(s,\cdot) + b_h'(s,a), \pi_h^*(\cdot|s)\rangle_{\mathcal{A}}$$

$$= \langle \mathcal{T}_h \widehat{f}_{h+1}'(s,\cdot) - \widetilde{f}_h'(s,\cdot) + b_h'(s,a), \pi_h^*(\cdot|s)\rangle_{\mathcal{A}}$$

$$\quad + \langle \mathcal{T}_h V_{h+1}^*(s,\cdot) - \mathcal{T}_h \widehat{f}_{h+1}'(s,\cdot), \pi_h^*(\cdot|s)\rangle_{\mathcal{A}}$$

$$\leq 2\langle b_h'(s,\cdot), \pi_h^*(\cdot,s)\rangle_{\mathcal{A}} + \epsilon + R_{h+1}$$

$$\leq R_{h+1} + \widetilde{O}\left(\frac{\sqrt{\log \mathcal{N}\mathcal{N}_b}H}{\sqrt{K\kappa}}\right)$$

$$\leq \widetilde{O}\left(\frac{\sqrt{\log \mathcal{N}\mathcal{N}_b}H}{\sqrt{K\kappa}}\right) \cdot (H - h + 1) = R_h.$$

The induction shows we can choose $\delta_h = h\delta/(2H^2)$. Thus, taking a union bound over all $h \in [H]$, we prove that with probability at least $1 - \delta/2$, the following inequality

$$0 \leq V_{h+1}^*(\cdot) - \widehat{f}_{h+1}'(\cdot) \leq \widetilde{O}\left(\frac{\sqrt{\log \mathcal{N}\mathcal{N}_b}H^2}{\sqrt{K\kappa}}\right) \tag{C.10}$$

628 holds for all $h \in [H]$ simultaneously.

629 Conditioned on this event, we can further bound the difference between $[\mathrm{Var}_h \widehat{f}'_{h+1}](s, a)$ and
630 $[\mathrm{Var}_h V^*_{h+1}](s, a)$.

$$
\begin{aligned}
\left| [\mathrm{Var}_h \widehat{f}'_{h+1}](s, a) - [\mathrm{Var}_h V^*_{h+1}](s, a) \right| &\leq \left| [\mathbb{P}_h \widehat{f}'^2_{h+1}](s, a) - [\mathbb{P}_h V^{*2}_{h+1}](s, a) \right| \\
&\quad + \left| \left( [\mathbb{P}_h \widehat{f}'_{h+1}](s, a) \right)^2 - \left( [\mathbb{P}_h V^*_{h+1}](s, a) \right)^2 \right| \\
&\leq O(H) \cdot \left\| V^*_{h+1} - \widehat{f}'_{h+1} \right\|_\infty \\
&\leq \widetilde{O} \left( \frac{\sqrt{\log \mathcal{N} \mathcal{N}_b} H^3}{\sqrt{K\kappa}} \right).
\end{aligned}
$$

631 The last inequality arises from (C.10). Therefore, for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, we have

$$
\begin{aligned}
\left| \mathbb{B}_h(s, a) - [\mathrm{Var}_h V^*_{h+1}](s, a) \right| &\leq \left| \mathbb{B}_h(s, a) - [\mathrm{Var}_h \widehat{f}'_{h+1}](s, a) \right| \\
&\quad + \left| [\mathrm{Var}_h \widehat{f}'_{h+1}](s, a) - [\mathrm{Var}_h V^*_{h+1}](s, a) \right| \\
&\leq \widetilde{O} \left( \frac{\sqrt{\log \mathcal{N} \mathcal{N}_b} H^3}{\sqrt{K\kappa}} \right).
\end{aligned}
$$

632 Thus, for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, we have

$$
\mathbb{B}_h(s, a) - \widetilde{O} \left( \frac{\sqrt{\log \mathcal{N} \mathcal{N}_b} H^3}{\sqrt{K\kappa}} \right) \leq [\mathrm{Var}_h V^*_{h+1}](s, a).
$$

633 Finally, using the fact that the function $\max\{1, \cdot\}$ is increasing and nonexpansive, we finish the proof
634 of Lemma B.5, which is

$$
[\mathbb{V}_h V^*_{h+1}](s, a) - \widetilde{O} \left( \frac{\sqrt{\log \mathcal{N} \mathcal{N}_b} H^3}{\sqrt{K\kappa}} \right) \leq \widehat{\sigma}^2_h(s, a) \leq [\mathbb{V}_h V^*_{h+1}](s, a).
$$

635 $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad\square$

## C.5 Proof of Lemma B.6

637 To prove this result, we need the following lemmas.

638 **Lemma C.3.** Based on the dataset $\mathcal{D} = \{s^k_h, a^k_h, r^k_h\}^{K,H}_{k,h=1}$, we define the filtration $\mathcal{H}^k_h =$
639 $\sigma(s^1_1, a^1_1, r^1_1, s^1_2, \ldots, r^1_H, s^1_{H+1}; x^2_1, a^2_1, r^2_1, s^2_2, \ldots, r^2_H, s^2_{H+1}; \cdots s^k_1, a^k_1, r^k_1, s^k_2, \ldots, r^k_h, s^k_{h+1})$. For
640 any fixed function $f, f' : \mathcal{S} \to\in [0, L]$, we define the following random variables:

$$
\eta^k_h := V^*_{h+1}(s^k_{h+1}) - [\mathbb{P}_h V^*_{h+1}](s^k_h, a^k_h)
$$

$$
D^s_h[f, f'] := 2 \frac{\eta^k_h}{(\widehat{\sigma}_h(z^k_h))^2} \left( f(z^k_h) - f'(z^k_h) \right),
$$

641 Suppose the variance function $\widehat{\sigma}_h$ satisfies the inequality in Lemma B.5, where

$$
[\mathbb{V}_h V^*_{h+1}](s, a) - \widetilde{O} \left( \frac{\sqrt{\log \mathcal{N} \mathcal{N}_b} H^3}{\sqrt{K\kappa}} \right) \leq \widehat{\sigma}^2_h(s, a) \leq [\mathbb{V}_h V^*_{h+1}](s, a).
$$

642 Then, with probability at least $1 - \delta/(4H^2 \mathcal{N}^2)$, the following inequality holds,

$$
\sum_{k \in [K]} D^k_h[f, f'] \leq \frac{4}{3} v(\delta) \sqrt{\lambda} + \sqrt{2} v(\delta) + 30 v^2(\delta) + \frac{1}{v(\delta)} \sqrt{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z^k_h))^2} \left( f(z^k_h) - f'(z^k_h) \right)^2 + \lambda},
$$

643 where $v(\delta) = \sqrt{2 \log \frac{H \mathcal{N}(2 \log(18LT) + 2)(\log(18L) + 2)}{\delta_h}}$.

23

644 *Proof.* We use Lemma D.1, with the following conditions:

$$D_h^k[f, f'] \text{ is adapted to the filtration } \mathcal{H}_h^k \text{ and } \mathbb{E}\left[D_h^k[f, f'] \mid \mathcal{H}_h^{k-1}\right] = 0.$$

$$\left|D_h^k[f, f']\right| \le 2\left|\eta_h^k\right| \max_z |f(z) - f'(z)| \le 8LH = M.$$

$$\sum_{k \in [K]} \mathbb{E}\left[\left(D_h^k[f, f']\right)^2 \Big| z_h^k\right] = 4 \sum_{k \in [K]} \frac{\mathbb{E}\left[(\eta_h^k)^2 \big| z_h^k\right]}{(\widehat{\sigma}_h(z_h^k))^4} \left(f(z_h^k) - f'(z_h^k)\right).$$

645 On the other hand,

$$\sum_{k \in [K]} \mathbb{E}\left[\left(D_h^k[f, f']\right)^2 \Big| z_h^k\right] = 4 \sum_{k \in [K]} \frac{\mathbb{E}\left[(\eta_h^k)^2 \big| z_h^k\right]}{(\widehat{\sigma}_h(z_h^k))^4} \left(f(z_h^k) - f'(z_h^k)\right)^2$$

$$\le 8 \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left(f(z_h^k) - f'(z_h^k)\right)^2,$$

646 where the last inequality holds because of the inequality in Lemma B.5:

$$\mathbb{E}\left[(\eta_h^k)^2 | z_h^k\right] = [\text{Var}_h V_{h+1}^*](s_h^k, a_h^k)$$

$$\le [\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k)$$

$$\le \left(\widehat{\sigma}_h(z_h^k)\right)^2 + \widetilde{O}\left(\frac{\sqrt{\log \mathcal{N}\mathcal{N}_b}H^3}{\sqrt{K\kappa}}\right)$$

$$\le 2\left(\widehat{\sigma}_h(z_h^k)\right)^2,$$

647 where we use the requirement that $K \ge \widetilde{\Omega}\left(\frac{\log \mathcal{N}\mathcal{N}_b H^6}{\kappa}\right)$.

648 Moreover, for any $k \in [K]$,

$$\left|D_h^k[f, f']\right| \le 2\left|\frac{\eta_h^k}{(\widehat{\sigma}_h(z_h^k))^2}\right| \left|f(z_h^k) - f'(z_h^k)\right|$$

$$\le 4H \sqrt{D_{\mathcal{F}_h}^2(z_h^k, \mathcal{D}_h, \widehat{\sigma}_h)\left(\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2}\left(f(z_h^k) - f'(z_h^k)\right)^2 + \lambda\right)}$$

$$\le \widetilde{O}\left(\frac{4H^2}{\sqrt{K\kappa}}\right)\sqrt{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2}\left(f(z_h^k) - f'(z_h^k)\right)^2 + \lambda}$$

$$\le \frac{1}{\upsilon(\delta)}\sqrt{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2}\left(f(z_h^k) - f'(z_h^k)\right)^2 + \lambda}.$$

649 The second inequality holds because of the definition of $D^2$ divergence (Definition 3.2). The
650 third inequality holds due to Lemma B.2. The last inequality holds because of the choice of
651 $K \ge \widetilde{\Omega}\left(\frac{\upsilon^2(\delta)H^4}{\kappa}\right)$.

652 Then using Lemma D.1 with $\upsilon = 1$, $m = 1$, we have

$$\sum_{k \in [K]} 2\frac{\eta_h^k}{(\widehat{\sigma}_h(z_h^k))^2}\left(f(z_h^k) - f'(z_h^k)\right) \le \upsilon(\delta)\sqrt{16 \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2}\left(f(z_h^k) - f'(z_h^k)\right)^2 + 2} + \frac{2}{3}\upsilon^2(\delta)$$

$$+ \frac{4}{3}\upsilon(\delta)\sqrt{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2}\left(f(z_h^k) - f'(z_h^k)\right)^2 + \lambda}$$

$$\le \frac{4}{3}\upsilon(\delta)\sqrt{\lambda} + \sqrt{2}\upsilon(\delta) + 30\upsilon^2(\delta)$$

$$+ \frac{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2}\left(f(z_h^k) - f'(z_h^k)\right)^2}{4}.$$

24

653 $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$

**Lemma C.4.** Based on the dataset $\mathcal{D} = \left\{ s_h^k, a_h^k, r_h^k \right\}_{k,h=1}^{K,H}$, we define the following filtration $\mathcal{H}_h^k = \sigma\left( s_1^1, a_1^1, r_1^1, s_2^1, \ldots, r_H^1, s_{H+1}^1; x_1^2, a_1^2, r_1^2, s_2^2, \ldots, r_H^2, s_{H+1}^2; \cdots s_1^k, a_1^k, r_1^k, s_2^k, \ldots, r_h^k, s_{h+1}^k \right)$. For any fixed functions $f, \widetilde{f} : \mathcal{S} \to [0, L]$ and $f' : \mathcal{S} \to [0, H]$, we define the following random variables

$$\xi_h^k[f'] := f'(s_{h+1}^k) - V_{h+1}^*(s_{h+1}^k) - \left[ \mathbb{P}_h (f' - V_{h+1}^*) \right](s_h^k, a_h^k),$$

$$\Delta_h^k \left[ f, \widetilde{f}, f' \right] := 2 \frac{\xi_h^k[f']}{(\widehat{\sigma}_h(z_h^k))^2} \left( f(z_h^k) - \widetilde{f}(z_h^k) \right),$$

Suppose the variance function $\widehat{\sigma}_h$ satisfies the inequality in Lemma B.5, where

$$[\mathbb{V}_h V_{h+1}^*](s, a) - \widetilde{O}\left( \frac{\sqrt{\log \mathcal{N} \mathcal{N}_b} H^3}{\sqrt{K\kappa}} \right) \le \widehat{\sigma}_h^2(s, a) \le [\mathbb{V}_h V_{h+1}^*](s, a).$$

Then, with probability at least $1 - \delta/(4H^2 \mathcal{N}^3 \mathcal{N}_b)$, the following inequality holds,

$$\sum_{k \in [K]} \Delta_h^k[f, \widetilde{f}, f'] \le \left( \frac{4}{3} \iota(\delta) \sqrt{\lambda} + \sqrt{2} \iota(\delta) \right) \| f' - V_{h+1}^* \|_\infty^2 + \frac{2}{3} \iota^2(\delta) / \log \mathcal{N}_b$$

$$+ 30 \iota^2(\delta) \| f' - V_{h+1}^* \|_\infty^2 + \frac{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} (f(z_h^k) - f'(z_h^k))^2}{4}.$$

where $\iota(\delta) = \sqrt{3 \log \frac{H \mathcal{N} \mathcal{N}_b (2 \log(18LT)+2)(\log(18L)+2)}{\delta}}$.

*Proof.* $\Delta_h^k[f, \widetilde{f}, f']$ is adapted to the filtration $\mathcal{H}_h^k$ and $\mathbb{E}\left[ \Delta_h^k[f, \widetilde{f}, f'] \mid \mathcal{H}_h^{k-1} \right] = 0$. We also have

$$\sum_{k \in [K]} \mathbb{E}\left[ (\Delta_h^k[f, \widetilde{f}, f'])^2 \Big| z_h^k \right] = 4 \sum_{k \in [K]} \frac{\mathbb{E}\left[ (\xi_h^k[f'])^2 \Big| z_h^k \right]}{(\widehat{\sigma}_h(z_h^k))^4} \left( f(z_h^k) - f'(z_h^k) \right)^2$$

$$\le 8 \sum_{k \in [K]} \frac{\| f' - V_{h+1}^* \|_\infty^2}{(\widehat{\sigma}_h(z_h^k))^2} \left( f(z_h^k) - f'(z_h^k) \right)^2.$$

Moreover, for any $k \in [K]$,

$$\left| \Delta_h^k[f, \widetilde{f}, f'] \right| \le 2 \left| \frac{\xi_h^k[f']}{(\widehat{\sigma}_h(z_h^k))^2} \right| \left| f(z_h^k) - f'(z_h^k) \right|$$

$$\le 4 \| f' - V_{h+1}^* \|_\infty \sqrt{D_{\mathcal{F}_h}^2(z_h^k, \mathcal{D}_h, \widehat{\sigma}_h) \left( \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( f(z_h^k) - f'(z_h^k) \right)^2 + \lambda \right)}$$

$$\le \widetilde{O}\left( \frac{H}{\sqrt{K\kappa}} \right) \cdot \| f' - V_{h+1}^* \|_\infty \sqrt{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( f(z_h^k) - f'(z_h^k) \right)^2 + \lambda}$$

$$\le \frac{\| f' - V_{h+1}^* \|_\infty}{\iota(\delta)} \sqrt{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( f(z_h^k) - f'(z_h^k) \right)^2 + \lambda}$$

The second inequality holds because of the definition of $D^2$ divergence (Definition 3.2). The third inequality holds due to Lemma B.2. The last inequality holds because of the choice of $K \ge \widetilde{\Omega}\left( \frac{\iota^2(\delta) H^4}{\kappa} \right)$.

25

Then using Lemma D.1 with $v = 1$, $m = 1/\log \mathcal{N}_b$, we have

$$\sum_{k \in [K]} 2 \frac{\xi_h^k[f']}{(\widehat{\sigma}_h(z_h^k))^2} \left( f(z_h^k) - \widetilde{f}(z_h^k) \right) \leq \iota(\delta) \sqrt{8 \sum_{k \in [K]} \frac{\|f' - V_{h+1}^*\|_\infty^2}{(\widehat{\sigma}_h(z_h^k))^2} \left( f(z_h^k) - f'(z_h^k) \right)^2 + 2}$$

$$+ \frac{2}{3} \iota^2(\delta) / \log \mathcal{N}_b + \frac{4}{3} \iota(\delta) \|f' - V_{h+1}^*\|_\infty \sqrt{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} (f(z_h^k) - f'(z_h^k))^2 + \lambda}$$

$$\leq \left( \frac{4}{3} \iota(\delta) \sqrt{\lambda} + \sqrt{2} \iota(\delta) \right) \|f' - V_{h+1}^*\|_\infty^2 + \frac{2}{3} \iota^2(\delta) / \log \mathcal{N}_b + 30 \iota^2(\delta) \|f' - V_{h+1}^*\|_\infty^2$$

$$+ \frac{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( f(z_h^k) - f'(z_h^k) \right)^2}{4}.$$

$\square$

*Proof of Lemma B.6.* We define the event $\mathcal{E}_h := \left\{ \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}(z_h^k))^2} \left( \bar{f}_h(z_h^k) - \widetilde{f}_h(z_h^k) \right)^2 > (\beta_h)^2 \right\}$.

The following inequality will be useful in our proof.

$$\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}(z_h^k))^2} \left( \bar{f}_h(z_h^k) - \widetilde{f}_h(z_h^k) \right)^2$$

$$= \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}(z_h^k))^2} \left[ \left( r_h^k + \widehat{f}_{h+1}(s_{h+1}^k) - \bar{f}_h(z_h^k) \right) + \left( \widetilde{f}_h(z_h^k) - r_h^k - \widehat{f}_{h+1}(s_{h+1}^k) \right) \right]^2$$

$$= \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}(z_h^k))^2} \left( r_h^k + \widehat{f}_{h+1}(s_{h+1}^k) - \bar{f}_h(z_h^k) \right)^2 + \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}(z_h^k))^2} \left( \widetilde{f}_h(z_h^k) - r_h^k - \widehat{f}_{h+1}(s_{h+1}^k) \right)^2$$

$$+ 2 \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}(z_h^k))^2} \left( r_h^k + \widehat{f}_{h+1}(s_{h+1}^k) - \bar{f}_h(z_h^k) \right) \left( \widetilde{f}_h(z_h^k) - r_h^k - \widehat{f}_{h+1}(s_{h+1}^k) \right)$$

$$\leq 2 \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}(z_h^k))^2} \left( r_h^k + \widehat{f}_{h+1}(s_{h+1}^k) - \bar{f}_h(z_h^k) \right)^2$$

$$+ 2 \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}(z_h^k))^2} \left( r_h^k + \widehat{f}_{h+1}(s_{h+1}^k) - \bar{f}_h(z_h^k) \right) \left( \widetilde{f}_h(z_h^k) - r_h^k - \widehat{f}_{h+1}(s_{h+1}^k) \right)$$

$$\leq 2 \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}(z_h^k))^2} \left( r_h^k + \widehat{f}_{h+1}(s_{h+1}^k) - \bar{f}_h(z_h^k) \right) \left( \widetilde{f}_h(z_h^k) - \bar{f}_h(z_h^k) \right). \tag{C.11}$$

In the third inequality, we use our choice of $\widetilde{f}_h$ in Algorithm 1 Line 10,

$$\widetilde{f}_h = \operatorname*{argmin}_{f_h \in \mathcal{F}_h} \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}(z_h^k))^2} \left( f_h(s_h^k, a_h^k) - r_h^k - \widehat{f}_{h+1}(s_{h+1}^k) \right)^2.$$

We first use Lemma C.3 at stage $H$. Let $f = \widetilde{f}_H \in \mathcal{F}_H$, $f' = \bar{f}_H \in \mathcal{F}_H$. We define

$$\eta_H^k := V_{H+1}^*(s_{H+1}^k) - [\mathbb{P}_H V_{H+1}^*](z_H^k)$$

$$D_H^k[f, f'] := 2 \frac{\eta_H^k}{(\widehat{\sigma}_H(z_H^k))^2} \left( f(z_H^k) - f'(z_H^k) \right).$$

Taking a union bound, we have with probability at least $1 - \delta/(4H^2)$, the following inequality holds,

$$\sum_{k \in [K]} 2 \frac{\eta_H^k}{(\widehat{\sigma}_H (z_H^k))^2} \left( \widetilde{f}(z_H^k) - \bar{f}(z_H^k) \right) \leq \frac{4}{3} v(\delta) \sqrt{\lambda} + \sqrt{2} v(\delta) + 30 v^2(\delta)$$

$$+ \frac{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_H(z_H^k))^2} (\widetilde{f}(z_H^k) - \bar{f}(z_H^k))^2}{4}. \tag{C.12}$$

Then we use Lemma C.4 at stage $H$. Let $f = \widetilde{f}_H \in \mathcal{F}_H$, $\widetilde{f} = \bar{f}_H \in \mathcal{F}_H$, $f' = \widehat{f}_{H+1} = 0$. We define:

$$\xi_H^k[f'] := f'(s_{H+1}^k) - V_{H+1}^*(s_{H+1}^k) - [\mathbb{P}_H(f' - V_{H+1}^*)](z_H^k)$$

$$\Delta_H^k[f, \widetilde{f}, f'] := 2\frac{\xi_H^k[f']}{(\widehat{\sigma}_H(z_H^k))^2}\left(f(z_H^k) - \widetilde{f}(z_H^k)\right).$$

Therefore, taking a union bound, we have with probability at least $1 - \delta/(4H^2)$, we have

$$\sum_{k\in[K]} 2\frac{\xi_H^k[\widehat{f}_{H+1}]}{(\widehat{\sigma}_H(z_H^k))^2}\left(\widetilde{f}_H(z_H^k) - \bar{f}_H(z_H^k)\right) \le \left(\frac{4}{3}\iota(\delta)\sqrt{\lambda} + \sqrt{2}\iota(\delta)\right)\|f' - V_{H+1}^*\|_\infty^2$$

$$+ \frac{2}{3}\iota^2(\delta)/\sqrt{\log\mathcal{N}_b} + 30\iota^2(\delta)\|\widehat{f}_{H+1} - V_{H+1}^*\|_\infty^2 + \frac{\sum_{k\in[K]}\frac{1}{(\widehat{\sigma}_H(z_H^k))^2}\left(\widetilde{f}_H(z_H^k) - \bar{f}_H(z_H^k)\right)^2}{4}.$$

$$\tag{C.13}$$

Combining (C.12) and (C.13), with probability at least $1 - \delta/(2H^2)$, we have

$$2\sum_{k\in[K]} \frac{1}{(\widehat{\sigma}_H(z_H^k))^2}\left(r_H^k + \widehat{f}_{H+1}(s_{H+1}^k) - \bar{f}_H(z_H^k)\right)\left(\widetilde{f}_H(z_H^k) - \bar{f}_H(z_H^k)\right)$$

$$= 2\sum_{k\in[K]} \frac{1}{(\widehat{\sigma}_H(z_H^k))^2}\left(r_H^k + \widehat{f}_{H+1}(s_{H+1}^k) - \mathcal{T}_H\widehat{f}_{H+1}(z_H^k)\right)\left(\widetilde{f}_H(z_H^k) - \bar{f}_H(z_H^k)\right)$$

$$+ 2\sum_{k\in[K]} \frac{1}{(\widehat{\sigma}_H(z_H^k))^2}\left(\mathcal{T}_H\widehat{f}_{H+1}(z_H^k) - \bar{f}_H(z_H^k)\right)\left(\widetilde{f}_H(z_H^k) - \bar{f}_H(z_H^k)\right)$$

$$\le 2\sum_{k\in[K]} \frac{1}{(\widehat{\sigma}_H(z_H^k))^2}\left(r_H^k + \widehat{f}_{H+1}(s_{H+1}^k) - \mathcal{T}_H\widehat{f}_{H+1}(z_H^k)\right)\left(\widetilde{f}_H(z_H^k) - \bar{f}_H(z_H^k)\right) + 4KL\epsilon$$

$$\le \frac{4}{3}v(\delta)\sqrt{\lambda} + \sqrt{2}v(\delta) + 30v^2(\delta) + \left(\frac{4}{3}\iota(\delta)\sqrt{\lambda} + \sqrt{2}\iota(\delta)\right)\|f' - V_{H+1}^*\|_\infty^2 + \frac{2}{3}\iota^2(\delta)/\log\mathcal{N}_b$$

$$+ 30\iota^2(\delta)\|\widehat{f}_{H+1} - V_{H+1}^*\|_\infty^2 + 8KL\epsilon + \frac{\sum_{k\in[K]}\frac{1}{(\widehat{\sigma}_H(z_H^k))^2}\left(\bar{f}_H(z_H^k) - \widetilde{f}_H(z_H^k)\right)^2}{2}$$

$$\le \frac{(\beta_H)^2}{2} + \frac{\sum_{k\in[K]}\frac{1}{(\widehat{\sigma}_H(z_H^k))^2}\left(\bar{f}_H(z_H^k) - \widetilde{f}_H(z_H^k)\right)^2}{2}.$$

In the last inequality, we use that fact $\widehat{f}_{H+1} = V_{H+1}^* = 0$ and our choice of $\beta_H$.

$$\beta_H = \sqrt{2\left(\frac{4}{3}v(\delta)\sqrt{\lambda} + \sqrt{2}v(\delta) + 30v^2(\delta) + \frac{2}{3}\iota^2(\delta)/\log\mathcal{N}_b + 8KL\epsilon\right)}$$

$$= \widetilde{O}(\sqrt{\log\mathcal{N}}).$$

But conditioned on the event $\mathcal{E}_H$, we have

$$\sum_{k\in[K]} \frac{1}{(\widehat{\sigma}_H(z_H^k))^2}\left(r_H^k + \widehat{f}_{H+1}(s_{H+1}^k) - \bar{f}_H(z_H^k)\right)\left(\widetilde{f}_H(z_H^k) - \bar{f}_H(z_H^k)\right)$$

$$\ge \sum_{k\in[K]} \frac{1}{(\widehat{\sigma}_H(z_H^k))^2}\left(\bar{f}_H(z_H^k) - \widetilde{f}_H(z_H^k)\right)^2$$

$$> \frac{(\beta_H)^2}{2} + \frac{\sum_{k\in[K]}\frac{1}{(\widehat{\sigma}(z_H^k))^2}\left(\bar{f}_H(z_H^k) - \widetilde{f}_H(z_H^k)\right)^2}{2}.$$

Here we use (C.11). We finally prove that $\mathbb{P}[\mathcal{E}_H] \ge 1 - \delta/2H^2$.

Suppose the event $\mathcal{E}_H$ holds, we can prove the following result.

$$
\begin{aligned}
Q_H^*(s,a) &= \mathcal{T}_H V_{H+1}^*(s,a) \\
&= \mathcal{T}_H \widehat{f}_{H+1}(s,a) \\
&\geq \widetilde{f}_H(s,a) - \left| \mathcal{T}_H \widehat{f}_{H+1}(s,a) - \widetilde{f}_H(s,a) \right| \\
&\geq \widetilde{f}_H(s,a) - \left( \epsilon + |\bar{f}_H(s,a) - \widetilde{f}_H(s,a)| \right) \\
&\geq \widetilde{f}_H(s,a) - b_H(s,a) - \epsilon \\
&= \widehat{f}_H(s,a).
\end{aligned}
$$

Here we use the definition of $\widehat{f}_H$ in Algorithm 1 Line 12. Lemma B.6 shows

$$
\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( \bar{f}_h(z_h^k) - \widetilde{f}_h(z_h^k) \right)^2 \leq (\beta_h)^2.
$$

Then the fifth inequality follows the bonus oracle assumption (Definition 4.1). Therefore, $V_H^*(s) \geq \widehat{f}_H(s)$ for all $s \in \mathcal{S}$.

We also have

$$
\begin{aligned}
V_H^*(s) - \widehat{f}_H(s) &= \langle Q_H^*(s,\cdot) - \widehat{f}_H(s,\cdot), \pi^*(\cdot|s) \rangle_{\mathcal{A}} + \langle \widehat{f}_H(s,\cdot), \pi_H^*(\cdot|s) - \widehat{\pi}_H(\cdot|s) \rangle_{\mathcal{A}} \\
&\leq \langle Q_H^*(s,\cdot) - \widehat{f}_H(s,\cdot), \pi^*(\cdot|s) \rangle_{\mathcal{A}} \\
&= \langle \mathcal{T}_H V_H^*(s,\cdot) - \widetilde{f}_H(s,\cdot) + b_H(s,a), \pi^*(\cdot|s) \rangle_{\mathcal{A}} \\
&= \langle \mathcal{T}_H \widehat{f}_{H+1}(s,\cdot) - \widetilde{f}_H(s,\cdot) + b_H(s,a), \pi^*(\cdot|s) \rangle_{\mathcal{A}} \\
&\quad + \langle \mathcal{T}_H V_{H+1}^*(s,\cdot) - \mathcal{T}_H \widehat{f}_{H+1}(s,\cdot), \pi_H^*(\cdot|s) \rangle_{\mathcal{A}} \\
&\leq 2 \langle b_H(s,\cdot), \pi_H^*(\cdot|s) \rangle_{\mathcal{A}} + \epsilon \\
&\leq \widetilde{O} \left( \frac{\sqrt{\log \mathcal{N}} H^2}{\sqrt{K \kappa}} \right).
\end{aligned}
$$

Here the second inequality holds because of the definition of $\widehat{\pi}$. The fifth inequality holds due to we the Bellman completeness assumption (Assumption 3.1):

$$
\left| \bar{f}_H(z) - \mathcal{T}_H \widehat{f}_{H+1}(z) \right| \leq \epsilon, \forall z \in \mathcal{S} \times \mathcal{A}.
$$

We also use Definition 4.1 and Lemma B.2.

Then we do the induction step. Let $R_h = \widetilde{O} \left( \frac{\sqrt{\log \mathcal{N}} H^2}{\sqrt{K \kappa}} \right) \cdot (H - h + 1)$, $\delta_h = h\delta/(4H^2)$. We define another event $\mathcal{E}_h^{\text{ind}}$ for induction.

$$
\mathcal{E}_h^{\text{ind}} = \{ 0 \leq V_h^*(s) - \widehat{f}_h(s) \leq R_h, \forall s \in \mathcal{S} \}.
$$

The above analysis shows that $\mathcal{E}_H \subseteq \mathcal{E}_H^{\text{ind}}$ and $\mathbb{P}[\mathcal{E}_H] \geq 1 - 2\delta_H$. Moreover, $\mathbb{P}[\mathcal{E}_H^{\text{ind}}] \geq 1 - 2\delta_H$

We conduct the induction in the following way. At stage $h$, if $\mathbb{P}[\mathcal{E}_{h+1}] \geq 1 - 2\delta_{h+1}$ and $\mathbb{P}[\mathcal{E}_{h+1}^{\text{ind}}] \geq 1 - 2\delta_{h+1}$, we prove that $\mathbb{P}[\mathcal{E}_h] \geq 1 - 2\delta_h$ and $\mathbb{P}[\mathcal{E}_h^{\text{ind}}] \geq 1 - 2\delta_h$.

Suppose at stage $h$, $\mathbb{P}[\mathcal{E}_{h+1}] \geq 1 - 2\delta_{h+1}$ and $\mathbb{P}[\mathcal{E}_{h+1}^{\text{ind}}] \geq 1 - 2\delta_{h+1}$. We first use Lemma C.3. Let $f = \widetilde{f}_h \in \mathcal{F}_h$, $f' = \bar{f}_h \in \mathcal{F}_h$. We define

$$
\eta_h^k := V_{h+1}^*(s_{h+1}^k) - [\mathbb{P}_h V_{h+1}^*](z_h^k)
$$

$$
D_h^k[f, f'] := 2 \frac{\eta_h^k}{(\widehat{\sigma}_h(z_h^k))^2} \left( f(z_h^k) - f'(z_h^k) \right).
$$

After taking a union bound, we have with probability at least $1 - \delta/(4H^2)$, the following inequality holds,

$$\sum_{k \in [K]} 2\frac{\eta_h^k}{(\widehat{\sigma}_h(z_h^k))^2} \left( \widetilde{f}(z_h^k) - \bar{f}(z_h^k) \right) \le \frac{4}{3}v(\delta)\sqrt{\lambda} + \sqrt{2}v(\delta) + 30v^2(\delta)$$

$$+ \frac{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( \widetilde{f}(z_h^k) - \bar{f}(z_h^k) \right)^2}{4}. \qquad (C.14)$$

Next, we use Lemma C.4 at stage $h$. Let $f = \widetilde{f}_h \in \mathcal{F}_h$, $\widetilde{f} = \bar{f}_h \in \mathcal{F}_h$, $f' = \widehat{f}_{h+1} = \{\widetilde{b}\}_{[0,H-h+1]}$, where $\widetilde{b} = \widetilde{f}_h - b_h \in \mathcal{F}_h - \mathcal{W}$. We define:

$$\xi_h^k[f'] := f'(s_{h+1}^k) - V_{h+1}^*(s_{h+1}^k) - \left[ \mathbb{P}_h(f' - V_{h+1}^*) \right](z_h^k)$$

$$\Delta_h^k[f, \widetilde{f}, f'] := 2\frac{\xi_h^k[f']}{(\widehat{\sigma}_h(z_h^k))^2} \left( f(z_h^k) - \widetilde{f}(z_h^k) \right),$$

After taking a union bound, we have with probability at least $1 - \delta/(4H^2)$, we have

$$\sum_{k \in [K]} 2\frac{\xi_h^k[\widehat{f}_{h+1}]}{(\widehat{\sigma}_h(z_h^k))^2} \left( \widetilde{f}_h(z_h^k) - \bar{f}_h(z_h^k) \right) \le \left( \frac{4}{3}\iota(\delta)\sqrt{\lambda} + \sqrt{2}\iota(\delta) \right) \|f' - V_{h+1}^*\|_\infty^2$$

$$+ \frac{2}{3}\iota^2(\delta)/\sqrt{\log \mathcal{N}_b} + 30\iota^2(\delta)\|\widehat{f}_{h+1} - V_{h+1}^*\|_\infty^2 + \frac{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2}(\widetilde{f}_h(z_h^k) - \bar{f}_h(z_h^k))^2}{4}. \qquad (C.15)$$

Let $U_h$ be the event that (C.14) and (C.15) holds simultaneously. On the event $U_h \cap \mathcal{E}_{h+1}^{\text{ind}}$, which satisfies $\mathbb{P}[U_h \cap \mathcal{E}_{h+1}^{\text{ind}}] \ge 1 - 2\delta_{h+1} - 2\delta/H^2 = 1 - 2\delta_h$, we have

$$2 \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( r_h^k + \widehat{f}_{h+1}(s_{h+1}^k) - \bar{f}_h(z_h^k) \right) \left( \widetilde{f}_h(z_h^k) - \bar{f}_h(z_h^k) \right)$$

$$= 2 \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( r_h^k + \widehat{f}_{h+1}(s_{h+1}^k) - \mathcal{T}_h \widehat{f}_{h+1}(z_h^k) \right) \left( \widetilde{f}_h(z_h^k) - \bar{f}_h(z_h^k) \right)$$

$$+ 2 \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( \mathcal{T}_h \widehat{f}_{h+1}(z_h^k) - \bar{f}_h(z_h^k) \right) \left( \widetilde{f}_h(z_h^k) - \bar{f}_h(z_h^k) \right)$$

$$\le 2 \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( r_h^k + \widehat{f}_{h+1}(s_{h+1}^k) - \mathcal{T}_h \widehat{f}_{h+1}(z_h^k) \right) \left( \widetilde{f}_h(z_h^k) - \bar{f}_h(z_h^k) \right) + 4KL\epsilon$$

$$\le \frac{4}{3}v(\delta)\sqrt{\lambda} + \sqrt{2}v(\delta) + 30v^2(\delta) + \left( \frac{4}{3}\iota(\delta)\sqrt{\lambda} + \sqrt{2}\iota(\delta) \right) \|\widehat{f}_{h+1} - V_{h+1}^*\|_\infty^2$$

$$+ \frac{2}{3}\iota^2(\delta)/\log \mathcal{N}_b + 30\iota^2(\delta)\|\widehat{f}_{h+1} - V_{h+1}^*\|_\infty^2 + 8KL\epsilon + \frac{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( \bar{f}_h(z_h^k) - \widetilde{f}_h(z_h^k) \right)^2}{2}$$

$$\le \frac{(\beta_h)^2}{2} + \frac{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}(z_h^k))^2} \left( \bar{f}_h(z_h^k) - \widetilde{f}_h(z_h^k) \right)^2}{2},$$

where the third inequality holds because of (C.14) and (C.15). The last inequality holds because on the event of $\mathcal{E}_{h+1}^{\text{ind}}$, $0 \le V_{h+1}^* - \widehat{f}_{h+1} \le R_{h+1} = \widetilde{O}\left( \frac{H^2}{\sqrt{K\kappa}} \right) \cdot (H - h)$ and the choice of $K \ge \widetilde{\Omega}\left( \frac{\iota(\delta)^2 H^6}{\kappa} \right)$. We also use our choice of $\beta_h = \widetilde{O}(\sqrt{\log \mathcal{N}})$.

29

However, on the event of $\mathcal{E}_h^c$, we have

$$2 \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}(z_h^k))^2} \left( r_h^k + \widehat{f}_{h+1}(s_{h+1}^k) - \bar{f}_h(z_h^k) \right) \left( \widetilde{f}_h(z_h^k) - \bar{f}_h(z_h^k) \right)$$

$$\geq \sum_{k \in [K]} \frac{1}{(\widehat{\sigma}(z_h^k))^2} \left( \bar{f}_h(z_h^k) - \widetilde{f}_h(z_h^k) \right)^2$$

$$> \frac{(\beta_h)^2}{2} + \frac{\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}(z_h^k))^2} \left( \bar{f}_h(z_h^k) - \widetilde{f}_h(z_h^k) \right)^2}{2}.$$

We conclude that $U_h \cap \mathcal{E}_{h+1}^{\text{ind}} \subseteq \mathcal{E}_h$, thus $\mathbb{P}[\mathcal{E}_h] \geq 1 - 2\delta_h$.

Next we prove $\mathbb{P}[\mathcal{E}_h^{\text{ind}}] \geq 1 - 2\delta_h$. Suppose the event $U_h \cap \mathcal{E}_{h+1}^{\text{ind}}$ holds, the above conclusion shows that

$$\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( \bar{f}_h(z_h^k) - \widetilde{f}_h(z_h^k) \right)^2 > (\beta_h)^2.$$

We can prove the following result.

$$\begin{aligned}
Q_h^*(s, a) = \mathcal{T}_h V_{h+1}^*(s, a) \\
\geq \mathcal{T}_h \widehat{f}_{h+1}(s, a) \\
\geq \widetilde{f}_h(s, a) - |\mathcal{T}_h \widehat{f}_{h+1}(s, a) - \widetilde{f}_h(s, a)| \\
\geq \widetilde{f}_h(s, a) - (\epsilon + |\bar{f}_h(s, a) - \widetilde{f}_h(s, a)|) \\
\geq \widetilde{f}_h(s, a) - b_h(s, a) \\
= \widehat{f}_h(s, a),
\end{aligned}$$

where the second inequality holds because the event $\mathcal{E}_h^{\text{ind}}$ holds. The fourth inequality holds because of the Bellman completeness assumption (Assumption 3.1). From Lemma B.6, we have

$$\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( \bar{f}_h(z_h^k) - \widetilde{f}_h(z_h^k) \right)^2 \leq (\beta_h)^2.$$

Then the fifth inequality holds due to the bonus oracle (Definition 4.1). Therefore, $V_h^*(s) \geq \widehat{f}_h(s)$ for all $s \in \mathcal{S}$.

We also have

$$\begin{aligned}
V_h^*(s) - \widehat{f}_h(s) &= \langle Q_h^*(s, \cdot) - \widehat{f}_h(s, \cdot), \pi^*(\cdot|s) \rangle_{\mathcal{A}} + \langle \widehat{f}_h(s, \cdot), \pi_h^*(\cdot|s) - \widehat{\pi}_h(\cdot|s) \rangle_{\mathcal{A}} \\
&\leq \langle Q_h^*(s, \cdot) - \widehat{f}_h(s, \cdot), \pi^*(\cdot|s) \rangle_{\mathcal{A}} \\
&= \langle \mathcal{T}_h V_h^*(s, \cdot) - \widetilde{f}_h(s, \cdot) + b_h(s, a), \pi^*(\cdot|s) \rangle_{\mathcal{A}} \\
&= \langle \mathcal{T}_h \widehat{f}_{h+1}(s, \cdot) - \widetilde{f}_h(s, \cdot) + b_h(s, a), \pi^*(\cdot|s) \rangle_{\mathcal{A}} \\
&\quad + \langle \mathcal{T}_h V_{h+1}^*(s, \cdot) - \mathcal{T}_h \widehat{f}_{h+1}(s, \cdot), \pi_h^*(\cdot|s) \rangle_{\mathcal{A}} \\
&\leq 2 \langle b_h(s, \cdot), \pi_h^*(\cdot, s) \rangle_{\mathcal{A}} + \epsilon + R_{h+1} \\
&\leq \widetilde{O}\left( \frac{\sqrt{\log \mathcal{N}} H^2}{\sqrt{K\kappa}} \right) \cdot (H - h + 1) = R_h.
\end{aligned}$$

The first equality holds because of our choice of the policy $\widehat{\pi}_h$. In the fifth inequality, we use Assumption 3.1

$$\left\| \bar{f}_h(\cdot, \cdot) - \mathcal{T}_h \widehat{f}_{h+1}(\cdot, \cdot) \right\|_{\infty} \leq \epsilon,$$

and the oracle of bonus function (Definition 4.1) with

$$\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left( \bar{f}_h(z_h^k) - \widetilde{f}_H(z_h^k) \right)^2 \leq (\beta_h)^2,$$

30

which holds by Lemma B.6. Therefore, we have $U_h \cap \mathcal{E}_{h+1}^{\text{ind}} \subseteq \mathcal{E}_h^{\text{ind}}$ and $\mathbb{P}[\mathcal{E}_h^{\text{ind}}] \geq 1 - 2\delta_h$. We also use the induction assumption. Thus we complete the proof of induction.

Finally, taking the union bound of all the $\mathcal{E}_h$, we get the result that with probability at least $1 - \delta/2$, the event $\cup_{h=1}^H \mathcal{E}_h$ holds, i.e for any $h \in [H]$ simultaneously, we have

$$\sum_{k \in [K]} \frac{1}{(\widehat{\sigma}_h(z_h^k))^2} \left(\bar{f}_h(z_h^k) - \widetilde{f}_h(z_h^k)\right)^2 \leq (\beta_h)^2.$$

Therefore, we complete the proof of Lemma B.6. $\square$

# D  Auxiliary lemmas

**Lemma D.1** (Agarwal et al. 2022)**.** Let $M > 0$, $V > v > 0$ be constants, and $\{x_i\}_{i \in [t]}$ be a stochastic process adapted to a filtration $\{\mathcal{H}_i\}_{i \in [t]}$. Suppose $\mathbb{E}[x_i | \mathcal{H}_{i-1}] = 0$, $|x_i| \leq M$ and $\sum_{i \in [t]} \mathbb{E}[x_i^2 | \mathcal{H}_{i-1}] \leq V^2$ almost surely. Then for any $\delta, \epsilon > 0$, let $\iota = \sqrt{\log \frac{(2\log(V/v)+2)\cdot(\log(M/m)+2)}{\delta}}$, we have

$$\mathbb{P}\left(\sum_{i \in [t]} x_i > \iota \sqrt{2\left(2\sum_{i \in [t]} \mathbb{E}[x_i^2 | \mathcal{H}_{i-1}] + v^2\right)} + \frac{2}{3}\iota^2 \left(2\max_{i \in [t]} |x_i| + m\right)\right) \leq \delta.$$

**Lemma D.2** (Regret Decomposition Property, Jin et al. 2021b)**.** Suppose the following inequality holds,

$$|\mathcal{T}_h \widehat{f}_{h+1}(z) - \widetilde{f}_h(z)| \leq b_h(z), \forall z = (s,a) \in \mathcal{S} \times \mathcal{A}, \forall h \in [H],$$

the regret of Algorithm 1 can be bounded as

$$V_1^*(s) - V_1^{\widehat{\pi}}(s) \leq 2\sum_{h=1}^H \mathbb{E}_{\pi^*}\left[b_h(s_h, a_h) \mid s_1 = s\right].$$

Here $\mathbb{E}_{\pi^*}$ is with respect to the trajectory induced by $\pi^*$ in the underlying MDP.

**Lemma D.3** (Azuma-Hoeffding inequality, Cesa-Bianchi and Lugosi 2006)**.** Let $\{x_i\}_{i=1}^n$ be a martingale difference sequence with respect to a filtration $\{\mathcal{G}_i\}$ satisfying $|x_i| \leq M$ for some constant $M$, $x_i$ is $\mathcal{G}_{i+1}$-measurable, $\mathbb{E}[x_i | \mathcal{G}_i] = 0$. Then for any $0 < \delta < 1$, with probability at least $1 - \delta$, we have

$$\sum_{i=1}^n x_i \leq M\sqrt{2n\log(1/\delta)}.$$