
Nearly Minimax Optimal Reinforcement Learning for Linear Markov Decision Processes

Jiafan He¹ Heyang Zhao¹ Dongruo Zhou¹ Quanquan Gu¹

Abstract

We study reinforcement learning (RL) with linear function approximation. For episodic time-inhomogeneous linear Markov decision processes (linear MDPs) whose transition probability can be parameterized as a linear function of a given feature mapping, we propose the first computationally efficient algorithm that achieves the nearly minimax optimal regret $\tilde{O}(d\sqrt{H^3K})$, where d is the dimension of the feature mapping, H is the planning horizon, and K is the number of episodes. Our algorithm is based on a weighted linear regression scheme with a carefully designed weight, which depends on a new variance estimator that (1) directly estimates the variance of the *optimal* value function, (2) monotonically decreases with respect to the number of episodes to ensure a better estimation accuracy, and (3) uses a rare-switching policy to update the value function estimator to control the complexity of the estimated value function class. Our work provides a complete answer to optimal RL with linear MDPs, and the developed algorithm and theoretical tools may be of independent interest.

1 Introduction

How to make reinforcement learning (RL) efficient with large state and action spaces has been a central research problem in the RL community. A widely used approach is *function approximation*, which approximates the value function in RL with a predefined function class for efficient exploration and exploitation. Although the intuition is simple, some basic questions about the function approximation approach still remain open. For instance, what is the optimal sample complexity (or regret) for RL algorithms

with function approximation to find the optimal policy? Such optimal sample complexity results have been widely-studied and established for tabular RL methods (e.g., Azar et al. 2017; Zhang and Ji 2019; Zhang et al. 2020), but are still understudied for RL with function approximation.

Some recent works have studied the optimal regret results for a special class of MDPs called *linear mixture Markov decision processes (linear mixture MDPs)* (Jia et al., 2020; Ayoub et al., 2020; Zhou et al., 2021a; Zhou and Gu, 2022), which assume that the transition probability of the MDP is a *linear* combination of several base models. More specifically, Zhou et al. (2021a) proposed the near-optimal algorithm for time-inhomogeneous linear mixture MDPs. Zhou and Gu (2022) further proposed near-optimal horizon-free algorithm for time-homogeneous linear mixture MDPs under the assumption that the total reward is bounded by 1. However, the computational efficiency of their algorithms highly depends on the *value-targeted regression* procedure (Jia et al., 2020; Ayoub et al., 2020), which relies on an integration or sampling oracle of the individual base model. Such an integration or sampling oracle exists for some special linear mixture MDPs but can be computationally expensive or even intractable in the general case.

Another line of works studies the *linear Markov decision processes* (linear MDPs) (Yang and Wang, 2019; Jin et al., 2020), which assumes that the transition probability and the reward of the environment enjoys a compact low-rank representation. The most appealing feature of linear MDPs is that they can induce a linear structure of the value function for any policy, which makes sample-efficient RL possible. Meanwhile, the algorithms for linear MDPs directly approximate the value function itself, which is computationally more efficient than the algorithms for linear mixture MDPs. In particular, Yang and Wang (2019) first proposed a near-optimal RL algorithm with the access to a generative model, which can generate any number of samples for any given state-action pairs. Without accessing the generative model, Jin et al. (2020) proposed an LSVI-UCB algorithm based on the principle of optimism in the face of uncertainty and achieved $\tilde{O}(\sqrt{d^3H^4K})$ regret, where d is the dimension of a linear MDP, H is the planning horizon and K is the number of episodes. Nevertheless, their algorithm is not optimal since there exists an $O(\sqrt{dH})$ gap between

¹Department of Computer Science, University of California, Los Angeles, CA 90095, USA. Correspondence to: Quanquan Gu <qgu@cs.ucla.edu>.

their regret upper bound, and the lower bound $O(d\sqrt{H^3K})$ proved in Zhou et al. (2021a). Zanette et al. (2020b) studied a more general MDP class called *low Bellman error class*, which contains linear MDPs as a special case, and they proposed a computationally inefficient algorithm with a near-optimal regret.

Therefore, a natural question arises¹:

Can we design a computationally efficient algorithm that achieves the minimax optimality for linear MDPs?

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

- We propose an algorithm LSVI-UCB++ which attains a near-optimal regret $\tilde{O}(d\sqrt{H^3K})$ when K is large, which matches the lower bound (Zhou et al., 2021a) up to logarithmic factors. To the best of our knowledge, this is the first computationally efficient RL algorithm that is nearly minimax-optimal for linear MDPs.
- The first key component of our algorithm is a variance-aware weighted ridge regression scheme, which is firstly introduced to achieve nearly minimax optimal regret for linear mixture MDPs in Zhou et al. (2021a) and later improved in Zhou and Gu (2022) to achieve horizon-free regret. Such a component reduces the variance of the estimators in our algorithm, which leads to a \sqrt{H} improvement in the regret over Jin et al. (2020).
- To improve the d dependence, inspired by previous works for tabular RL (Azar et al., 2017), our algorithm utilizes a new strategy to estimate the variance of the estimated value function. Unlike the previous approach for linear mixture MDPs (Zhou et al., 2021a), our new estimator directly estimates the variance of the *true* value function and computes the difference between the variances of the true value function and the estimated one. Such a strategy allows the variance estimator to focus on a simpler function class that only includes the true value function, and therefore gives a tighter confidence set than that in Jin et al. (2020).
- To obtain a uniform variance upper bound, we construct our value function estimator as a monotonically decreasing estimator with a “rare-switching” update strategy,

¹We are aware of a recently published work (Hu et al., 2022), which claims to achieve the nearly minimax optimal regret for linear MDPs. However, a closer examination of their proof can find a technical error, which makes their result invalid. We will discuss it in more detail and show why our algorithm and proof can get around the issue in Appendix A. Using the techniques proposed by our paper, Hu et al. (2022) recently fixed the technical flaw by using the “rare-switching” update strategy and also abandoning the over-optimistic estimator. This is acknowledged in the updated arXiv version of Hu et al. (2022).

Model	Algorithm	Regret
Linear MDP	LSVI-UCB (Jin et al., 2020)	$\tilde{O}(\sqrt{d^3H^4K})$
	LSVI-UCB++ (Our work)	$\tilde{O}(d\sqrt{H^3K})$
Lower bound	Zhou et al. (2021a)	$\Omega(d\sqrt{H^3K})$

Table 1. Comparison of RL with linear function approximation in terms of regret guarantee.

which makes the estimated value function decrease with respect to the episodes and being updated rarely. Together with our new variance estimator, we can remove the additional \sqrt{d} dependency from the previous regret, which makes our algorithm nearly minimax optimal. Notably, our algorithm only needs to update the policy $O(\log K)$ times instead of K times, and therefore enjoys a low-switching cost.

For the ease of comparison, we summarize the regret bounds of our algorithm and previous algorithms for linear MDPs in Table 1.

Recently, an independent concurrent work (Agarwal et al., 2022) proposed a different algorithm that can also achieve near-optimal regret for linear MDPs. Their algorithm follows the algorithm design in Hu et al. (2022), which introduces an additional over-optimistic value function to construct a monotonic variance estimator, and a non-Markovian policy to fix the technical flaw in Hu et al. (2022). In contrast, our algorithm takes a neat approach and constructs the monotonic variance estimator with a simple “rare-switching” update strategy, which enjoys low-switching cost. Agarwal et al. (2022) also studied RL with nonlinear function approximation, which is beyond the scope of this work.

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 $\Sigma \in \mathbb{R}^{d \times d}$, we denote by $\|\mathbf{x}\|_2$ the Euclidean norm and $\|\mathbf{x}\|_\Sigma = \sqrt{\mathbf{x}^\top \Sigma \mathbf{x}}$. For two sequences $\{a_n\}$ and $\{b_n\}$, we write $a_n = O(b_n)$ if there exists an absolute constant C such that $a_n \leq Cb_n$, and we write $a_n = \Omega(b_n)$ if there exists an absolute constant C such that $a_n \geq Cb_n$. We use $\tilde{O}(\cdot)$ and $\tilde{\Omega}(\cdot)$ 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 $[n] = \{1, 2, \dots, n\}$ to denote the set of integers from 1 to n .

2 Related Work

Near-optimal tabular reinforcement learning There is a voluminous amount of works developing nearly minimax optimal algorithms for tabular MDPs under different settings (Azar et al., 2017; Zanette and Brunskill, 2019; Zhang and Ji, 2019; Simchowitz and Jamieson, 2019; Zhang et al., 2020; 2021a; He et al., 2021b). A key idea behind these works is to exploit the $O(H^2)$ total variance of the value functions for each episode (Azar et al., 2017; Jin et al., 2018). Azar et al. (2017) first proposed this idea to design Bernstein-type bonuses in tabular MDPs and provided an $\tilde{O}(\sqrt{H^2 SAK})$ regret upper bound, which matches the lower bound for the tabular setting. In their analysis, Azar et al. (2017) also introduced a new value function decomposition scheme which mainly focuses on the variances of the optimal value function rather than the estimated value function. Zhang et al. (2021a) further improved the dependence on H for the constant terms and achieved a nearly minimax optimal horizon-free regret (nearly independent of H) under the assumption that the total reward is bounded by 1. Our algorithm extends the idea of Bernstein-type bonuses and value function decomposition in Azar et al. (2017) to RL with linear function approximation.

Reinforcement Learning with Linear Function Approximation. There exists a large body of literature on RL with linear function approximation (Jiang et al., 2017; Dann et al., 2018; Yang and Wang, 2019; Jin et al., 2020; Wang et al., 2020; Du et al., 2019; Sun et al., 2019; Zanette et al., 2020a; Yang and Wang, 2020; Modi et al., 2020; Ayoub et al., 2020; Zhou et al., 2021a; He et al., 2021a; Zhou and Gu, 2022). All these works assume certain linear structures of the underlying MDP. The most related work to ours is initiated by Yang and Wang (2019), which assumes that the reward function and the transition probability are linear in the feature mapping $\phi(s, a)$ for each state-action pair (s, a) . Jin et al. (2020) further considered *Linear MDPs* and proposed LSVI-UCB which achieves an $\tilde{O}(\sqrt{d^3 H^4 K})$ regret bound. Zanette et al. (2020a) proposed a Thompson sampling based algorithm for linear MDPs, which attains a regret upper bound of order $\tilde{O}(\sqrt{d^4 H^5 K})$. Another popular MDP model for RL with linear function approximation is *linear mixture Markov Decision Processes* (Modi et al., 2020; Yang and Wang, 2020; Jia et al., 2020; Ayoub et al., 2020), or *Linear Kernel MDPs* (Zhou et al., 2021b), where the transition probability is a linear combination of several base models. For linear mixture MDPs, Zhou et al. (2021a) is the first to achieve a nearly minimax optimal regret bound. There are also works achieving horizon-free regret bounds for time-homogeneous linear mixture MDPs (Zhang et al., 2021b; Zhou and Gu, 2022). Compared with Zhou et al. (2021a), our algorithm is the first to achieve the near-optimality for linear MDPs.

3 Preliminaries

In this work, we consider the episodic Markov Decision Processes (MDP), where the MDP can be denoted by a tuple of $M(\mathcal{S}, \mathcal{A}, H, \{r_h\}_{h=1}^H, \{\mathbb{P}_h\}_{h=1}^H)$. Here, \mathcal{S} is the state space, \mathcal{A} is the finite action space, H is the length of each episode (i.e., planning horizon), $r_h : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]^2$ is the reward function at stage h and $\mathbb{P}_h(s'|s, a)$ is the transition probability function at stage h which denotes the probability for state s to transfer to next state s' with current action a . Following Jin et al. (2020), we assume that \mathcal{S} is a measurable space with possibly infinite number of states and \mathcal{A} is a finite set. A policy $\pi : \mathcal{S} \times [H] \rightarrow \mathcal{A}$ is a function that maps a state s and the stage number h to an action a . For any stage $h \in [H]$ and policy π , we define the value function $V_h^\pi(s)$ and the action-value function $Q_h^\pi(s, a)$ as follows:

$$\begin{aligned} Q_h^\pi(s, a) &= r_h(s, a) \\ &+ \mathbb{E} \left[\sum_{h'=h+1}^H r_{h'}(s_{h'}, a_{h'}) \mid s_h = s, a_h = a \right], \\ V_h^\pi(s) &= Q_h^\pi(s, \pi(s, h)), \end{aligned}$$

where $s_{h'+1} \sim \mathbb{P}_h(\cdot | s_{h'}, a_{h'})$ denotes the state at stage $h' + 1$ and $a_{h'} = \pi(s_{h'}, h')$ denotes the action at stage h' . Furthermore, we can 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 $Q_h^*(s, a) = \max_\pi Q_h^\pi(s, a)$. By this definition, the value function $V_h^\pi(s)$ and action-value function $Q_h^\pi(s, a)$ are bounded in $[0, H]$. For any function $V : \mathcal{S} \rightarrow \mathbb{R}$, we denote $[\mathbb{P}_h V](s, a) = \mathbb{E}_{s' \sim \mathbb{P}_h(\cdot | s, a)} V(s')$ and $[\mathbb{V}_h V](s, a) = [\mathbb{P}_h V^2](s, a) - ([\mathbb{P}_h V](s, a))^2$ for simplicity. Thus, for every stage $h \in [H]$ and policy π , we 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 $Q_h^*(s, a)$ and $V_h^*(s)$:

$$\begin{aligned} Q_h^\pi(s, a) &= r_h(s, a) + [\mathbb{P}_h V_{h+1}^\pi](s, a), \\ Q_h^*(s, a) &= r_h(s, a) + [\mathbb{P}_h V_{h+1}^*](s, a), \end{aligned}$$

where $V_{H+1}^\pi(s) = V_{H+1}^*(s) = 0$. At the beginning of each episode $k \in [K]$, the agent selects a policy π_k to be followed in this episode. At each stage $h \in [H]$, the agent first observes the current state s_h^k , chooses an action following the policy π_k and then observes the next state with $s_{h+1}^k \sim \mathbb{P}_h(\cdot | s_h^k, a_h^k)$. Based on these definitions, we further define the regret in the first K episodes as follows:

Definition 3.1. For any algorithm Alg, we define its regret on learning an MDP $M(\mathcal{S}, \mathcal{A}, H, r, \mathbb{P})$ in the first K

²In this work, we study the deterministic and known reward functions for simplicity, and it is not difficult to generalize our results to stochastic and unknown linear reward functions in (Jin et al., 2019), where $r_h(s, a) = \langle \phi(s, a), \boldsymbol{\mu}_h \rangle$.

episodes as the sum of the sub-optimality gaps for episode $k = 1, \dots, K$, i.e.,

$$\text{Regret}(K) = \sum_{k=1}^K V_1^*(s_1^k) - V_1^{\pi_k}(s_1^k),$$

where π_k is the agent's policy in the k -th episode.

Linear Markov Decision Process In this work, we focus on the linear Markov decision Process (Jin et al., 2020; Yang and Wang, 2019), which is formally defined as follows:

Definition 3.2. An MDP $\mathcal{M}(\mathcal{S}, \mathcal{A}, H, \{r_h\}_{h=1}^H, \{\mathbb{P}_h\}_{h=1}^H)$ is a linear MDP if for any stage $h \in [H]$, there exists an unknown measure $\theta_h(\cdot) : \mathcal{S} \rightarrow \mathbb{R}^d$ and a known feature mapping $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$, such that for each state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$ and state $s' \in \mathcal{S}$, we have

$$\mathbb{P}_h(s'|s, a) = \langle \phi(s, a), \theta_h(s') \rangle. \quad (3.1)$$

For simplicity, we assume that the norms of $\theta_n(\cdot)$ and $\phi(\cdot, \cdot)$ are upper bounded by $\|\phi(s, a)\|_2 \leq 1$ and $\|\theta_h(\mathcal{S})\|_2 \leq \sqrt{d}$. For linear MDPs, we have the following property:

Proposition 3.3 (Proposition 2.3, Jin et al. 2020). *For any policy π , there exist weights $\{\mathbf{w}_h^\pi\}_{h=1}^H$ such that for any state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$ and stage $h \in [H]$, we have $[\mathbb{P}V_{h+1}^\pi](s, a) = \langle \phi(s, a), \mathbf{w}_h^\pi \rangle$.*

4 The Proposed Algorithm

In this section, we propose a new algorithm LSVI-UCB++ to learn the linear MDPs (See Definition 3.2). The main algorithm is illustrated in Algorithm 1. In the sequel, we introduce the key ideas of the proposed algorithm one by one.

4.1 Weighted Ridge Regression

The basic framework of our algorithm follows the LSVI-UCB algorithm proposed by Jin et al. (2020). Based on Proposition 3.3 that the expected value function $[\mathbb{P}_h V_{h+1}^\pi](s, a) = \langle \phi(s, a), \mathbf{w}_h^\pi \rangle$, Algorithm 1 reduces the learning of the optimal action-value function into a series of linear regression problems. In order to have a good estimation for the vector \mathbf{w}_h^π and achieve the minimax-optimal regret result, Algorithm 1 adapts the weighted ridge regression method (Henderson, 1975), which was used in heteroscedastic linear bandits (Lattimore et al., 2015; Kirschner and Krause, 2018) and more recently RL with linear function approximation (Zhou et al., 2021a) for linear mixture MDPs. In detail, for each stage $h \in [H]$ and episode $k \in [K]$, we construct the estimator $\widehat{\mathbf{w}}_{k,h}$ by solving the following weighted ridge regression

$$\widehat{\mathbf{w}}_{k,h} \leftarrow \underset{\mathbf{w} \in \mathbb{R}^d}{\text{argmin}} \lambda \|\mathbf{w}\|_2^2$$

Algorithm 1 LSVI-UCB++

Require: Regularization parameter $\lambda > 0$, confidence radius $\beta, \bar{\beta}, \tilde{\beta}$

- 1: Initialize $k_{\text{last}} = 0$ and for each stage $h \in [H]$ set $\Sigma_{0,h}, \Sigma_{1,h} \leftarrow \lambda \mathbf{I}$
- 2: For each stage $h \in [H]$ and state-action $(s, a) \in \mathcal{S} \times \mathcal{A}$, set $Q_{0,h}(s, a) \leftarrow H, \check{Q}_{0,h}(s, a) \leftarrow 0$
- 3: **for** episodes $k = 1, \dots, K$ **do**
- 4: Received the initial state s_1^k .
- 5: **for** stage $h = H, \dots, 1$ **do**
- 6: $\widehat{\mathbf{w}}_{k,h} = \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) V_{k,h+1}(s_{h+1}^i)$
- 7: $\check{\mathbf{w}}_{k,h} = \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) \check{V}_{k,h+1}(s_{h+1}^i)$
- 8: **if** there exists a stage $h' \in [H]$ such that $\det(\Sigma_{k,h'}) \geq 2 \det(\Sigma_{k_{\text{last}},h'})$ **then**
- 9: $Q_{k,h}(s, a) = \min \left\{ r_h(s, a) + \widehat{\mathbf{w}}_{k,h}^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)}, Q_{k-1,h}(s, a), H \right\}$
- 10: $\check{Q}_{k,h}(s, a) = \max \left\{ r_h(s, a) + \check{\mathbf{w}}_{k,h}^\top \phi(s, a) - \bar{\beta} \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)}, \check{Q}_{k-1,h}(s, a), 0 \right\}$
- 11: Set the last updating episode $k_{\text{last}} = k$
- 12: **else**
- 13: $Q_{k,h}(s, a) = Q_{k-1,h}(s, a)$
- 14: $\check{Q}_{k,h}(s, a) = \check{Q}_{k-1,h}(s, a)$
- 15: **end if**
- 16: $V_{k,h}(s) = \max_a Q_{k,h}(s, a)$
- 17: $\check{V}_{k,h}(s) = \max_a \check{Q}_{k,h}(s, a)$
- 18: **end for**
- 19: **for** stage $h = 1, \dots, H$ **do**
- 20: Take action $a_h^k \leftarrow \text{argmax}_a Q_{k,h}(s_h^k, a)$
- 21: Set the estimated variance $\sigma_{k,h}$ as in (4.1)
- 22: $\bar{\sigma}_{k,h} \leftarrow \max \left\{ \sigma_{k,h}, H, 2d^3 H^2 \|\phi(s_h^k, a_h^k)\|_{\Sigma_{k,h}^{-1}}^{1/2} \right\}$
- 23: $\Sigma_{k+1,h} = \Sigma_{k,h} + \bar{\sigma}_{k,h}^{-2} \phi(s_h^k, a_h^k) \phi(s_h^k, a_h^k)^\top$
- 24: Receive next state s_{h+1}^k
- 25: **end for**
- 26: **end for**

$$+ \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} (\mathbf{w}^\top \phi(s_h^i, a_h^i) - V_{k,h+1}(s_{h+1}^i))^2.$$

Here, we take the inverse of the estimated variances $\sigma_{k,h}^2$ as the weights for the regression problem and set $\sigma_{k,h}$ as

$$\bar{\sigma}_{k,h} = \max \left\{ \sigma_{k,h}, H, 2d^3 H^2 \|\phi(s_h^k, a_h^k)\|_{\Sigma_{k,h}^{-1}}^{1/2} \right\}$$

in Line 22 of Algorithm 1, which depends on the uncertainty term $\|\phi(s_h^k, a_h^k)\|_{\Sigma_{k,h}^{-1}}$. Note that the uncertainty-dependent weight has also been used in He et al. (2022) to defend the adversarial corruption in the linear bandits problem. The reason why we want to use an uncertainty-dependent weight can be explained by the following lemma.

Lemma 4.1 (Theorem 4.3, Zhou and Gu 2022). *Let $\{\mathcal{G}_k\}_{k=1}^\infty$ be a filtration, and $\{\mathbf{x}_k, \eta_k\}_{k \geq 1}$ be a stochastic process such that $\mathbf{x}_k \in \mathbb{R}^d$ is \mathcal{G}_k -measurable and $\eta_k \in \mathbb{R}$ is \mathcal{G}_{k+1} -measurable. Let $L, \sigma > 0$, $\boldsymbol{\mu}^* \in \mathbb{R}^d$. For $k \geq 1$, let $y_k = \langle \boldsymbol{\mu}^*, \mathbf{x}_k \rangle + \eta_k$ and suppose that η_k, \mathbf{x}_k also satisfy*

$$\mathbb{E}[\eta_k | \mathcal{G}_k] = 0, \mathbb{E}[\eta_k^2 | \mathcal{G}_k] \leq \sigma^2, |\eta_k| \leq R, \|\mathbf{x}_k\|_2 \leq L.$$

For $k \geq 1$, let $\mathbf{Z}_k = \lambda \mathbf{I} + \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^\top$, $\mathbf{b}_k = \sum_{i=1}^k y_i \mathbf{x}_i$, $\boldsymbol{\mu}_k = \mathbf{Z}_k^{-1} \mathbf{b}_k$, and $\beta_k = \tilde{O}(\sigma \sqrt{d} + \max_{1 \leq i \leq k} |\eta_i| \min\{1, \|\mathbf{x}_i\|_{\mathbf{Z}_{i-1}^{-1}}\})$. Then, for any $0 < \delta < 1$, with probability at least $1 - \delta$, for all $k \in [K]$, we have

$$\|\sum_{i=1}^k \mathbf{x}_i \eta_i\|_{\mathbf{Z}_k^{-1}} \leq \beta_k, \|\boldsymbol{\mu}_k - \boldsymbol{\mu}^*\|_{\mathbf{Z}_k} \leq \beta_k + \sqrt{\lambda} \|\boldsymbol{\mu}^*\|_2.$$

By Lemma 4.1, one can easily verify that $|\langle \widehat{\mathbf{w}}_{k,h}, \phi(s, a) \rangle - \mathbb{P}_h V_{k,h+1}(s, a)| = O(\beta \|\Sigma_{k,h}^{-1/2} \phi(s, a)\|_2)$, where $\beta = \tilde{O}(\sqrt{d})$. Such an $\tilde{O}(\sqrt{d})$ dependence is similar to that in Zhou and Gu (2022), which allows our algorithm to use a tighter confidence set than Jin et al. (2020). Therefore, we can construct the optimistic value function $Q_{k,h}$ with the linear function and an additional exploration bonus term (Line 7 in Algorithm 1), i.e.,

$$Q_{k,h}(s, a) \approx r_h(s, a) + \widehat{\mathbf{w}}_{k,h}^\top \phi(s, a) + \beta \|\Sigma_{k,h}^{-1/2} \phi(s, a)\|_2.$$

With the help of the exploration bonus, we can show that the optimistic value function $Q_{k,h}(s, a)$ is an upper bound of the optimal value function $Q_h^*(s, a)$ and the summation of the sub-optimality gaps can be upper bounded by the summation of exploration bonus $\sum_{h=1}^H \sum_{k=1}^K \beta \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)}$. By adapting the weighted ridge regression, Zhou and Gu (2022) proposed HF-UCRL-VTR+, which is able to achieve a nearly minimax optimal regret for linear mixture MDPs. However, their algorithm and approach cannot be directly applied to linear MDPs, and we need to construct a pessimistic value function $\check{V}_{k,h}$ for the optimal value function $Q_h^*(s, a)$ to estimate the gap between $V_{k,h}(s)$ and $V_h^*(s)$, where we have $V_{k,h}(s) - V_h^*(s) \leq V_{k,h}(s) - \check{V}_{k,h}(s)$. Similar to the optimistic value function, we construct the vector $\check{\mathbf{w}}_{k,h}$ by solving the following weighted ridge regression,

$$\begin{aligned} \check{\mathbf{w}}_{k,h} \leftarrow \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \lambda \|\mathbf{w}\|_2^2 \\ + \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} (\mathbf{w}^\top \phi(s_h^i, a_h^i) - \check{V}_{k,h+1}(s_{h+1}^i))^2, \end{aligned}$$

and compute the pessimistic value function $\check{Q}_{k,h}$ as:

$$\begin{aligned} \check{Q}_{k,h}(s, a) \approx r_h(s, a) \\ + \check{\mathbf{w}}_{k,h}^\top \phi(s, a) - \bar{\beta} \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)}, \end{aligned}$$

where $\bar{\beta} = \tilde{O}(\sqrt{d^3 H^2})$. We can show that the pessimistic value function $\check{V}_{k,h}(s)$ is a lower bound for the optimal value function $V_h^*(s)$.

4.2 Variance Estimator

We compare our variance estimator and its counterparts in Zhou et al. (2021a). Zhou et al. (2021a) first introduced variance estimators into RL with linear function approximation. They studied linear mixture MDPs, and their algorithm estimates the variance of the *optimistic* value function $V_{k,h+1}(s)$ directly. In comparison, for linear MDPs, estimating the variance of the *optimistic* value function $V_{k,h+1}(s)$ will encounter the dependence issue, which is discussed in Jin et al. (2020) and will introduce an additional \sqrt{d} factor in the regret due to the covering-based decoupling argument. Inspired by the previous works (Azar et al., 2017; Hu et al., 2022), we decompose the *optimistic* value function $V_{k,h+1}(s)$ into the *optimal* value function $V_{h+1}^*(s)$ and the sub-optimality gap $V_{k,h+1}(s) - V_{h+1}^*(s)$, then estimate their variances $[\mathbb{V}_h V_{h+1}^*](s, a)$ and $[\mathbb{V}_h (V_{k,h+1} - V_{h+1}^*)](s, a)$ separately.

For the variance of *optimal* value function $[\mathbb{V}_h V_{h+1}^*](s, a)$, since neither the variance operator \mathbb{V}_h nor the optimal value function V_{h+1}^* is observable, Algorithm 1 takes several steps to estimate these two quantities. In detail, Algorithm 1 uses the optimistic value function $V_{k,h+1}$ to estimate the optimal value function V_{h+1}^* and introduce an error term $D_{k,h}$ to bound the difference between $\mathbb{V}_h V_{k,h+1}$ and $\mathbb{V}_h V_{h+1}^*$. For the variance operator, Algorithm 1 follows Zhou et al. (2021a) to write the variance as the difference between the second-order moment and the square of the first-order moment of $V_{k,h}$, which is upper bounded by the bonus term $E_{k,h}$. More specifically, the variance of function $V_{k,h}$ can be denoted by

$$[\mathbb{V}_h V_{k,h}](s, a) = [\mathbb{P}_h V_{k,h}^2](s, a) - ([\mathbb{P}_h V_{k,h}](s, a))^2.$$

According to the Proposition 3.3, the expectation $\mathbb{P}_h V_{k,h}(s, a)$ and $\mathbb{P}_h V_{k,h}^2(s, a)$ are linear in the feature mapping $\phi(s, a)$ and can be approximated as follows,

$$\begin{aligned} [\mathbb{V}_h V_{k,h}](s, a) &\approx \bar{\mathbb{V}}_h V_{k,h+1}(s_h^k, a_h^k) \\ &:= [\tilde{\mathbf{w}}_{k,h}^\top \phi(s_h^k, a_h^k)]_{[0, H^2]} - [\widehat{\mathbf{w}}_{k,h}^\top \phi(s_h^k, a_h^k)]_{[0, H]}^2, \end{aligned}$$

where

$$\begin{aligned} \tilde{\mathbf{w}}_{k,h} &:= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \lambda \|\mathbf{w}\|_2^2 \\ &+ \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} (\mathbf{w}^\top \phi(s_h^i, a_h^i) - V_{k,h+1}^2(s_{h+1}^i))^2 \end{aligned}$$

is the solution to the weighted ridge regression problem for the squared value function. To summarize, LSVI-UCB++

constructs the estimated variance $\sigma_{k,h}$ as follows:

$$\sigma_{k,h} = \sqrt{[\tilde{V}_{k,h} V_{k,h+1}](s_h^k, a_h^k) + E_{k,h} + D_{k,h} + H}, \quad (4.1)$$

where $E_{k,h}$ and $D_{k,h}$ are defined as follows

$$\begin{aligned} E_{k,h} &= \min \left\{ \tilde{\beta} \|\Sigma_{k,h}^{-1/2} \phi(s_h^k, a_h^k)\|_2, H^2 \right\} \\ &\quad + \min \left\{ 2H\tilde{\beta} \|\Sigma_{k,h}^{-1/2} \phi(s_h^k, a_h^k)\|_2, H^2 \right\}, \\ D_{k,h} &= \min \left\{ 4d^3 H^2 \left(\hat{\mathbf{w}}_{k,h}^\top \phi(s_h^k, a_h^k) - \check{\mathbf{w}}_{k,h}^\top \phi(s_h^k, a_h^k) \right) \right. \\ &\quad \left. + 2\tilde{\beta} \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)}, d^3 H^3 \right\}. \end{aligned}$$

Here $E_{k,h}$ is the error between the estimated variance and the true variance of $V_{k,h+1}$, and $D_{k,h}$ is the error between the variance of $V_{k,h+1}$ and the variance of the optimal value function V_h^* . For term $D_{k,h}$, we use the difference between the optimistic value function $V_{k,h}$ and the pessimistic value function $\check{V}_{k,h}$ to bound the difference between $V_{k,h}$ and V_h^* . More discussions on the decomposition and the variance estimator can be found in Section 6.

5 Main Results

In this section, we provide the regret bound for our algorithm LSVI-UCB++.

Theorem 5.1. *For any linear MDP M , if we set the parameters $\lambda = 1/H^2$ and confidence radius $\beta, \bar{\beta}, \tilde{\beta}$ as*

$$\begin{aligned} \beta &= O\left(H\sqrt{d\lambda} + \sqrt{d \log^2(1 + dKH/(\delta\lambda))}\right), \\ \bar{\beta} &= O\left(H\sqrt{d\lambda} + \sqrt{d^3 H^2 \log^2(dHK/(\delta\lambda))}\right), \\ \tilde{\beta} &= O\left(H^2\sqrt{d\lambda} + \sqrt{d^3 H^4 \log^2(dHK/(\delta\lambda))}\right), \end{aligned}$$

then with high probability of at least $1 - 7\delta$, the regret of LSVI-UCB++ is upper bounded as follows:

$$\text{Regret}(K) \leq \tilde{O}\left(d\sqrt{H^3 K} + d^7 H^8\right).$$

In addition, the number of updates for $Q_{k,h}, \check{Q}_{k,h}$ is upper bounded by $O(dH \log(1 + K/\lambda))$.

Remark 5.2. When the number of episodic K satisfies that K is large, the regret can be simplified as $\tilde{O}(d\sqrt{H^3 K})$. Compared with the lower bound $\Omega(d\sqrt{H^3 K})$ proved in Zhou et al. (2021a), our regret bound matches the lower bound up to logarithmic factors, which suggests that LSVI-UCB++ is near-optimal for linear MDPs.

Remark 5.3. For LSVI-UCB++, based on the optimistic property $Q_{k,h}(s, a) \geq V_h^*(s)$ and the pessimistic property $V_{k,h}^{\pi^k}(s) \geq \check{V}_{k,h}(s)$, the sub-optimality at episode k

is upper bounded by $V_{k,1}(s_1) - \check{V}_{k,1}(s_1)$. Thus, the total regret for the first K episodes can be roughly upper bounded by $\text{Regret}(K) \leq \sum_{k=1}^K (V_{k,1}(s_1) - \check{V}_{k,1}(s_1)) = O(d^2 \sqrt{H^5/K} + d^8 H^9/K)$. When the initial state s_1^k is fixed across all episodes $k \in [K]$, according to the monotonic property of the optimistic value function $V_{k,1}$ and the pessimistic value function $\check{V}_{k,1}$, the sub-optimality gap $(V_{k,1}(s_1) - \check{V}_{k,1}(s_1))$ is decreasing. As a result, the cumulative regret up to episode K satisfies $V_{K,1}(s_1) - \check{V}_{K,1}(s_1) \leq 1/K \times \sum_{k=1}^K (V_{k,1}(s_1) - \check{V}_{k,1}(s_1)) = O(1/\sqrt{K})$, which implies a Probably Approximately Correct (PAC) guarantee. Therefore, LSVI-UCB++ will converge to the optimal policy and enjoys an (ϵ, δ) -PAC guarantee with sample complexity $\tilde{O}(1/\epsilon^2)$ without any modification of the algorithm. In contrast, to obtain the (ϵ, δ) -PAC guarantee, the LSVI-UCB algorithm in Jin et al. (2020) needs to randomly select a policy uniformly from the previous K policies.

Computational Complexity As shown in Jin et al. (2020), the computational complexity of the original LSVI-UCB is $O(d^2 |\mathcal{A}| H K^2)$, where \mathcal{A} is a finite action space and $|\mathcal{A}|$ is the size of the action set. Compared with the LSVI-UCB algorithm, Algorithm 1 uses the ‘‘rare-switching’’ technique, where the algorithm only updates the estimated value functions if the determinant of the covariance matrix $\det(\Sigma_{i,h})$ doubles (Line 8). According to Lemma F.1, the number of episodes that triggers the updating criterion is at most $dH \log(1 + K/\lambda)$ and the action-value function $Q_{k,h}(s, a)$ can be represented as a minimum over $dH \log(1 + K/\lambda)$ quadratic functions. Therefore, given all previous optimistic weight vectors $\mathbf{w}_{i,h}$ and covariance matrices $\Sigma_{i,h}$, computing the optimistic value function $Q_{k,h}(s, a)$ needs $\tilde{O}(d^3 H)$ computational complexity. Thus, for each episode $k \in [K]$, calculating the value function $Q_{k,h}(s_h^k, a)$, choosing the action $a_h^k \leftarrow \arg \max_a Q_{k,h}(s_h^k, a)$ and estimating the variance $\bar{\sigma}_{k,h}$ will only lead to $\tilde{O}(d^3 H^2 |\mathcal{A}|)$ computational complexity.

For computing the linear regression weight vectors (Line 6 to Line 7), if the updating criterion is not triggered in episode k , then LSVI-UCB++ only needs to update the weight vectors $\hat{\mathbf{w}}_{k,h}$ and $\check{\mathbf{w}}_{k,h}$. Since the value functions $V_{k,h+1}$ and $\check{V}_{k,h+1}$ remain unchanged, we only need to compute the new terms $\sigma_{k-1,h}^{-2} \phi(s_h^{k-1}, a_h^{k-1}) V_{k,h+1}(s_h^{k-1})$ and $\sigma_{k-1,h}^{-2} \phi(s_h^{k-1}, a_h^{k-1}) \check{V}_{k,h+1}(s_h^{k-1})$, which has an $O(d^3 H |\mathcal{A}|)$ computational complexity. On the other hand, if the updating criterion is triggered in episode k , then LSVI-UCB++ needs to update the value function and recalculate the weight vectors $\hat{\mathbf{w}}_{k,h}, \check{\mathbf{w}}_{k,h}$, which incurs an $\tilde{O}(d^4 H^2 |\mathcal{A}| K)$ computational complexity. Combining the computational complexity for all episodes and noticing that the number of episodes that trigger the updating criterion

is at most $\tilde{O}(dH)$, the total computational complexity of LSVI-UCB++ is $\tilde{O}(d^4 H^3 |\mathcal{A}|K)$, which improves the original LSVI-UCB algorithm by a factor of K .

6 Overview of Key Techniques

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

6.1 Decompose $V_{k,h+1}$ to V_{h+1}^* and $V_{k,h+1} - V_{h+1}^*$

We start with estimating the Bellman backup $[\mathbb{P}_h V_{k,h+1}](s, a)$, which is the main difficulty in almost all existing analyses of algorithms for linear MDPs. According to Proposition 3.3, for any value function V and state-action pair (s, a) , the Bellman backup $[\mathbb{P}_h V](s, a)$ is always a linear function of the feature mapping $\phi(s, a)$ can be approximated as follows

$$\begin{aligned} & [\widehat{\mathbb{P}}_{k,h} V_{k,h+1}](s, a) \\ & \approx \phi(s, a)^\top \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \sigma_{i,h}^{-2} \phi(s_h^i, a_h^i) V_{k,h+1}(s_{h+1}^i), \end{aligned}$$

which utilizes all the past observations $\phi(s_h^i, a_h^i)$ and the associated values $V(s_{h+1}^i)$. In addition, the estimation error for this estimator can be measured by

$$\begin{aligned} & [\widehat{\mathbb{P}}_{k,h} V_{k,h+1}](s, a) - [\mathbb{P}_h V_{k,h+1}](s, a) \\ & \approx \phi(s, a)^\top \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \sigma_{i,h}^{-2} \phi(s_h^i, a_h^i) \eta_{i,h}(V_{k,h+1}), \quad (6.1) \end{aligned}$$

where $\eta_{i,h}(V) = V(s_{h+1}^i) - [\mathbb{P}_h V](s_{h+1}^i)$ denotes the stochastic transition noise at episode i with value function V . According to the Bernstein-type self-normalized martingale inequality (Lemma 4.1), the summation of stochastic noise can be bounded by a small value (e.g., $|\widehat{\mathbb{P}}_{k,h} V](s, a) - [\mathbb{P}_h V](s, a)| \leq \beta \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)}$). However, Jin et al. (2020) noticed that the estimation of the optimistic value function $V_{k,h+1}$ depends on the past observations $(s_h^k, a_h^k, s_{h+1}^k)$, which violates the conditional independence required by the martingale concentration inequality, i.e., $\mathbb{E}[\eta_{i,h}(V_{k,h+1})] \neq 0$.

To deal with this problem, Jin et al. (2020) applied the uniform convergence argument based on covering number for all possible value functions and introduce a fixed covering set to replace $V_{k,h}$ in their analysis. In detail, the function class considered in Jin et al. (2020) is denoted by

$$\mathcal{V} = \left\{ V \left| V(\cdot) = \max_a \min \left(H, \mathbf{w}_i^\top \phi(\cdot, a) + \beta \sqrt{\phi(\cdot, a)^\top \Sigma_i^{-1} \phi(\cdot, a)} \right), \|\mathbf{w}_i\| \leq L, \Sigma_i \succeq \lambda \mathbf{I} \right. \right\}.$$

We can cover \mathcal{V} by an ϵ -net denoted by \mathcal{N}_ϵ , and its covering entropy $\log \mathcal{N}_\epsilon$ satisfies $\log \mathcal{N}_\epsilon = \tilde{O}(d^2)$. Such an approach, although fixing the dependency issue in (6.1), introduces an additional \sqrt{d} factor to their final regret since each V belongs to a quadratic function class by their optimistic construction, which prevents them from achieving the optimal d dependency in the regret.

In comparison, our approach gets around the covering number issue by decomposing the value function $V_{k,h}$ into the optimal value function V_{h+1}^* and the sub-optimality gap $V_{k,h+1} - V_{h+1}^*$. Such an analysis approach has been firstly considered in the tabular MDPs Azar et al. (2017); Zhang et al. (2021a) and later in the linear MDP by Hu et al. (2022). More specifically, we have

$$\begin{aligned} & [\widehat{\mathbb{P}}_{k,h} V_{k,h+1}](s, a) - [\mathbb{P}_h V_{k,h+1}](s, a) \\ & \approx \underbrace{\phi(s, a)^\top \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \sigma_{i,h}^{-2} \phi(s_h^i, a_h^i) \eta_{i,h}(V_{k,h+1})}_{I_1} \\ & = \underbrace{\phi(s, a)^\top \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \sigma_{i,h}^{-2} \phi(s_h^i, a_h^i) \eta_{i,h}(V_{h+1}^*)}_{I_1} \\ & \quad + \underbrace{\phi(s, a)^\top \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \sigma_{i,h}^{-2} \phi(s_h^i, a_h^i) \eta_{i,h}(\Delta V_{k,h+1})}_{I_2}, \quad (6.2) \end{aligned}$$

where $\Delta V_{k,h+1} = V_{k,h+1} - V_{h+1}^*$ denotes the estimation error for value function $V_{k,h+1}$. For the first term I_1 , as discussed in Section 4.2, we can use $V_{k,h+1}$ to approximate the optimal value function V_{h+1}^* and the estimation error for the variance can be bounded by:

$$|[\widehat{\mathbb{V}}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k)| \leq E_{k,h} + D_{k,h}.$$

Since the optimal value function V_{h+1}^* is fixed across all episodes $k \in [K]$ and does not depend on the past observations, such an approach can prevent the covering number argument and save a \sqrt{d} factor in the regret compared with Jin et al. (2020). For the second term I_2 , the sub-optimality gap $\Delta V_{k,h+1} = V_{k,h+1} - V_{h+1}^*$ depends on the past observations and we still need to use the covering number argument. However, the magnitude of the sub-optimality gap $\Delta V_{k,h+1}$ is small provided that $V_{k,h+1}$ is an accurate estimate for V_{h+1}^* . In this case, term I_2 will be dominated by term I_1 even with the extra factors from the covering number argument. With the help of the decomposition, we have the following Bernstein-type error bound between the estimated $\widehat{\mathbb{P}}_{k,h} V_{k,h+1}$ and its true value:

$$|\widehat{\mathbf{w}}_{k,h}^\top \phi(s, a) - \mathbb{P}_h V_{k,h+1}(s, a)| \leq \beta \|\phi(s, a)\|_{\Sigma_{k,h}^{-1}}.$$

This analysis also explains why Algorithm 1 needs to estimate the variance of V_{h+1}^* and $\Delta V_{k,h+1}$ instead of $V_{k,h+1}$.

6.2 Monotonic Variance Estimator

Here we provide more details about the variance estimator $\sigma_{k,h}$. According to our previous discussion, we decompose the value function $V_{k,h}$, and only need to control the estimation errors I_1 and I_2 in (6.2) separately. In order to derive a Bernstein-type error bound, we use Lemma 4.1 for both the optimal value function V_{h+1}^* and $\Delta V_{k,h+1}$, which require an estimation for the variance $\mathbb{V}_h V_{h+1}^*$ and $\mathbb{V}_h[V_{k,h+1} - V_{h+1}^*]$. For the variance $\mathbb{V}_h V_{h+1}^*$, as we discussed in Section 4.2, we approximate it with the following empirical variance:

$$\begin{aligned} & \bar{\mathbb{V}}_h V_{k,h+1}(s_h^k, a_h^k) - [\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k) \\ &= \underbrace{\bar{\mathbb{V}}_h V_{k,h+1}(s_h^k, a_h^k) - [\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k)}_{J_1} \\ & \quad + \underbrace{[\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k)}_{J_2}, \end{aligned}$$

where the estimation error J_1 can be controlled by a Hoeffding-type bound (term $E_{k,h}$) and the estimation error J_2 can be upper bounded by

$$\begin{aligned} & |[\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k)| \\ & \leq 4H[\mathbb{P}_h(V_{k,h+1} - V_{h+1}^*)](s_h^k, a_h^k). \end{aligned}$$

For this error bound, the optimal value function V_{h+1}^* is not observable and we replace it by the *pessimistic* value function $\check{V}_{k,h}$, which gives an upper bound

$$\begin{aligned} & 4H[\mathbb{P}_h(V_{k,h+1} - V_{h+1}^*)](s_h^k, a_h^k) \\ & \leq 4H[\mathbb{P}_h(V_{k,h+1} - \check{V}_{k,h+1})](s_h^k, a_h^k). \end{aligned}$$

The above term is further dominated by the term $D_{k,h}$. With a similar approach, for the variance $\mathbb{V}_h[V_{k,h+1} - V_{h+1}^*]$, it can be upper bound by

$$\begin{aligned} & [\mathbb{V}_h(V_{k,h+1} - V_{h+1}^*)](s_h^k, a_h^k) \\ & \leq 2H[\mathbb{P}_h(V_{k,h+1} - V_{h+1}^*)](s_h^k, a_h^k) \\ & \leq 2H[\mathbb{P}_h(V_{k,h+1} - \check{V}_{k,h+1})](s_h^k, a_h^k) \approx D_{k,h}/(d^3 H), \end{aligned} \tag{6.3}$$

where we approximate the optimal value function V_{h+1}^* by the *pessimistic* value function $\check{V}_{k,h}$ and introduce an extra $d^3 H$ -factor in the $D_{k,h}$ to offset the error caused by the covering number argument.

However, there exists another difficulty in our algorithm and analysis when we extend the result in (6.3) to future episode $i > k$. In detail, while the value function $V_{i,h+1}(s_{h+1}^k)$ and corresponding variance $[\mathbb{V}_h(V_{i,h+1} - V_{h+1}^*)](s_h^k, a_h^k)$ will change across different episodes, the estimated variance $\sigma_{k,h}$ is chosen at episode k and cannot be changed in the subsequent episode. Therefore, $\sigma_{k,h}$

should be a uniform variance upper bound for all subsequent episodes. To achieve such a uniform variance upper bound, it suffices to have the sub-optimality gap $V_{k,h+1} - V_{h+1}^*$ to be monotonically decreasing. Our solution is to set $V_{k,h+1}(s)$ to be a monotonically decreasing sequence in k given any state s , by setting it as the minimum between its current estimate and its predecessor $V_{k-1,h+1}$ (Line 9 of Algorithm 1). A similar approach is also applied to $\check{V}_{k,h+1}$ to guarantee the estimate sequence is monotonically increasing in k . Then, the following property shows that the estimated variance of the sub-optimality at episode k $D_{k,h}$ holds for all the subsequent episodes.

$$[\mathbb{V}_h(V_{i,h+1} - V_{h+1}^*)](s_h^k, a_h^k) \leq D_{k,h}/(d^3 H), \forall i > k.$$

This idea was firstly introduced by Azar et al. (2017) for tabular MDPs. Hu et al. (2022) adopted a similar idea to guarantee the monotonicity for linear MDPs, while their approach is to construct another sequence of ‘‘over-optimistic’’ value functions, which turns out to be flawed as we will discuss in Appendix A.

6.3 Rare-Switching Value Function Update

As we discussed in Section 6.2, we ensure the monotonicity and construct the variance estimation, by taking minimization with its predecessor $V_{k-1,h+1}$. However, this approach will introduce an extra issue for the augmented value function class. In detail, the optimistic value function $V_{k,h}$ can be denoted by the minimum over several quadratic functions and belongs to the following function class,

$$\begin{aligned} \mathcal{V}_h = \left\{ V \mid V(\cdot) = \max_a \min_{1 \leq i \leq l} \min \left(H, r_h(\cdot, a) + \mathbf{w}_i^\top \phi(\cdot, a) \right. \right. \\ \left. \left. + \beta \sqrt{\phi(\cdot, a)^\top \Sigma_i^{-1} \phi(\cdot, a)} \right), \|\mathbf{w}_i\| \leq L, \Sigma_i \succeq \lambda \mathbf{I} \right\}, \end{aligned}$$

where l is the number of quadratic functions and equals to the number of policy updates in Algorithm 1. Here, we denote the covering number of that function class by \mathcal{N} , and the covering number of a quadratic function class by \mathcal{N}_q . We are specifically interested in the covering entropy $\log \mathcal{N}$, which is a standard complexity measure of the function class, and it will directly affect the regret of our algorithm. The standard approach to computing the covering entropy suggests that $\log \mathcal{N} = l \log \mathcal{N}_q$. In this case, if we update the value function at each episode $k \in [K]$ and minimize with its predecessor, then there will be an extra K factor in the covering number, which is unacceptable. Inspired by the ‘‘rare-switching’’ technique (Abbasi-Yadkori et al., 2011; Wang et al., 2021), it is not necessary or efficient to update the value functions $V_{k,h+1}$ and $\check{V}_{k,h+1}$ at each episode. Instead, we only need to update the value function when the determinant of the covariance matrix grows much larger than before (Line 8 in Algorithm 1), which requires at most $\tilde{O}(dH)$

updates. Such an update strategy reduces the covering entropy from $\log \mathcal{N} = K \log \mathcal{N}_q$ to $\log \mathcal{N} = \tilde{O}(dH) \log \mathcal{N}_q$, which makes the regret of Algorithm 1 tight. In addition, the “rare-switching” nature also reduces the computational complexity from $\Omega(K^2)$ to $\Omega(K)$, which makes LSVI-UCB++ more efficient.

7 Conclusions and Future Work

In this paper, we propose a near-optimal algorithm LSVI-UCB++ for linear MDPs. LSVI-UCB++ is based on weighted ridge regression, where the weights are constructed from a novel variance estimator that comes from a direct estimation of the variance of the true value function, and a “rare-switching” updating rule to update the value function estimator. We prove that with high probability, LSVI-UCB++ obtains an $\tilde{O}(d\sqrt{H^3K})$ regret, which matches the lower bound in Zhou et al. (2021a) up to logarithmic factors. Our algorithm is also computationally efficient. In the future, we will study how to design computationally efficient near-optimal RL algorithms with general nonlinear function approximation with misspecification.

Acknowledgements

We thank the anonymous reviewers for their helpful comments. JH, HZ, DZ and QG are supported in part by the National Science Foundation CAREER Award 1906169 and the Sloan Research Fellowship. The views and conclusions contained in this paper are those of the authors and should not be interpreted as representing any funding agencies.

References

ABBASI-YADKORI, Y., PÁL, D. and SZEPESVÁRI, C. (2011). Improved algorithms for linear stochastic bandits. In *NIPS*, vol. 11.

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

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.

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**.

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*.

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

HE, J., ZHOU, D. and GU, Q. (2021b). Nearly minimax optimal reinforcement learning for discounted mdps. *Advances in Neural Information Processing Systems* **34** 22288–22300.

HE, J., ZHOU, D., ZHANG, T. and GU, Q. (2022). Nearly optimal algorithms for linear contextual bandits with adversarial corruptions. *arXiv preprint arXiv:2205.06811*.

HENDERSON, C. R. (1975). Best linear unbiased estimation and prediction under a selection model. *Biometrics* 423–447.

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

JIA, Z., YANG, L., SZEPESVARI, C. and WANG, M. (2020). Model-based reinforcement learning with value-targeted regression. In *Learning for Dynamics and Control*. 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., ALLEN-ZHU, Z., BUBECK, S. and JORDAN, M. I. (2018). Is q-learning provably efficient? *Advances in neural information processing systems* **31**.

JIN, C., YANG, Z., WANG, Z. and JORDAN, M. I. (2019). Provably efficient reinforcement learning with linear function approximation. *arXiv preprint arXiv:1907.05388*.

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.

KIRSCHNER, J. and KRAUSE, A. (2018). Information directed sampling and bandits with heteroscedastic noise. In *Conference On Learning Theory*. PMLR.

- LATTIMORE, T., CRAMMER, K. and SZEPESVÁRI, C. (2015). Linear multi-resource allocation with semi-bandit feedback. *Advances in Neural Information Processing Systems* **28**.
- 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.
- SIMCHOWITZ, M. and JAMIESON, K. G. (2019). Non-asymptotic gap-dependent regret bounds for tabular mdps. *Advances in Neural Information Processing Systems* **32**.
- 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.
- WANG, T., ZHOU, D. and GU, Q. (2021). Provably efficient reinforcement learning with linear function approximation under adaptivity constraints. *Advances in Neural Information Processing Systems* **34** 13524–13536.
- WANG, Y., WANG, R., DU, S. S. and KRISHNAMURTHY, A. (2020). Optimism in reinforcement learning with generalized linear function approximation. In *International Conference on Learning Representations*.
- 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.
- 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. and BRUNSKILL, E. (2019). Tighter problem-dependent regret bounds in reinforcement learning without domain knowledge using value function bounds. In *International Conference on Machine Learning*. 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.
- ZHANG, Z. and JI, X. (2019). Regret minimization for reinforcement learning by evaluating the optimal bias function. *Advances in Neural Information Processing Systems* **32**.
- ZHANG, Z., JI, X. and DU, S. (2021a). Is reinforcement learning more difficult than bandits? a near-optimal algorithm escaping the curse of horizon. In *Conference on Learning Theory*. PMLR.
- ZHANG, Z., YANG, J., JI, X. and DU, S. S. (2021b). Improved variance-aware confidence sets for linear bandits and linear mixture mdp. *Advances in Neural Information Processing Systems* **34** 4342–4355.
- ZHANG, Z., ZHOU, Y. and JI, X. (2020). Almost optimal model-free reinforcement learning via reference-advantage decomposition. *Advances in Neural Information Processing Systems* **33** 15198–15207.
- ZHOU, D. and GU, Q. (2022). Computationally efficient horizon-free reinforcement learning for linear mixture mdps. *arXiv preprint arXiv:2205.11507*.
- ZHOU, D., GU, Q. and SZEPESVARI, C. (2021a). Nearly minimax optimal reinforcement learning for linear mixture markov decision processes. In *Conference on Learning Theory*. PMLR.
- ZHOU, D., HE, J. and GU, Q. (2021b). Provably efficient reinforcement learning for discounted mdps with feature mapping. In *International Conference on Machine Learning*. PMLR.

A Comparison with Hu et al. (2022)

In this section, we give a detailed comparison with Hu et al. (2022). We first elaborate on the importance of the monotonic property, then discuss the issue on the over-optimistic value function $\hat{V}_{k,h}(s)$ proposed in Hu et al. (2022) and finally illustrate the difference in the algorithm design between the algorithm in Hu et al. (2022) and our algorithm.

As we discuss in the Section 6, both our LSVI-UCB++ algorithm and LSVI-UCB+ algorithm in Hu et al. (2022) get rid of the covering number issue by decomposing the value function $V_{k,h+1}(s)$ to $V_{h+1}^*(s)$ and $V_{k,h+1}(s) - V_{h+1}^*(s)$. In detail, from the proof of Lemma B.5, we have shown in (D.18) that the estimation error can be decomposed as

$$\begin{aligned} & \left\| \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (V_{k,h+1}(s_{h+1}^i) - \mathbb{P}_h V_{k,h+1}(s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}} \\ &= \left\| \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (V_{k,h+1}(s_{h+1}^i) - \mathbb{P}_h V_{k,h+1}(s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}^{-1}} \\ &\leq \underbrace{\left\| \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (V_{h+1}^*(s_{h+1}^i) - \mathbb{P}_h V_{h+1}^*(s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}^{-1}}}_{J_1} \\ &\quad + \underbrace{\left\| \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (\Delta V_{k,h+1}(s_{h+1}^i) - [\mathbb{P}_h(\Delta V_{k,h+1})](s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}^{-1}}}_{J_2}. \end{aligned}$$

To control the concentration error on the term J_2 , we use Lemma 4.1 and only need to estimate the variance $\mathbb{V}_h[\Delta V_{k,h+1}](s, a) = \mathbb{V}_h[(V_{k,h+1} - V_{h+1}^*)](s, a)$, which is trivially upper bounded by $2H \cdot [\mathbb{P}_h \Delta V_{k,h+1}](s, a)$. In order to make the upper bound of variance at episode i hold for all subsequent episode $k > i$, we need to guarantee that the trivial upper bound $2H \cdot [\mathbb{P}_h \Delta V_{k,h+1}](s, a)$ is decreasing in k , which requires the optimistic value function $V_{k,h+1}(s)$ to be monotonically decreasing.

To satisfy this requirement, Hu et al. (2022) constructs an over-optimistic value function $\hat{V}_{k,h}(s)$, which has the following monotonicity property.

Lemma A.1 (Lemma D.2, Hu et al. 2022). *For any stage $h \in [H]$ and episodes $i < j$, the over-optimistic value function $\hat{V}_{j,h}(s)$ satisfies:*

$$\hat{V}_{i,h}(s) \geq \hat{V}_{j,h}(s),$$

where $\hat{V}_{j,h}(s)$ is the optimistic value function.

Based on this monotonically decreasing property, the estimation error $2H \cdot [\mathbb{P}_h(\hat{V}_{i,h+1} - V_{h+1}^*)](s, a)$ is a uniform variance upper bound for all subsequent episodes. Unfortunately, the last inequality in the proof of Lemma A.1 claims that $[\mathbb{P}_h V_{i,h+1}](s, a) - [\mathbb{P}_h V_{j,h+1}](s, a)$ holds due to $\hat{V}_{i,h}(s) \geq \hat{V}_{j,h}(s)$, which is not true. Thus, the monotonic property of over-optimistic value function $\hat{V}_{k,h}(s)$ does not hold and the estimated variance $\sigma_{i,h}$ may no longer be a variance upper bound for the subsequent episodes.

In comparison, our LSVI-UCB++ ensures the monotonic property by choosing the minimum of the optimistic value functions in the first k episodes (See Line 9). As the cost of ensuring monotonic property, the resulting value function class can be regarded as a minimum over K quadratic function classes, and the covering number grows exponentially in K . To overcome this problem, we utilize the ‘rare-switching’ technique from previous works (Abbasi-Yadkori et al., 2011; Wang et al., 2021), which reduces the number of updates to $\tilde{O}(dH)$ and thus controls the complexity growth of the resulting value function class.

B Proof Sketch

This section is devoted to provide a proof sketch of Theorem 5.1.

B.1 High-Probability Events

We first introduce the following high-probability events:

1. We define \mathcal{E} as the event that the following inequalities hold for all $s, a, k, h \in \mathcal{S} \times \mathcal{A} \times [K] \times [H]$.

$$\begin{aligned} |\widehat{\mathbf{w}}_{k,h}^\top \boldsymbol{\phi}(s, a) - [\mathbb{P}_h V_{k,h+1}](s, a)| &\leq \bar{\beta} \sqrt{\boldsymbol{\phi}(s, a)^\top \boldsymbol{\Sigma}_{k,h}^{-1} \boldsymbol{\phi}(s, a)}, \\ |\widetilde{\mathbf{w}}_{k,h}^\top \boldsymbol{\phi}(s, a) - [\mathbb{P}_h V_{k,h+1}^2](s, a)| &\leq \widetilde{\beta} \sqrt{\boldsymbol{\phi}(s, a)^\top \boldsymbol{\Sigma}_{k,h}^{-1} \boldsymbol{\phi}(s, a)}, \\ |\check{\mathbf{w}}_{k,h}^\top \boldsymbol{\phi}(s, a) - [\mathbb{P}_h \check{V}_{k,h+1}](s, a)| &\leq \bar{\beta} \sqrt{\boldsymbol{\phi}(s, a)^\top \boldsymbol{\Sigma}_{k,h}^{-1} \boldsymbol{\phi}(s, a)}, \end{aligned}$$

where $\widetilde{\beta} = O\left(H^2 \sqrt{d\lambda} + \sqrt{d^3 H^4 \log^2(dHK/(\delta\lambda))}\right)$ and $\bar{\beta} = O\left(H \sqrt{d\lambda} + \sqrt{d^3 H^2 \log^2(dHK/(\delta\lambda))}\right)$.

2. We define $\widetilde{\mathcal{E}}_h$ as the event such that for all episode $k \in [K]$, stage $h \leq h' \leq H$ and state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$, the weight vector $\widehat{\mathbf{w}}_{k,h}$ satisfies that

$$|\widehat{\mathbf{w}}_{k,h'}^\top \boldsymbol{\phi}(s, a) - [\mathbb{P}_h V_{k,h'+1}](s, a)| \leq \beta \sqrt{\boldsymbol{\phi}(s, a)^\top \boldsymbol{\Sigma}_{k,h'}^{-1} \boldsymbol{\phi}(s, a)}, \quad (\text{B.1})$$

where $\beta = O\left(H \sqrt{d\lambda} + \sqrt{d \log^2(1 + dKH/(\delta\lambda))}\right)$. For simplicity, we further define events $\widetilde{\mathcal{E}} = \widetilde{\mathcal{E}}_1$ that (B.1) holds for all stage $h \in [H]$.

Our ultimate goal is to show that $\widetilde{\mathcal{E}}$ holds with high probability. Intuitively speaking, \mathcal{E} serves as a ‘coarse’ event where the concentration results hold with a larger confidence radius $\bar{\beta}$ and $\widetilde{\beta}$, and $\widetilde{\mathcal{E}}$ serves as a ‘refined’ event where the confidence radius β is tighter than $\bar{\beta}$ and $\widetilde{\beta}$. To start with, the following lemma shows that \mathcal{E} holds with high probability.

Lemma B.1. *Event \mathcal{E} holds with probability at least $1 - 7\delta$.*

Next, we prove $\widetilde{\mathcal{E}} = \widetilde{\mathcal{E}}_1$ holds with high probability. Since $\widehat{\mathbf{w}}_{k,h}$ ’s are obtained from weighted linear regression whose weights depend on the variances of $V_{k,h+1}$, the key technical challenge is to show that our adapted weights $\sigma_{k,h}$ ’s are indeed upper bounds of these variances for all $h \in [H]$. We use backward induction to prove such a statement. In detail, the following two lemmas provide estimation error bounds at stage h conditioned on $\widetilde{\mathcal{E}}_{h+1}$.

Lemma B.2. *On the event \mathcal{E} and $\widetilde{\mathcal{E}}_{h+1}$, for each episode $k \in [K]$ and stage h , the estimated variance satisfies*

$$\begin{aligned} |[\bar{\mathbb{V}}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k)| &\leq E_{k,h}, \\ |[\bar{\mathbb{V}}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k)| &\leq E_{k,h} + D_{k,h}. \end{aligned}$$

Lemma B.3. *On the event \mathcal{E} and $\widetilde{\mathcal{E}}_{h+1}$, for any episode k and $i > k$, we have*

$$|\mathbb{V}_h(V_{i,h+1} - V_{h+1}^*)|(s_h^k, a_h^k) \leq D_{k,h}/(d^3 H).$$

We also have the following lemma, which shows that our constructed value functions Q, V , and \check{Q}, \check{V} are optimistic and pessimistic estimators of the true value functions under the events we defined before.

Lemma B.4. *On the event \mathcal{E} and $\widetilde{\mathcal{E}}_h$, for all episode $k \in [K]$ and stage $h \leq h' \leq H$, we have $Q_{k,h}(s, a) \geq Q_h^*(s, a) \geq \check{Q}_{k,h}(s, a)$. In addition, we have $V_{k,h}(s) \geq V_h^*(s) \geq \check{V}_{k,h}(s)$.*

Equipped with Lemmas B.2, B.3 and B.4, one can easily prove $\sigma_{k,h}$ indeed serves as an upper bound of the true variance of $V_{k,h+1}$ at stage h . Therefore, by the backward induction, we can prove the following lemma.

Lemma B.5. *On the events \mathcal{E} , event $\widetilde{\mathcal{E}}$ holds with probability at least $1 - \delta$.*

B.2 Regret Decomposition

Now, we prove the regret bound based on the high-probability events defined before. Based on Lemma B.4, for all stage $h \in [H]$ and episode $k \in [K]$, we have $Q_{k,h}(s_h^k, a_h^k) = V_{k,h}(s_h^k) \geq V_h^*(s_h^k)$. Thus, we can bound the regret as follows,

$$\mathbf{Regret}(K) = \sum_{k=1}^K (V_1^*(s_1^k) - V_{k,1}^\pi(s_1^k))$$

$$\begin{aligned}
 &\leq \sum_{k=1}^K (V_{k,1}(s_1^k) - V_{k,1}^{\pi^k}(s_1^k)) \\
 &\leq \sum_{k=1}^K \sum_{h=1}^H \left\{ [\mathbb{P}_h(V_{k,h+1} - V_{k,h+1}^{\pi^k})](s_h^k, a_h^k) - (V_{k,h+1}(s_{h+1}^k) - V_{k,h+1}^{\pi^k}(s_{h+1}^k)) \right\} \\
 &\quad + O \left(\sum_{k=1}^K \sum_{h=1}^H \min \left(\beta \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)}, H \right) \right),
 \end{aligned}$$

where the last inequality holds due to the decomposition of the difference of value functions and the high-probability events defined in Section B.1. Using standard regret analysis, we can bound the first term as the sum of a martingale difference sequence. Then it remains to bound the summation of bonus terms, $\sum_{k=1}^K \sum_{h=1}^H \beta \|\phi(s_h^k, a_h^k)\|_{\Sigma_{k,h}^{-1}}$. By Cauchy-Schwartz inequality, this summation can be bounded by

$$\sum_{k=1}^K \sum_{h=1}^H \beta \|\phi(s_h^k, a_h^k)\|_{\Sigma_{k,h}^{-1}} \leq \tilde{O} \left(d^4 H^8 + \beta d^7 H^5 + \beta \sqrt{dHK + dH \sum_{k=1}^K \sum_{h=1}^H \sigma_{k,h}^2} \right),$$

where the calculation details are deferred to Lemma E.1. According to the definition of $\sigma_{k,h}^2$, we have $\sigma_{k,h}^2 \leq O([\tilde{V}_{k,h} V_{k,h+1}](s_h^k, a_h^k) + E_{k,h} + D_{k,h} + H)$. By carefully bounding the summation of $[\tilde{V}_{k,h} V_{k,h+1}](s_h^k, a_h^k)$ by relating them to the summation of $[\mathbb{V}_h V_{k,h+1}^{\pi^k}](s_h^k, a_h^k)$ and using the total variance lemma (Lemma C.5, Jin et al. 2018), we have

$$\sum_{k=1}^K \sum_{h=1}^H \sigma_{k,h}^2 \leq \tilde{O}(H^2 K + d^{10.5} H^{16}).$$

Putting all pieces together, we can obtain the high-probability regret bound

$$\text{Regret}(K) \leq \tilde{O}(d\sqrt{H^3 K} + d^7 H^8).$$

C Detailed Proof of Theorem 5.1

In this section, we provide the proof of Theorem 5.1. Firstly, for the stochastic transition noises, we define the following high-probability events:

$$\begin{aligned}
 \mathcal{E}_1 &= \left\{ \forall h \in [H], \sum_{k=1}^K \sum_{h'=h}^H [\mathbb{P}_h(V_{k,h+1} - V_{k,h+1}^{\pi^k})](s_h^k, a_h^k) \right. \\
 &\quad \left. - \sum_{k=1}^K \sum_{h'=h}^H (V_{k,h+1}(s_{h+1}^k) - V_{k,h+1}^{\pi^k}(s_{h+1}^k)) \leq 2\sqrt{2H^3 K \log(H/\delta)} \right\}, \\
 \mathcal{E}_2 &= \left\{ \forall h \in [H], \sum_{k=1}^K \sum_{h'=h}^H [\mathbb{P}_h(V_{k,h+1} - \tilde{V}_{k,h+1})](s_h^k, a_h^k) \right. \\
 &\quad \left. - \sum_{k=1}^K \sum_{h'=h}^H (V_{k,h+1}(s_{h+1}^k) - \tilde{V}_{k,h+1}(s_{h+1}^k)) \leq 2\sqrt{2H^3 K \log(H/\delta)} \right\}.
 \end{aligned}$$

Then according to the Azuma–Hoeffding inequality (Lemma G.2), we have $\Pr(\mathcal{E}_1) \geq 1 - \delta$ and $\Pr(\mathcal{E}_2) \geq 1 - \delta$. Based on the definition of events $\mathcal{E}_1, \mathcal{E}_2$ and events $\mathcal{E}, \tilde{\mathcal{E}}$ in Section B.1, the regret in the first K episodes can be upper bounded by the summation of estimated variance $\sum_{k=1}^K \sum_{h=1}^H \sigma_{k,h}^2$ and we have the following lemma.

Lemma C.1. *On the events $\tilde{\mathcal{E}}, \mathcal{E}$ and \mathcal{E}_1 , for all stage $h \in [H]$, the regret in the first K episodes is upper bounded by:*

$$\sum_{k=1}^K (V_{k,h}(s_h^k) - V_{k,h}^{\pi^k}(s_h^k)) \leq 16d^4 H^8 \iota + 40\beta d^7 H^5 \iota + 8\beta \sqrt{2dH\iota \sum_{h=1}^H \sum_{k=1}^K (\sigma_{k,h}^2 + H)} + 4\sqrt{H^3 K \log(H/\delta)},$$

and for all stage $h \in [H]$, we further have

$$\begin{aligned} & \sum_{k=1}^K \sum_{h=1}^H [\mathbb{P}_h(V_{k,h+1} - V_{k,h+1}^{\pi^k})](s_h^k, a_h^k) \\ & \leq 16d^4 H^9 \iota + 40\beta d^7 H^6 \iota + 8H\beta \sqrt{2dH\iota \sum_{h=1}^H \sum_{k=1}^K (\sigma_{k,h}^2 + H) + 4\sqrt{H^5 K \log(H/\delta)}}, \end{aligned}$$

where $\iota = \log(1 + K/(d\lambda))$.

In addition, for the sub-optimality gap between the optimistic value function $V_{k,h}(s)$ and pessimistic value function $\check{V}_{k,h}(s)$, we have the following lemma.

Lemma C.2. *On the events $\tilde{\mathcal{E}}$, \mathcal{E} and \mathcal{E}_2 , the difference between the optimistic value function $V_{k,h}$ and the pessimistic value function $\check{V}_{k,h}$ is upper bounded by:*

$$\begin{aligned} & \sum_{k=1}^K \sum_{h=1}^H [\mathbb{P}_h(V_{k,h+1} - \check{V}_{k,h+1})](s_h^k, a_h^k) \\ & \leq 32d^4 H^9 \iota + 40(\beta + \bar{\beta})d^7 H^6 \iota + 8H(\beta + \bar{\beta}) \sqrt{2dH\iota \sum_{h=1}^H \sum_{k=1}^K (\sigma_{k,h}^2 + H) + 4\sqrt{H^5 K \log(H/\delta)}}, \end{aligned}$$

where $\iota = \log(1 + K/(d\lambda))$.

For the summation of variance $\sum_{k=1}^K \sum_{h=1}^H [\mathbb{V}_h V_{k,h+1}^{\pi^k}](s_h^k, a_h^k)$, we denote the following high probability events \mathcal{E}_3 :

$$\mathcal{E}_3 = \left\{ \sum_{k=1}^K \sum_{h=1}^H [\mathbb{V}_h V_{k,h+1}^{\pi^k}](s_h^k, a_h^k) \leq 3H^2 K + 3H^3 \log(1/\delta) \right\}.$$

Then Lemma C.5 in Jin et al. (2018) shows that the probability of events \mathcal{E}_3 is lower bounded by $\Pr(\mathcal{E}_3) \geq 1 - \delta$. Furthermore, on the event $\mathcal{E} \cap \tilde{\mathcal{E}} \cap \mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3$, the following lemma gives an upper bound of the total estimated variance $\sum_{h=1}^H \sum_{k=1}^K \sigma_{k,h}^2$.

Lemma C.3. *On the event $\mathcal{E} \cap \tilde{\mathcal{E}} \cap \mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3$, the total estimated variance is upper bounded by:*

$$\sum_{k=1}^K \sum_{h=1}^H \sigma_{k,h}^2 \leq O(H^2 K + d^{10.5} H^{16} \log^{1.5}(1 + dKH/\delta)).$$

With all previous lemma, we start to prove our main Theorem 5.1.

Proof of Theorem 5.1. On the event $\mathcal{E} \cap \tilde{\mathcal{E}} \cap \mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3$, the regret is upper bounded by:

$$\begin{aligned} \mathbf{Regret}(K) &= \sum_{k=1}^K (V_1^*(s_1^k) - V_{k,1}^{\pi^k}(s_1^k)) \\ &\leq \sum_{k=1}^K (V_{k,1}(s_1^k) - V_{k,1}^{\pi^k}(s_1^k)) \\ &\leq 16d^4 H^8 \iota + 40\beta d^7 H^5 \iota + 8\beta \sqrt{2dH\iota \sum_{h=1}^H \sum_{k=1}^K (\sigma_{k,h}^2 + H) + 4\sqrt{H^3 K \log(H/\delta)}} \\ &= \tilde{O}\left(d^7 H^8 + d\sqrt{H^3 K \log^2(1 + dKH/\delta)}\right), \end{aligned} \tag{C.1}$$

where $\iota = \log(1 + K/(d\lambda))$, the first inequality holds due to Lemma B.4, the second inequality holds due to Lemma C.2 and the last inequality holds due to Lemma C.3. Since the event $\mathcal{E} \cap \tilde{\mathcal{E}} \cap \mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3$ holds with probability at least $1 - 7\delta$, (C.1) holds. In addition, according to Lemma F.1, the number of updates for $Q_{k,h}, \tilde{Q}_{k,h}$ is upper bounded by $O(dH \log(1 + K/\lambda))$. Thus, we complete the proof of Theorem 5.1. updates for $Q_{k,h}, \tilde{Q}_{k,h}$ is upper bounded by $O(dH \log(1 + K/\lambda))$. \square

D Proof of Lemmas in Section B

In this section, we provide the proof of Lemmas in Section B and we need the following lemma, which extends Lemma D.4 in Jin et al. (2020) to weighted ridge regression.

Lemma D.1. (Lemma D.4, Jin et al. 2020 with weighted linear regression) *Let $\{x_k\}_{k=1}^\infty$ be a real-valued stochastic process on state space \mathcal{S} with corresponding filtration $\{\mathcal{F}_k\}_{k=1}^\infty$. Let $\{\phi_k\}_{k=1}^\infty$ be an \mathbb{R}^d -valued stochastic process, where $\phi_k \in \mathcal{F}_{k-1}$ and $\|\phi_k\|_2 \leq 1$. Let $\{w_k\}_{k=1}^\infty$ be an real-valued stochastic process where $w_k \in \mathcal{F}_{k-1}$ and $0 \leq w_k \leq C$. For any $k \geq 0$, we define $\Sigma_k = \lambda \mathbf{I} + \sum_{i=1}^k w_i^2 \phi_i \phi_i^\top$. Then with probability at least $1 - \delta$, for all $k \in \mathbb{N}$ and all function $V \in \mathcal{V}$ with $\max_s |V(x)| \leq H$, we have*

$$\left\| \sum_{i=1}^k w_i^2 \phi_i \left\{ V(x_i) - \mathbb{E}[V(x_i) | \mathcal{F}_{i-1}] \right\} \right\|_{\Sigma_k^{-1}}^2 \leq 4C^2 H^2 \left[\frac{d}{2} \log(1 + kC^2/\lambda) + \log \frac{\mathcal{N}_\epsilon}{\delta} \right] + 8k^2 C^4 \epsilon^2 / \lambda,$$

where \mathcal{N}_ϵ is the ϵ -covering number of the function class \mathcal{V} with respect to the distance function $\text{dist}(V_1, V_2) = \max_s |V_1(s) - V_2(s)|$.

Proof of Lemma D.1. For any function $V \in \mathcal{V}$, based on the definition of ϵ -covering number, there exists a function \tilde{V} in the ϵ -net, such that

$$\text{dist}(V, \tilde{V}) \leq \epsilon. \quad (\text{D.1})$$

Therefore, the concentration error for the value function can be decomposed as

$$\begin{aligned} & \left\| \sum_{i=1}^k w_i^2 \phi_i \left\{ V(x_i) - \mathbb{E}[V(x_i) | \mathcal{F}_{i-1}] \right\} \right\|_{\Sigma_k^{-1}}^2 \\ & \leq 2 \underbrace{\left\| \sum_{i=1}^k w_i^2 \phi_i \left\{ \tilde{V}(x_i) - \mathbb{E}[\tilde{V}(x_i) | \mathcal{F}_{i-1}] \right\} \right\|_{\Sigma_k^{-1}}^2}_{I_1} + 2 \underbrace{\left\| \sum_{i=1}^k w_i^2 \phi_i \left\{ \Delta_V(x_i) - \mathbb{E}[\Delta_V(x_i) | \mathcal{F}_{i-1}] \right\} \right\|_{\Sigma_k^{-1}}^2}_{I_2}, \end{aligned} \quad (\text{D.2})$$

where $\Delta_V = V - \tilde{V}$ and the inequality holds due to $\|\mathbf{a} + \mathbf{b}\|_{\Sigma}^2 \leq 2\|\mathbf{a}\|_{\Sigma}^2 + 2\|\mathbf{b}\|_{\Sigma}^2$. For any fixed value function \tilde{V} , we apply Lemma G.5 with $\mathbf{x}_i = w_i \phi_i$, $\eta_t = w_i \tilde{V}(x_i) - w_i \mathbb{E}[\tilde{V}(x_i)]$. According to the definition of \mathbf{x}_i, η_i , we have following property

$$\begin{aligned} \|\mathbf{x}_i\|_2 &= w_i \|\phi_i\|_2 \leq C, \\ \mathbb{E}[\eta_i | \mathcal{F}_i] &= 0, |\eta_i| = |w_i \tilde{V}(x_i) - w_i \mathbb{E}[\tilde{V}(x_i)]| \leq HC. \end{aligned}$$

Therefore, according to Lemma G.5, after taking an union bound over the ϵ -net of the function class \mathcal{V} , with probability at least $1 - \delta/H$, the first term I_1 is upper bounded by:

$$\begin{aligned} I_1 &= \left\| 2 \sum_{i=1}^k w_i^2 \phi_i \left\{ \tilde{V}(x_i) - \mathbb{E}[\tilde{V}(x_i) | \mathcal{F}_{i-1}] \right\} \right\|_{\Sigma_k^{-1}}^2 \\ &\leq 4C^2 H^2 \left[\frac{d}{2} \log(1 + kC^2/\lambda) + \log \frac{\mathcal{N}_\epsilon}{\delta} \right]. \end{aligned} \quad (\text{D.3})$$

For the second term, it can be upper bounded by

$$I_2 = 2 \left\| \sum_{i=1}^k w_i^2 \phi_i \left\{ \Delta_V(x_i) - \mathbb{E}[\Delta_V(x_i) | \mathcal{F}_{i-1}] \right\} \right\|_{\Sigma_k^{-1}}^2$$

$$\begin{aligned}
 &\leq 2k \sum_{i=1}^k \left\| w_i^2 \phi_i \left\{ \Delta_V(x_i) - \mathbb{E}[\Delta_V(x_i) | \mathcal{F}_{i-1}] \right\} \right\|_{\Sigma_k^{-1}}^2 \\
 &\leq 8k^2 C^4 \epsilon^2 / \lambda,
 \end{aligned} \tag{D.4}$$

where the first inequality holds due to the Cauchy-Schwartz inequality and the last inequality holds due to the facts that $|\Delta_V(x_i)| \leq \epsilon$, $0 \leq w_i \leq C$, $\|\phi_i\|_2 \leq 1$, $\Sigma_k \succeq \lambda \mathbf{I}$. Substituting the results in (D.3) and (D.4) into (D.2), we finish the proof of Lemma D.1. \square

D.1 Proof of Lemma B.1

In this subsection, we provide the proof of Lemma B.1, which suggests a Hoeffding-type upper bound for the estimation error.

Proof of Lemma B.1. Firstly, for any fixed stage $h \in [H]$ and the optimistic value function $V_{k,h+1}$, according to Lemma G.1, there exists a vector $\mathbf{w}_{k,h}$ such that $\mathbb{P}_h V_{k,h+1}(s, a)$ can be represented by $\mathbf{w}_{k,h}^\top \phi(s, a)$ and $\|\mathbf{w}_{k,h}\|_2 \leq H\sqrt{d}$. Therefore, the estimation error can be decomposed as

$$\begin{aligned}
 &\|\widehat{\mathbf{w}}_{k,h} - \mathbf{w}_{k,h}\|_{\Sigma_{k,h}} \\
 &= \left\| \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) V_{k,h+1}(s_{h+1}^i) - \Sigma_{k,h}^{-1} \left(\lambda \mathbf{I} + \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top \right) \mathbf{w}_{k,h} \right\|_{\Sigma_{k,h}} \\
 &= \left\| \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (V_{k,h+1}(s_{h+1}^i) - [\mathbb{P}_h V_{k,h+1}](s_h^i, a_h^i)) - \lambda \Sigma_{k,h}^{-1} \mathbf{w}_{k,h} \right\|_{\Sigma_{k,h}} \\
 &\leq \underbrace{\|\lambda \Sigma_{k,h}^{-1} \mathbf{w}_{k,h}\|_{\Sigma_{k,h}}}_{I_1} + \underbrace{\left\| \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (V_{k,h+1}(s_{h+1}^i) - [\mathbb{P}_h V_{k,h+1}](s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}}}_{I_2},
 \end{aligned} \tag{D.5}$$

where the first inequality holds due to the fact that $\|\mathbf{a} + \mathbf{b}\|_{\Sigma} \leq \|\mathbf{a}\|_{\Sigma} + \|\mathbf{b}\|_{\Sigma}$. For the first term I_1 , since $\Sigma_{k,h} \succeq \lambda \mathbf{I}$ and $\|\mathbf{w}_{k,h}\|_2 \leq H\sqrt{d}$, it is upper bounded by

$$I_1 = \|\lambda \mathbf{w}_{k,h}\|_{\Sigma_{k,h}^{-1}} \leq \sqrt{\lambda} \cdot \|\mathbf{w}_{k,h}\|_2 \leq H\sqrt{d\lambda}. \tag{D.6}$$

For the second term I_2 , we apply Lemma D.1 with the optimistic value function class \mathcal{V}_h and $\epsilon = H\sqrt{\lambda}/K$, then for any fixed stage $h \in [H]$, with probability at least $1 - \delta/H$, for all episode $k \in [K]$, we have

$$\begin{aligned}
 I_2 &= \left\| \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (V_{k,h+1}(s_{h+1}^i) - [\mathbb{P}_h V_{k,h+1}](s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}} \\
 &\leq \sqrt{4C^2 H^2 \left[\frac{d}{2} \log(1 + kC^2/\lambda) + \log \frac{HN_\epsilon}{\delta} \right] + 8k^2 C^4 \epsilon^2 / \lambda} \\
 &\leq \sqrt{4H \left[\frac{d}{2} \log(1 + k/(\lambda H)) + \log \frac{HN_\epsilon}{\delta} \right] + 8k^2 \epsilon^2 / (\lambda H^2)} \\
 &\leq \sqrt{4H \left[\frac{d}{2} \log(1 + k/(\lambda H)) + \log \frac{HN_\epsilon}{\delta} \right] + 8} \\
 &= O\left(\sqrt{d^3 H^2 \log^2(dHK/(\delta\lambda))}\right),
 \end{aligned} \tag{D.7}$$

where the first inequality holds due to Lemma D.1, the second inequality holds due to $0 \leq \bar{\sigma}_{i,h}^{-1} \leq 1/\sqrt{H}$, the third inequality holds due to Lemma F.6 and $\epsilon = H\sqrt{\lambda}/K$. Substituting (D.6) and (D.7) into (D.5), we have

$$\|\widehat{\mathbf{w}}_{k,h} - \mathbf{w}_{k,h}\|_{\Sigma_{k,h}} \leq I_1 + I_2 = O\left(H\sqrt{d\lambda} + \sqrt{d^3 H^2 \log^2(dHK/(\delta\lambda))}\right) = \bar{\beta}. \tag{D.8}$$

Therefore, the estimation error is upper bounded by

$$\begin{aligned}
 |\widehat{\mathbf{w}}_{k,h}^\top \phi(s, a) - [\mathbb{P}_h V_{k,h+1}](s, a)| &= |\widehat{\mathbf{w}}_{k,h}^\top \phi(s, a) - \mathbf{w}_{k,h}^\top \phi(s, a)| \\
 &\leq \|\widehat{\mathbf{w}}_{k,h} - \mathbf{w}_{k,h}\|_{\Sigma_{k,h}} \cdot \|\phi(s, a)\|_{\Sigma_{k,h}^{-1}} \\
 &\leq \bar{\beta} \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)},
 \end{aligned} \tag{D.9}$$

where the first inequality holds due to Cauchy-Schwartz inequality and the last inequality holds due to (D.8). Replacing the value function class by the pessimistic value function class \check{V}_h (or squared value function class \check{V}_h^2) and following the same proof of (D.9), we can derive the following upper bound for the estimation errors:

$$\begin{aligned}
 |\check{\mathbf{w}}_{k,h}^\top \phi(s, a) - [\mathbb{P}_h V_{k,h+1}^2](s, a)| &\leq \tilde{\beta} \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)}, \\
 |\check{\mathbf{w}}_{k,h}^\top \phi(s, a) - [\mathbb{P}_h \check{V}_{k,h+1}](s, a)| &\leq \bar{\beta} \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)},
 \end{aligned}$$

where $\tilde{\beta} = O\left(H\sqrt{d\lambda} + \sqrt{d^3 H^4 \log^2(dHK/(\delta\lambda))}\right)$ and $\bar{\beta} = O\left(H\sqrt{d\lambda} + \sqrt{d^3 H^2 \log^2(dHK/(\delta\lambda))}\right)$. Thus, we finish the proof of Lemma B.1. \square

D.2 Proof of Lemma B.2

In this subsection, we provide the proof of Lemma B.2 for the variance estimator.

Proof of Lemma B.2. Firstly, according to Lemma B.1, we have

$$\begin{aligned}
 &|[\bar{\mathbb{V}}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k)| \\
 &= \left| [\check{\mathbf{w}}_{k,h}^\top \phi(s_h^k, a_h^k)]_{[0,H^2]} - [\widehat{\mathbf{w}}_{k,h}^\top \phi(s_h^k, a_h^k)]_{[0,H]}^2 - [\mathbb{P}_h V_{k,h+1}^2](s_h^k, a_h^k) + ([\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k))^2 \right| \\
 &\leq \left| [\check{\mathbf{w}}_{k,h}^\top \phi(s_h^k, a_h^k)]_{[0,H^2]} - [\mathbb{P}_h V_{k,h+1}^2](s_h^k, a_h^k) \right| + \left| [\widehat{\mathbf{w}}_{k,h}^\top \phi(s_h^k, a_h^k)]_{[0,H]}^2 - ([\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k))^2 \right| \\
 &= \left| [\check{\mathbf{w}}_{k,h}^\top \phi(s_h^k, a_h^k)]_{[0,H^2]} - [\mathbb{P}_h V_{k,h+1}^2](s_h^k, a_h^k) \right| \\
 &\quad + \left| [\widehat{\mathbf{w}}_{k,h}^\top \phi(s_h^k, a_h^k)]_{[0,H]} + [\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k) \right| \cdot \left| [\widehat{\mathbf{w}}_{k,h}^\top \phi(s_h^k, a_h^k)]_{[0,H]} - [\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k) \right| \\
 &\leq \min \left\{ \tilde{\beta}_k \|\Sigma_{k,h}^{-1/2} \phi(s_h^k, a_h^k)\|_2, H^2 \right\} + \min \left\{ 2H\bar{\beta}_k \|\Sigma_{k,h}^{-1/2} \phi(s_h^k, a_h^k)\|_2, H^2 \right\} \\
 &= E_{k,h},
 \end{aligned} \tag{D.10}$$

where the first inequality holds due to $|a + b| \leq |a| + |b|$ and the last inequality holds due to Lemma B.1 with the fact that $0 \leq [\widehat{\mathbf{w}}_{k,h}^\top \phi(s_h^k, a_h^k)]_{[0,H]} + [\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k) \leq 2H$. In addition, for the variance $[\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k)$, we have

$$\begin{aligned}
 &|[\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k)| \\
 &= \left| [\mathbb{P}_h V_{k,h+1}^2](s_h^k, a_h^k) - ([\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k))^2 - [\mathbb{P}_h (V_{h+1}^*)^2](s_h^k, a_h^k) + ([\mathbb{P}_h V_{h+1}^*](s_h^k, a_h^k))^2 \right| \\
 &\leq \left| [\mathbb{P}_h V_{k,h+1}^2](s_h^k, a_h^k) - [\mathbb{P}_h (V_{h+1}^*)^2](s_h^k, a_h^k) \right| + \left| ([\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k))^2 - ([\mathbb{P}_h V_{h+1}^*](s_h^k, a_h^k))^2 \right| \\
 &= \left| [\mathbb{P}_h (V_{k,h+1} - V_{h+1}^*)(V_{k,h+1} + V_{h+1}^*)](s_h^k, a_h^k) \right| \\
 &\quad + \left| ([\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{P}_h V_{h+1}^*](s_h^k, a_h^k)) \cdot ([\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k) + [\mathbb{P}_h V_{h+1}^*](s_h^k, a_h^k)) \right| \\
 &\leq 4H([\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{P}_h V_{h+1}^*](s_h^k, a_h^k)),
 \end{aligned} \tag{D.11}$$

where the first inequality holds due to $|a + b| \leq |a| + |b|$ and the last inequality holds due to Lemma B.4 ($V_{k,h+1}(s') \geq V_{h+1}^*(s')$) with the fact that $0 \leq V_{h+1}^*(s'), V_{k,h+1}(s') \leq H$. Based on the event \mathcal{E} and $\tilde{\mathcal{E}}_{h+1}$, (D.11) can be further bounded by

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

$$\begin{aligned}
 &\leq ([\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{P}_h \check{V}_{k,h+1}](s_h^k, a_h^k)) \\
 &\leq \widehat{\mathbf{w}}_{k,h}^\top \phi(s, a) + \bar{\beta} \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)} - \check{\mathbf{w}}_{k,h}^\top \phi(s, a) + \bar{\beta} \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)}, \tag{D.12}
 \end{aligned}$$

where the first inequality holds due to Lemma B.4 ($V_{h+1}^*(s') \geq \check{V}_{k,h+1}(s')$) and the last inequality holds due to the definition of events \mathcal{E} . Combining the results in (D.10), (D.11) and (D.10), we have

$$\begin{aligned}
 &|[\bar{\mathbb{V}}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k)| \\
 &\leq |[\bar{\mathbb{V}}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k)| + |[\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k)| \\
 &\leq E_{k,h} + 4H \left(\widehat{\mathbf{w}}_{k,h}^\top \phi(s, a) + \bar{\beta} \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)} - \check{\mathbf{w}}_{k,h}^\top \phi(s, a) + \bar{\beta} \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)} \right).
 \end{aligned}$$

In addition, since the value functions $V_{k,h}(s)$ and $V_{h+1}^*(s)$ is upper bounded by H , we have $|[\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k)|$ which implies that

$$\begin{aligned}
 &|[\bar{\mathbb{V}}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k)| \\
 &\leq |[\bar{\mathbb{V}}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k)| + |[\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k)| \\
 &\leq E_{k,h} + H^2.
 \end{aligned}$$

Thus, we finish the proof of Lemma B.2. \square

D.3 Proof of Lemma B.3

In this subsection, we provide the proof of Lemma B.3 for the variance estimator.

Proof of Lemma B.3. On the event \mathcal{E} and $\tilde{\mathcal{E}}_{h+1}$, we have

$$\begin{aligned}
 &[\mathbb{V}_h (V_{i,h+1} - V_{h+1}^*)](s_h^k, a_h^k) \\
 &\leq [\mathbb{P}_h (V_{i,h+1} - V_{h+1}^*)^2](s_h^k, a_h^k) \\
 &\leq 2H [\mathbb{P}_h (V_{i,h+1} - V_{h+1}^*)](s_h^k, a_h^k) \\
 &\leq 2H ([\mathbb{P}_h V_{i,h+1}](s_h^k, a_h^k) - [\mathbb{P}_h \check{V}_{k,h+1}](s_h^k, a_h^k)) \\
 &\leq 2H ([\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{P}_h \check{V}_{k,h+1}](s_h^k, a_h^k)) \\
 &\leq 2H \left(\widehat{\mathbf{w}}_{k,h}^\top \phi(s_h^k, a_h^k) + \bar{\beta} \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)} - \check{\mathbf{w}}_{k,h}^\top \phi(s_h^k, a_h^k) + \bar{\beta} \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)} \right),
 \end{aligned}$$

where the first inequality holds due to $\text{Var}(x) \leq \mathbb{E}[x^2]$, the second and third inequality holds due to Lemma B.4 with the fact that $0 \leq V_{i,h+1}(s'), V_{h+1}^*(s') \leq H$, the fourth inequality the fact $V_{k,h+1} \geq V_{i,h+1}$ from the update-rule in Algorithm 1 and the fifth inequality holds due to Lemma B.1. On the other hand, since value function $0 \leq V_{i,h+1}(s'), V_{h+1}^*(s') \leq H$, we have

$$[\mathbb{V}_h (V_{i,h+1} - V_{h+1}^*)](s_h^k, a_h^k) \leq H^2 = (d^3 H^3)/(d^3 H).$$

Thus, we finish the proof of Lemma B.3. \square

D.4 Proof of Lemma B.4

In this subsection, we provide proof of optimistic property.

Proof of Lemma B.4. We prove this lemma by induction. First, we prove the base case for the last stage $H + 1$. Under this situation, for all state $s \in \mathcal{S}$ and action $a \in \mathcal{A}$, we have $Q_{k,H+1}(s, a) = Q_h^*(s, a) = \check{Q}_{k,h}(s, a) = 0$ and $V_{k,h}(s) \geq V_h^*(s) \geq \check{V}_{k,h}(s) = 0$. Thus, the results in Lemma B.4 holds for stage $H + 1$.

Now, we focus on stage $h + 1$. Since events $\tilde{\mathcal{E}}_h$ directly implies the events $\tilde{\mathcal{E}}_{h+1}$, according to the reduction assumption, we have

$$V_{k,h+1}(s) \geq V_{h+1}^*(s) \geq \check{V}_{k,h}(s). \tag{D.13}$$

Thus, for all episode $k \in [K]$, we have

$$r_h(s, a) + \widehat{\mathbf{w}}_{k,h}^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)} - Q_h^*(s, a) \geq [\mathbb{P}_h(V_{k,h+1} - V_{h+1}^*)](s, a) \geq 0,$$

where the first inequality holds due to the definition of events $\tilde{\mathcal{E}}_h$ and the second inequality holds due to (D.13). Furthermore, the optimal value function is upper bounded by $Q_h^*(s, a) \leq H$ and it implies that

$$Q_h^*(s, a) \leq \min \left\{ \min_{1 \leq i \leq k} r_h(s, a) + \widehat{\mathbf{w}}_{i,h}^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \Sigma_{i,h}^{-1} \phi(s, a)}, H \right\} \leq Q_{k,h}(s, a). \quad (\text{D.14})$$

With a similar argument, for the pessimistic action-value function $\check{Q}_{k,h}(s, a)$, we have

$$r_h(s, a) + \check{\mathbf{w}}_{k,h}^\top \phi(s, a) - \bar{\beta} \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)} - Q_h^*(s, a) \leq [\mathbb{P}_h(\check{V}_{k,h+1} - V_{h+1}^*)](s, a) \leq 0.$$

Since the optimal value function is lower bounded by $Q_h^*(s, a) \geq 0$, the result further implies that

$$Q_h^*(s, a) \geq \max \left\{ \max_{1 \leq i \leq k} r_h(s, a) + \widehat{\mathbf{w}}_{\text{klast},h}^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \Sigma_{\text{klast},h}^{-1} \phi(s, a)}, 0 \right\} \geq \check{Q}_{k,h}(s, a). \quad (\text{D.15})$$

In addition, for the value function V , we have

$$\begin{aligned} V_{k,h}(s) &= \max_a Q_{i,h}(s, a) \geq \min_{1 \leq i \leq k} \max_a Q_h^*(s, a) = V_h^*(s), \\ \check{V}_{k,h}(s) &= \max_a Q_{i,h}(s, a) \leq \max_{1 \leq i \leq k} \max_a Q_h^*(s, a) = V_h^*(s), \end{aligned}$$

where the first inequality holds due to (D.14) and the second inequality holds due to (D.15). Thus, by induction, we finish the proof of Lemma B.4. \square

D.5 Proof of Lemma B.5

In this subsection, we provide the proof of Lemma B.5, which suggests a Bernstein-type upper bound for the estimation error.

Proof of Lemma B.5. We prove Lemma B.5 by induction. First, we prove the base case for the last stage H . Under this situation, the weight vector $\widehat{\mathbf{w}}_{k,h} = 0$ and $V_{k,h+1}(s, a) = 0$. Thus, the result in Lemma B.5 holds for stage H .

For stage $h \in [H]$ and $k \in [K]$, according to Lemma G.1, there exists a vector $\mathbf{w}_{k,h}$ such that $\mathbb{P}_h V_{k,h+1}(s, a)$ can be represented by $\mathbf{w}_{k,h}^\top \phi(s, a)$ and $\|\mathbf{w}_{k,h}\|_2 \leq H\sqrt{d}$. Conditioned on the event $\tilde{\mathcal{E}}_{h+1}$, the estimation error can be decomposed as

$$\begin{aligned} & \|\widehat{\mathbf{w}}_{k,h} - \mathbf{w}_{k,h}\|_{\Sigma_{k,h}} \\ &= \left\| \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) V_{k,h+1}(s_{h+1}^i) - \Sigma_{k,h}^{-1} \left(\lambda \mathbf{I} + \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top \right) \mathbf{w}_{k,h} \right\|_{\Sigma_{k,h}} \\ &= \left\| \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (V_{k,h+1}(s_{h+1}^i) - [\mathbb{P}_h V_{k,h+1}](s_h^i, a_h^i)) - \lambda \Sigma_{k,h}^{-1} \mathbf{w}_{k,h} \right\|_{\Sigma_{k,h}} \\ &\leq \underbrace{\|\lambda \Sigma_{k,h}^{-1} \mathbf{w}_{k,h}\|_{\Sigma_{k,h}}}_{I_1} + \underbrace{\left\| \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (V_{k,h+1}(s_{h+1}^i) - [\mathbb{P}_h V_{k,h+1}](s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}}}_{I_2}, \end{aligned} \quad (\text{D.16})$$

where the first inequality holds due to the fact that $\|\mathbf{a} + \mathbf{b}\|_{\Sigma} \leq \|\mathbf{a}\|_{\Sigma} + \|\mathbf{b}\|_{\Sigma}$. For the first term I_1 , since $\Sigma_{k,h} \succeq \lambda \mathbf{I}$ and $\|\mathbf{w}_{k,h}\|_2 \leq H\sqrt{d}$, it is upper bounded by

$$I_1 = \|\lambda \mathbf{w}_{k,h}\|_{\Sigma_{k,h}^{-1}} \leq \sqrt{\lambda} \cdot \|\mathbf{w}_{k,h}\|_2 \leq H\sqrt{d\lambda}. \quad (\text{D.17})$$

For the second term I_2 , we have

$$\begin{aligned}
 I_2 &= \left\| \sum_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (V_{k,h+1}(s_{h+1}^i) - [\mathbb{P}_h V_{k,h+1}](s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}} \\
 &= \left\| \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (V_{k,h+1}(s_{h+1}^i) - [\mathbb{P}_h V_{k,h+1}](s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}^{-1}} \\
 &\leq \underbrace{\left\| \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (V_{h+1}^*(s_{h+1}^i) - [\mathbb{P}_h V_{h+1}^*](s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}^{-1}}}_{J_1} \\
 &\quad + \underbrace{\left\| \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (\Delta V_{k,h+1}(s_{h+1}^i) - [\mathbb{P}_h(\Delta V_{k,h+1})](s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}^{-1}}}_{J_2}, \tag{D.18}
 \end{aligned}$$

where $\Delta V_{k,h+1} = V_{k,h+1} - V_{h+1}^*$.

For the term J_1 , we apply Lemma 4.1 with $\mathbf{x}_i = \bar{\sigma}_{i,h}^{-1} \phi(s_h^i, a_h^i)$ and $\eta_i = \mathbb{1} \{ [\mathbb{V}_h V_{h+1}^*](s_h^i, a_h^i) \leq \bar{\sigma}_{i,h}^2 \} \cdot \bar{\sigma}_{i,h}^{-1} (V_{h+1}^*(s_{h+1}^i) - [\mathbb{P}_h V_{h+1}^*](s_h^i, a_h^i))$. For \mathbf{x}_t, η_t , we have the following property:

$$\begin{aligned}
 \|\mathbf{x}_i\|_2 &= \|\bar{\sigma}_{i,h}^{-1} \phi(s_h^i, a_h^i)\|_2 \leq \|\phi(s_h^i, a_h^i)\|_2 / \sqrt{H} \leq 1/\sqrt{H}, \\
 \mathbb{E}[\eta_i | \mathcal{F}_i] &= 0, |\eta_t| \leq \left| \bar{\sigma}_{i,h}^{-1} (V_{h+1}^*(s_{h+1}^i) - [\mathbb{P}_h V_{h+1}^*](s_h^i, a_h^i)) \right| \leq \sqrt{H}, \\
 \mathbb{E}[\eta_i^2 | \mathcal{F}_i] &= \mathbb{E} \left[\mathbb{1} \{ [\mathbb{V}_h V_{h+1}^*](s_h^i, a_h^i) \leq \bar{\sigma}_{i,h}^2 \} \cdot \bar{\sigma}_{i,h}^{-2} [\mathbb{V}_h V_{h+1}^*](s_h^i, a_h^i) \right] \leq 1, \\
 \max_i \{ |\eta_i| \cdot \min\{1, \|\mathbf{x}_i\|_{\Sigma_{i,h}^{-1}}\} \} &\leq 2H \bar{\sigma}_{i,h}^{-1} \|\mathbf{x}_i\|_{\Sigma_{i,h}^{-1}} \leq \sqrt{d}.
 \end{aligned}$$

Thus, with probability at least $1 - \delta/H$, for all $k \in [K]$, we have

$$\left\| \sum_{i=1}^{k-1} \mathbf{x}_i \eta_i \right\|_{\Sigma_{k,h}^{-1}} \leq O\left(\sqrt{d \log^2(1 + dKH/(\delta\lambda))}\right).$$

In addition, on the event $\tilde{\mathcal{E}}_{h+1}$ and \mathcal{E} , according to Lemma B.2, we have

$$\bar{\sigma}_{k,h}^2 \geq [\bar{\mathbb{V}}_{k,h} V_{k,h+1}](s_h^k, a_h^k) + E_{k,h} + D_{k,h} \geq [\mathbb{V}_h V_{h+1}^*](s_h^k, a_h^k),$$

which further implies that

$$\begin{aligned}
 J_1 &= \left\| \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (V_{h+1}^*(s_{h+1}^i) - [\mathbb{P}_h V_{h+1}^*](s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}^{-1}} \\
 &= \left\| \sum_{i=1}^{k-1} \mathbf{x}_i \eta_i \right\|_{\Sigma_{k,h}^{-1}} \\
 &\leq O\left(\sqrt{d \log^2(1 + dKH/(\delta\lambda))}\right). \tag{D.19}
 \end{aligned}$$

For the term J_2 , we can not directly use Lemma 4.1, Since the stochastic noise $(\Delta V_{k,h+1}(s_{h+1}^i) - [\mathbb{P}_h(\Delta V_{k,h+1})](s_h^i, a_h^i))$ is not \mathcal{F}_{i+1} measurable. Thus, we need to use the ϵ -net covering argument. In detail, for each episode, i , the value function $V_{i,h}$ belongs to the optimistic value function class \mathcal{V} . If we set $\epsilon = \sqrt{\lambda}/(4H^2 d^2 K)$, then according to Lemma F.6, the covering entropy for function class $\mathcal{V} - V_{h+1}^*$ is upper bounded by

$$\log \mathcal{N}_\epsilon \leq O(d^3 H^2 \log^2(dHK/\lambda)). \tag{D.20}$$

Then for function $V_{k,h}$, there must exist a function \tilde{V} in the ϵ -net, such that

$$\text{dist}(\Delta V_{k,h}, \tilde{V}) \leq \epsilon. \quad (\text{D.21})$$

Therefore, the variance of function \tilde{V} is upper bounded by

$$\begin{aligned} & [\mathbb{V}_h \tilde{V}](s_h^i, a_h^i) - [\mathbb{V}_h(\Delta V_{k,h+1})](s_h^i, a_h^i) \\ &= [\mathbb{P}_h \tilde{V}^2](s_h^i, a_h^i) - [\mathbb{P}_h(\Delta V_{k,h+1})^2](s_h^i, a_h^i) + \left([\mathbb{P}_h(\Delta V_{k,h+1})](s_h^i, a_h^i) \right)^2 - (\mathbb{P}_h \tilde{V}(s_h^i, a_h^i))^2 \\ &\leq 2\text{dist}(\Delta V_{k,h}, \tilde{V}) \cdot \max_{s'} |\Delta V_{k,h+1} + \tilde{V}(s')| \\ &\leq 4H \cdot \text{dist}(\Delta V_{k,h}, \tilde{V}) \\ &\leq 1/d^2, \end{aligned} \quad (\text{D.22})$$

where the first inequality holds due to the definition of distance between different functions, the third inequality holds since $|\Delta V_{k,h+1}(s') + \tilde{V}(s')| \leq 2H$ and the last inequality holds due to the definition of ϵ -net. Thus, we apply Lemma 4.1 with $\mathbf{x}_i = \bar{\sigma}_{i,h}^{-1} \phi(s_h^i, a_h^i)$ and $\eta_i = \mathbb{1} \{ [\mathbb{V}_h \tilde{V}](s_h^i, a_h^i) \leq \bar{\sigma}_{i,h}^2 / (d^3 H) \} \cdot \bar{\sigma}_{i,h}^{-1} (\tilde{V}(s_{h+1}^i) - [\mathbb{P}_h \tilde{V}](s_h^i, a_h^i))$. Therefore, For \mathbf{x}_t, η_t , we have the following property:

$$\begin{aligned} \|\mathbf{x}_i\|_2 &= \|\bar{\sigma}_{i,h}^{-1} \phi(s_h^i, a_h^i)\|_2 \leq \|\phi(s_h^i, a_h^i)\|_2 / \sqrt{H} \leq 1/\sqrt{H}, \\ \mathbb{E}[\eta_i | \mathcal{F}_i] &= 0, |\eta_t| \leq \left| \bar{\sigma}_{i,h}^{-1} (V_{h+1}^*(s_{h+1}^i) - [\mathbb{P}_h \tilde{V}_{h+1}](s_h^i, a_h^i)) \right| \leq \sqrt{H}, \\ \mathbb{E}[\eta_i^2 | \mathcal{F}_i] &= \mathbb{E} \left[\mathbb{1} \{ [\mathbb{V}_h \tilde{V}](s_h^i, a_h^i) \leq \bar{\sigma}_{i,h}^2 / (d^3 H) \} \cdot \bar{\sigma}_{i,h}^{-2} [\mathbb{V}_h \tilde{V}](s_h^i, a_h^i) \right] \leq 1/(d^3 H), \\ \max_i \{ |\eta_i| \cdot \min\{1, \|\mathbf{x}_i\|_{\Sigma_{i,h}^{-1}}\} \} &\leq 2H \bar{\sigma}_{i,h}^{-1} \|\mathbf{x}_i\|_{\Sigma_{i,h}^{-1}} \leq 1/(d^3 H). \end{aligned}$$

After taking a union bound over the ϵ -net, with probability at least $1 - \delta$, we have

$$\left\| \sum_{i=1}^{k-1} \mathbf{x}_i \eta_i \right\|_{\Sigma_{k,h}^{-1}} \leq O\left(\sqrt{d \log^2(1 + dKH/(\delta\lambda))}\right). \quad (\text{D.23})$$

In addition, on the event $\tilde{\mathcal{E}}_{h+1}$ and \mathcal{E} , according to Lemmas B.2 and B.3, we have

$$\begin{aligned} \bar{\sigma}_{i,h}^2 &\geq [\bar{\mathbb{V}}_{i,h} V_{i,h+1}](s_h^k, a_h^k) + E_{i,h} + D_{i,h} + H \\ &\geq D_{i,h} + H \\ &\geq d^3 H [\mathbb{V}_h(\Delta V_{k,h+1})](s_h^i, a_h^i) + H \\ &\geq d^3 H [\mathbb{V}_h \tilde{V}](s_h^i, a_h^i), \end{aligned}$$

For simplicity, we denote $\bar{V} = \Delta V_{k,h+1} - \tilde{V}$ and it further implies that

$$\begin{aligned} J_2 &= \left\| \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (\Delta V(s_{h+1}^i) - [\mathbb{P}_h(\Delta V_{k,h+1})](s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}^{-1}} \\ &\leq 2 \left\| \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (\tilde{V}(s_{h+1}^i) - [\mathbb{P}_h \tilde{V}](s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}^{-1}} \\ &\quad + 2 \left\| \sum_{i=1}^{k-1} \bar{\sigma}_{i,h}^{-2} \phi(s_h^i, a_h^i) (\bar{V}(s_{h+1}^i) - [\mathbb{P}_h \bar{V}](s_h^i, a_h^i)) \right\|_{\Sigma_{k,h}^{-1}} \\ &\leq 2 \left\| \sum_{i=1}^{k-1} \mathbf{x}_i \eta_i \right\|_{\Sigma_{k,h}^{-1}} + 8\epsilon^2 k^2 / \lambda \\ &\leq O\left(\sqrt{d \log^2(1 + dKH/(\delta\lambda))}\right). \end{aligned} \quad (\text{D.24})$$

where the first inequality holds due to $\|\mathbf{a} + \mathbf{b}\|_{\Sigma}^2 \leq 2\|\mathbf{a}\|_{\Sigma}^2 + 2\|\mathbf{b}\|_{\Sigma}^2$, the second inequality holds due to the fact that $|\bar{V}(s')| \leq \epsilon$, $\|\phi(s, a)\|_2 \leq 1$, $\Sigma_{k,h} \geq \lambda \mathbf{I}$, $\bar{\Sigma}_{k,h}^{-1} \leq 1$ and the last inequality holds due to (D.23) with $\epsilon = \sqrt{\lambda}/(4H^2d^2K)$. Substituting the results in (D.17), (D.18), (D.19) and (D.23) into (D.16), we obtain

$$\|\widehat{\mathbf{w}}_{k,h} - \mathbf{w}_{k,h}\|_{\Sigma_{k,h}} \leq I_1 + J_1 + J_2 \leq O\left(H\sqrt{d\lambda} + \sqrt{d\log^2(1 + dKH/(\delta\lambda))}\right) = \beta, \quad (\text{D.25})$$

Therefore, the estimation error is upper bounded by

$$\begin{aligned} |\widehat{\mathbf{w}}_{k,h}^\top \phi(s, a) - \mathbb{P}_h V_{k,h+1}(s, a)| &= |\widehat{\mathbf{w}}_{k,h}^\top \phi(s, a) - \mathbf{w}_{k,h}^\top \phi(s, a)| \\ &\leq \|\widehat{\mathbf{w}}_{k,h} - \mathbf{w}_{k,h}\|_{\Sigma_{k,h}} \cdot \|\phi(s, a)\|_{\Sigma_{k,h}^{-1}} \\ &\leq \beta \sqrt{\phi(s, a)^\top (\Sigma_{k,h})^{-1} \phi(s, a)}, \end{aligned}$$

where the first inequality holds due to Cauchy-Schwartz inequality and the last inequality holds due to (D.25), which implies the results in Lemma B.5 holds for stage h . Therefore, by induction, we finish the proof of Lemma B.5. \square

E Proof of Lemmas in Appendix C

In this section, we provide the proof of Lemmas in Appendix C and we need the following auxiliary Lemma, which is modified from Lemma 4.4 in Zhou and Gu (2022)

Lemma E.1. *For any parameters $\beta' \geq 1$ and $C \geq 1$, the summation of bonuses is upper bounded by*

$$\sum_{k=1}^K \min\left(\beta' \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)}, C\right) \leq 4d^4 H^6 C \iota + 10\beta' d^5 H^4 \iota + 2\beta' \sqrt{2d \sum_{k=1}^K (\sigma_{k,h}^2 + H)},$$

where $\iota = \log(1 + K/(d\lambda))$.

Proof of Lemma E.1. For each stage $h \in [H]$, the summation of bonuses is upper bounded by

$$\begin{aligned} &\sum_{k=1}^K \min\left(\beta' \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)}, C\right) \\ &\leq \sum_{k=1}^K \beta' \min\left(\sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)}, 1\right) + C \sum_{k=1}^K \mathbb{1}\left\{\sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)} \geq 1\right\} \\ &\leq C \sum_{k=1}^K \mathbb{1}\left\{\sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)} \geq 1\right\} + 10\beta' d^5 H^4 \iota + 2\beta' \sqrt{2d \sum_{k=1}^K (\sigma_{k,h}^2 + H)}, \end{aligned} \quad (\text{E.1})$$

where $\iota = \log(1 + K/(d\lambda))$ and the last inequality holds due to Lemma G.6. Now, we only need to estimate the number of episodes where the bonus is larger than 1 and we denote these episodes as $\{k_1, \dots, k_m\}$. For simplicity, we denote

$$\Sigma'_i = \lambda \mathbf{I} + \sum_{j=1}^i \bar{\sigma}_{k_j, h}^2 \phi(s_h^{k_j}, a_h^{k_j}) \phi(s_h^{k_j}, a_h^{k_j})^\top,$$

and we have

$$\sum_{i=1}^m \phi(s_h^{k_i}, a_h^{k_i})^\top \Sigma'_{i-1} \phi(s_h^{k_i}, a_h^{k_i}) \geq \sum_{i=1}^m \phi(s_h^{k_i}, a_h^{k_i})^\top \Sigma_{k_i, h}^{-1} \phi(s_h^{k_i}, a_h^{k_i}) \geq m. \quad (\text{E.2})$$

On the other hand, notice that the estimated variance $\bar{\sigma}_{k,h}^2$ is upper bounded by $4d^4 H^4 / \lambda$, we have

$$\sum_{i=1}^m \phi(s_h^{k_i}, a_h^{k_i})^\top \Sigma'_{i-1} \phi(s_h^{k_i}, a_h^{k_i}) \leq 4d^4 H^4 / \lambda \cdot \sum_{i=1}^m \bar{\sigma}_{k_i, h}^{-2} \phi(s_h^{k_i}, a_h^{k_i})^\top \Sigma'_{i-1} \phi(s_h^{k_i}, a_h^{k_i}) \leq 4d^4 H^6 \iota, \quad (\text{E.3})$$

where $\iota = \log(1 + K/(d\lambda))$ and the last inequality holds due to Lemma G.3. Combining the results in (E.2) and (E.3), we have $m \leq 4d^4 H^6 \beta^2 \iota$, and it further implies that

$$\begin{aligned}
 & \sum_{k=1}^K \min \left(\beta' \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)}, C \right) \\
 & \leq C \sum_{k=1}^K \mathbb{1} \left\{ \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)} \geq 1 \right\} + 10\beta' d^5 H^4 \iota + 2\beta' \sqrt{2d \sum_{k=1}^K (\sigma_{k,h}^2 + H)} \\
 & \leq 4d^4 H^6 C \iota + 10\beta' d^5 H^4 \iota + 2\beta' \sqrt{2d \sum_{k=1}^K (\sigma_{k,h}^2 + H)}.
 \end{aligned}$$

Thus, we finish the proof of Lemma E.1. \square

E.1 Proof of Lemma C.1

Proof of Lemma C.1. For all stage $h \in [H]$ and episode $k \in [K]$, we have

$$\begin{aligned}
 & V_{k,h}(s_h^k) - V_{k,h}^{\pi^k}(s_h^k) \\
 & = Q_{k,h}(s_h^k, a_h^k) - Q_{k,h}^{\pi^k}(s_h^k, a_h^k) \\
 & \leq \min \left(\widehat{\mathbf{w}}_{k_{\text{last},h}}^\top \phi(s, a) + \beta \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k_{\text{last},h}}^{-1} \phi(s_h^k, a_h^k)}, H \right) - [\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k) \\
 & \quad + [\mathbb{P}_h (V_{k,h+1} - V_{k,h+1}^{\pi^k})](s_h^k, a_h^k) \\
 & \leq [\mathbb{P}_h (V_{k,h+1} - V_{k,h+1}^{\pi^k})](s_h^k, a_h^k) + 2 \min \left(\beta \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k_{\text{last},h}}^{-1} \phi(s_h^k, a_h^k)}, H \right) \\
 & \leq [\mathbb{P}_h (V_{k,h+1} - V_{k,h+1}^{\pi^k})](s_h^k, a_h^k) + 4 \min \left(\beta \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)}, H \right) \\
 & = V_{k,h+1}(s_{h+1}^k) - V_{k,h+1}^{\pi^k}(s_{h+1}^k) + [\mathbb{P}_h (V_{k,h+1} - V_{k,h+1}^{\pi^k})](s_h^k, a_h^k) - (V_{k,h+1}(s_{h+1}^k) - V_{k,h+1}^{\pi^k}(s_{h+1}^k)) \\
 & \quad + 4 \min \left(\beta \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)}, H \right), \tag{E.4}
 \end{aligned}$$

where the first inequality holds due to the definition of value function $Q_{k,h}(s_h^k, a_h^k)$, the second inequality holds due to Lemma B.5 and the last inequality holds due to Lemma G.4 with the updating rule (Line 8). Furthermore, for all stage $h \in [H]$, we have

$$\begin{aligned}
 & \sum_{k=1}^K (V_{k,h}(s_h^k) - V_{k,h}^{\pi^k}(s_h^k)) \\
 & \leq \sum_{k=1}^K \sum_{h'=h}^H 4 \min \left(\beta \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)}, H \right) \\
 & \quad + \sum_{k=1}^K \sum_{h'=h}^H \left([\mathbb{P}_h (V_{k,h+1} - V_{k,h+1}^{\pi^k})](s_h^k, a_h^k) - (V_{k,h+1}(s_{h+1}^k) - V_{k,h+1}^{\pi^k}(s_{h+1}^k)) \right) \\
 & \leq \sum_{k=1}^K \sum_{h'=h}^H 4 \min \left(\beta \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)}, H \right) + 4\sqrt{H^3 K \log(H/\delta)} \\
 & \leq 16d^4 H^8 \iota + 40\beta d^7 H^5 \iota + 8\beta \sum_{h'=h}^H \sqrt{2d \sum_{k=1}^K (\sigma_{k,h'}^2 + H)} + 4\sqrt{H^3 K \log(H/\delta)} \\
 & \leq 16d^4 H^8 \iota + 40\beta d^7 H^5 \iota + 8\beta \sqrt{2dH \sum_{h=1}^H \sum_{k=1}^K (\sigma_{k,h}^2 + H)} + 4\sqrt{H^3 K \log(H/\delta)}, \tag{E.5}
 \end{aligned}$$

where the first inequality holds by taking the summation of (E.4) for $k \in [K]$ and $h \leq h' \leq H$, the second inequality holds due to the definition of events \mathcal{E}_1 , the third inequality holds due to Lemma E.1 and the last inequality holds due to Cauchy-Schwartz inequality. Furthermore, taking the summation of (E.5), we have

$$\begin{aligned}
 & \sum_{k=1}^K \sum_{h=1}^H [\mathbb{P}_h(V_{k,h+1} - V_{k,h+1}^{\pi^k})](s_h^k, a_h^k) \\
 &= \sum_{k=1}^K \sum_{h=1}^H (V_{k,h+1}(s_{h+1}^k) - V_{k,h+1}^{\pi^k}(s_{h+1}^k)) \\
 & \quad + \sum_{k=1}^K \sum_{h=1}^H \left([\mathbb{P}_h(V_{k,h+1} - V_{k,h+1}^{\pi^k})](s_h^k, a_h^k) - (V_{k,h+1}(s_{h+1}^k) - V_{k,h+1}^{\pi^k}(s_{h+1}^k)) \right) \\
 &\leq \sum_{k=1}^K \sum_{h=1}^H (V_{k,h+1}(s_{h+1}^k) - V_{k,h+1}^{\pi^k}(s_{h+1}^k)) + 2\sqrt{2H^3K \log(H/\delta)} \\
 &\leq 16d^4H^9\iota + 40\beta d^7H^6\iota + 8H\beta \sqrt{2dH\iota \sum_{h=1}^H \sum_{k=1}^K (\sigma_{k,h}^2 + H) + 4\sqrt{H^5K \log(H/\delta)}},
 \end{aligned}$$

where the first inequality holds due to Lemma G.2 and the last inequality holds due (E.5). Therefore, we finish the proof of Lemma C.1. \square

E.2 Proof of Lemma C.2

Proof of Lemma C.2. For each stage $h \in [H]$ and episode $k \in [K]$, we have

$$\begin{aligned}
 & V_{k,h}(s_h^k) - \check{V}_{k,h}(s_h^k) \\
 &\leq Q_{k,h}(s_h^k, a_h^k) - \check{Q}_{k,h}(s_h^k, a_h^k) \\
 &\leq \min \left(\check{\mathbf{w}}_{k_{\text{last}},h}^\top \phi(s, a) + \beta \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k_{\text{last}},h}^{-1} \phi(s_h^k, a_h^k), H} \right) - [\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k) \\
 & \quad - \max \left(\widehat{\mathbf{w}}_{k_{\text{last}},h}^\top \phi(s, a) - \bar{\beta} \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k_{\text{last}},h}^{-1} \phi(s_h^k, a_h^k), 0} \right) + [\mathbb{P}_h \check{V}_{k,h+1}](s_h^k, a_h^k) \\
 & \quad + [\mathbb{P}_h(V_{k,h+1} - \check{V}_{k,h+1})](s_h^k, a_h^k) \\
 &\leq [\mathbb{P}_h(V_{k,h+1} - \check{V}_{k,h+1})](s_h^k, a_h^k) + 2 \min \left(\beta \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k_{\text{last}},h}^{-1} \phi(s_h^k, a_h^k), H} \right) \\
 & \quad + 2 \min \left(\bar{\beta} \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k_{\text{last}},h}^{-1} \phi(s_h^k, a_h^k), H} \right) \\
 &\leq [\mathbb{P}_h(V_{k,h+1} - \check{V}_{k,h+1})](s_h^k, a_h^k) + 4 \min \left(\beta \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k), H} \right) \\
 & \quad + 4 \min \left(\bar{\beta} \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k), H} \right) \\
 &= V_{k,h+1}(s_{h+1}^k) - \check{V}_{k,h+1}(s_{h+1}^k) + [\mathbb{P}_h(V_{k,h+1} - \check{V}_{k,h+1})](s_h^k, a_h^k) - (V_{k,h+1}(s_{h+1}^k) - \check{V}_{k,h+1}(s_{h+1}^k)) \\
 & \quad + 4 \min \left(\beta \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k), H} \right) + 4 \min \left(\bar{\beta} \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k), H} \right), \tag{E.6}
 \end{aligned}$$

where the first inequality holds due to the fact that $\check{V}_{k,h}(s_h^k) = \max_a \check{Q}_{k,h}(s_h^k, a) \geq \check{Q}_{k,h}(s_h^k, a_h^k)$, the second inequality holds due to the definition of value functions $Q_{k,h}$ and $\check{Q}_{k,h}$, the third inequality holds due to Lemma B.5 and Lemma B.1, and the last inequality holds due to Lemma G.4 with the updating rule (Line 8). Furthermore, for all stage $h \in [H]$, we have

$$\begin{aligned}
 & \sum_{k=1}^K (V_{k,h}(s_h^k) - \check{V}_{k,h}(s_h^k)) \\
 &\leq \sum_{k=1}^K \sum_{h'=h}^H 4 \min \left(\beta \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k), H} \right) + 4 \min \left(\bar{\beta} \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k), H} \right)
 \end{aligned}$$

$$\begin{aligned}
 & + \sum_{k=1}^K \sum_{h'=h}^H \left([\mathbb{P}_h(V_{k,h+1} - \check{V}_{k,h+1})](s_h^k, a_h^k) - (V_{k,h+1}(s_{h+1}^k) - \check{V}_{k,h+1}(s_{h+1}^k)) \right) \\
 \leq & \sum_{k=1}^K \sum_{h'=h}^H 4 \min \left(\beta \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)}, H \right) + 4 \min \left(\bar{\beta} \sqrt{\phi(s_h^k, a_h^k)^\top \Sigma_{k,h}^{-1} \phi(s_h^k, a_h^k)}, H \right) \\
 & + 4\sqrt{H^3 K \log(H/\delta)} \\
 \leq & 32d^4 H^8 \iota + 40(\beta + \bar{\beta})d^7 H^5 \iota + 8(\beta + \bar{\beta}) \sum_{h'=h}^H \sqrt{2d \sum_{k=1}^K (\sigma_{k,h'}^2 + H)} + 4\sqrt{H^3 K \log(H/\delta)} \\
 \leq & 32d^4 H^8 \iota + 40(\beta + \bar{\beta})d^7 H^5 \iota + 8(\beta + \bar{\beta}) \sqrt{2dH \iota \sum_{h=1}^H \sum_{k=1}^K (\sigma_{k,h}^2 + H)} + 4\sqrt{H^3 K \log(H/\delta)}, \tag{E.7}
 \end{aligned}$$

where the first inequality holds by taking the summation of (E.6) for $k \in [K]$ and $h \leq h' \leq H$, the second inequality holds due to the definition of event \mathcal{E}_2 , the third inequality holds due to Lemma E.1 and the last inequality holds due to Cauchy-Schwartz inequality. Furthermore, taking the summation of (E.7), we have

$$\begin{aligned}
 & \sum_{k=1}^K \sum_{h=1}^H [\mathbb{P}_h(V_{k,h+1} - \check{V}_{k,h+1})](s_h^k, a_h^k) \\
 = & \sum_{k=1}^K \sum_{h=1}^H (V_{k,h+1}(s_{h+1}^k) - \check{V}_{k,h+1}(s_{h+1}^k)) \\
 & + \sum_{k=1}^K \sum_{h=1}^H \left([\mathbb{P}_h(V_{k,h+1} - \check{V}_{k,h+1})](s_h^k, a_h^k) - (V_{k,h+1}(s_{h+1}^k) - \check{V}_{k,h+1}(s_{h+1}^k)) \right) \\
 \leq & \sum_{k=1}^K \sum_{h=1}^H (V_{k,h+1}(s_{h+1}^k) - \check{V}_{k,h+1}(s_{h+1}^k)) + 2\sqrt{2H^3 K \log(H/\delta)} \\
 \leq & 32d^4 H^9 \iota + 40(\beta + \bar{\beta})d^7 H^6 \iota + 8H(\beta + \bar{\beta}) \sqrt{2dH \iota \sum_{h=1}^H \sum_{k=1}^K (\sigma_{k,h}^2 + H)} + 4\sqrt{H^5 K \log(H/\delta)},
 \end{aligned}$$

where the first inequality holds due to Lemma G.2 and the last inequality holds due (E.7). Therefore, we finish the proof of Lemma C.2. \square

E.3 Proof of Lemma C.3

Proof of Lemma C.3. According to the definition of estimated variance $\sigma_{k,h}$, we have

$$\begin{aligned}
 \sum_{k=1}^K \sum_{h=1}^H \sigma_{k,h}^2 & = \sum_{k=1}^K \sum_{h=1}^H [\bar{\mathbb{V}}_{k,h} V_{k,h+1}](s_h^k, a_h^k) + E_{k,h} + D_{k,h} + H \\
 & = H^2 K + \underbrace{\sum_{k=1}^K \sum_{h=1}^H ([\bar{\mathbb{V}}_{k,h} V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k))}_{I_1} + \underbrace{\sum_{k=1}^K \sum_{h=1}^H E_{k,h}}_{I_2} + \underbrace{\sum_{k=1}^K \sum_{h=1}^H D_{k,h}}_{I_3} \\
 & \quad + \underbrace{\sum_{k=1}^K \sum_{h=1}^H ([\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{k,h+1}^\pi](s_h^k, a_h^k))}_{I_4} + \underbrace{\sum_{k=1}^K \sum_{h=1}^H [\mathbb{V}_h V_{k,h+1}^\pi](s_h^k, a_h^k)}_{I_5}. \tag{E.8}
 \end{aligned}$$

For the term I_1 , according to Lemma B.2, it is upper bounded by:

$$I_1 = \sum_{k=1}^K \sum_{h=1}^H ([\bar{\mathbb{V}}_{k,h} V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k)) \leq \sum_{k=1}^K \sum_{h=1}^H E_{k,h} = I_2. \tag{E.9}$$

For the term I_2 , it is upper bounded by

$$\begin{aligned}
 I_2 &= \sum_{k=1}^K \sum_{h=1}^H E_{k,h} \\
 &= \sum_{k=1}^K \sum_{h=1}^H \min \left\{ \tilde{\beta} \|\Sigma_{k,h}^{-1/2} \phi(s_h^k, a_h^k)\|_2, H^2 \right\} + \min \left\{ 2H\bar{\beta} \|\Sigma_{k,h}^{-1/2} \phi(s_h^k, a_h^k)\|_2, H^2 \right\} \\
 &\leq 8d^4 H^9 \iota + (10\tilde{\beta} + 20\bar{\beta})d^5 H^5 \iota + (2\tilde{\beta} + 4\bar{\beta})H \sqrt{2d\iota \sum_{h=1}^H \sum_{k=1}^K (\sigma_{k,h}^2 + H)}, \tag{E.10}
 \end{aligned}$$

where $\iota = \log(1 + K/(d\lambda))$ and the inequality holds due to Lemma E.1.

For the term I_3 , it is upper bounded by

$$\begin{aligned}
 I_3 &= \sum_{k=1}^K \sum_{h=1}^H D_{k,h} \\
 &= \sum_{k=1}^K \sum_{h=1}^H \min \left\{ 4d^3 H^2 \left(\widehat{\mathbf{w}}_{k,h}^\top \phi(s, a) - \check{\mathbf{w}}_{k,h}^\top \phi(s, a) + 2\bar{\beta} \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)} \right), d^3 H^3 \right\} \\
 &\leq \sum_{k=1}^K \sum_{h=1}^H \min \left\{ 4d^3 H^2 \left([\mathbb{P}_h(V_{k,h+1} - \check{V}_{k,h+1})](s_h^k, a_h^k) + 4\bar{\beta} \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)} \right), d^3 H^3 \right\} \\
 &\leq \sum_{k=1}^K \sum_{h=1}^H 4d^3 H^2 [\mathbb{P}_h(V_{k,h+1} - \check{V}_{k,h+1})](s_h^k, a_h^k) + \sum_{k=1}^K \sum_{h=1}^H \min \left\{ 16d^3 H^2 \bar{\beta} \sqrt{\phi(s, a)^\top \Sigma_{k,h}^{-1} \phi(s, a)}, d^3 H^3 \right\} \\
 &\leq 4d^7 H^9 \iota + 160\bar{\beta}d^8 H^7 \iota + 32d^3 H^3 \bar{\beta} \sqrt{2d\iota \sum_{h=1}^H \sum_{k=1}^K (\sigma_{k,h}^2 + H)} \\
 &\quad + \sum_{k=1}^K \sum_{h=1}^H 4d^3 H^2 [\mathbb{P}_h(V_{k,h+1} - \check{V}_{k,h+1})](s_h^k, a_h^k) \\
 &\leq 132d^7 H^{11} \iota + 320(\beta + \bar{\beta})d^{10} H^8 \iota + 64d^3 H^3 (\beta + \bar{\beta}) \sqrt{2dH\iota \sum_{h=1}^H \sum_{k=1}^K (\sigma_{k,h}^2 + H) + 4d^3 \sqrt{H^9 K \log(H/\delta)}}, \tag{E.11}
 \end{aligned}$$

where $\iota = \log(1 + K/(d\lambda))$, the first inequality holds due to Lemma B.1, the second inequality holds due to the fact that $V_{k,h+1}(s) \geq V_{h+1}^*(s) \geq \check{V}_{k,h+1}(s)$, the third inequality holds due to Lemma E.1 and the last inequality holds due to Lemma C.2.

For the term I_4 , it is upper bounded by

$$\begin{aligned}
 I_4 &= \sum_{k=1}^K \sum_{h=1}^H ([\mathbb{V}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{V}_h V_{k,h+1}^{\pi^k}](s_h^k, a_h^k)) \\
 &= \sum_{k=1}^K \sum_{h=1}^H \left([\mathbb{P}_h V_{k,h+1}^2](s_h^k, a_h^k) - ([\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k))^2 - [\mathbb{P}_h (V_{k,h+1}^{\pi^k})^2](s_h^k, a_h^k) + ([\mathbb{P}_h V_{k,h+1}^{\pi^k}](s_h^k, a_h^k))^2 \right) \\
 &\leq \sum_{k=1}^K \sum_{h=1}^H ([\mathbb{P}_h V_{k,h+1}^2](s_h^k, a_h^k) - [\mathbb{P}_h (V_{k,h+1}^{\pi^k})^2](s_h^k, a_h^k)) \\
 &\leq 2H \sum_{k=1}^K \sum_{h=1}^H ([\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{P}_h V_{k,h+1}^{\pi^k}](s_h^k, a_h^k))
 \end{aligned}$$

$$\leq 32d^4H^{10}\iota + 80\beta d^7H^7\iota + 16H^2\beta \sqrt{2dH\iota \sum_{h=1}^H \sum_{k=1}^K (\sigma_{k,h}^2 + H)} + 8\sqrt{H^7K \log(H/\delta)}, \quad (\text{E.12})$$

where the first inequality holds due to the fact that $V_{k,h+1}^{\pi^k}(s') \leq V_{k,h+1}(s')$, the second inequality holds due to $0 \leq V_{k,h+1}(s'), V_{k,h+1}^{\pi^k}(s') \leq H$ and the last inequality holds due to Lemma C.1.

Based on the definition of events \mathcal{E}_3 , for the term I_5 , we have

$$I_5 = \sum_{k=1}^K \sum_{h=1}^H [\mathbb{V}_h V_{k,h+1}^{\pi^k}](s_h^k, a_h^k) \leq 3(H^2K + H^3 \log(1/\delta)). \quad (\text{E.13})$$

Substituting the results in (E.9), (E.10), (E.11), (E.12) and (E.13) into (E.8), we have

$$\begin{aligned} & \sum_{k=1}^K \sum_{h=1}^H \sigma_{k,h}^2 \\ &= I_1 + I_2 + I_3 + I_4 + I_5 \\ &\leq 3H^2K + 183d^7H^{11}\iota + 460(\beta + \tilde{\beta} + \bar{\beta})d^{10}H^8\iota + 12d^3\sqrt{H^9K \log(H/\delta)} \\ &\quad + 92d^3H^3(\beta + \tilde{\beta} + \bar{\beta}) \sqrt{2dH\iota \sum_{h=1}^H \sum_{k=1}^K (\sigma_{k,h}^2 + H)} \\ &\leq 3H^2K + 183d^7H^{11}\iota + 460(\beta + \tilde{\beta} + \bar{\beta})d^{10}H^8\iota + 12d^3\sqrt{H^9K \log(H/\delta)} \\ &\quad + 92d^3H^3(\beta + \tilde{\beta} + \bar{\beta})\sqrt{2dH^2K\iota} + 92d^3H^3(\beta + \tilde{\beta} + \bar{\beta}) \sqrt{2dH\iota \sum_{h=1}^H \sum_{k=1}^K \sigma_{k,h}^2}, \end{aligned}$$

where $\iota = \log(1 + K/(d\lambda))$,

$$\begin{aligned} \beta &= O\left(H\sqrt{d\lambda} + \sqrt{d \log^2(1 + dKH/(\delta\lambda))}\right) \\ \tilde{\beta} &= O\left(H^2\sqrt{d\lambda} + \sqrt{d^3H^4 \log^2(dHK/(\delta\lambda))}\right) \\ \bar{\beta} &= O\left(H\sqrt{d\lambda} + \sqrt{d^3H^2 \log^2(dHK/(\delta\lambda))}\right), \end{aligned}$$

and the last inequality holds due to the fact that $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$. Therefore, by the fact that $x \leq a\sqrt{x} + b$ implies $x \leq a^2 + 2b$ and $\lambda = 1/H^2$, we have

$$\sum_{k=1}^K \sum_{h=1}^H \sigma_{k,h}^2 \leq O(H^2K + d^{10.5}H^{16} \log^{1.5}(1 + dKH/\delta)).$$

Thus, we complete the proof of Lemma C.3. \square

F Covering Number Arguments

F.1 Number of Value Function Updating

According to the determinant-based criterion in Algorithm 1 (Line 8), the number of episodes where the algorithm updates the value function is upper bounded by:

Lemma F.1. *The number of episodes where the algorithm updates the value function in Algorithm 1 is upper bounded by $dH \log(1 + K/\lambda)$.*

Proof. We denote $k_0 = 0$ and suppose that $\{k_1, \dots, k_m\}$ be the episodes where the algorithm updates the value function. Then according to the determinant-based criterion (Line 8), for each episodes k_i , there exists a stage $h \in [H]$ such that

$$\det(\Sigma_{k_i,h}) \geq 2 \det(\Sigma_{k_{i-1},h}).$$

According to the update rule of $\Sigma_{k,h}$ (Line 24), for other stage $h' \neq h$, we have $\Sigma_{k_i,h'} \succeq \Sigma_{k_i,h'}$, which implies $\det(\Sigma_{k_i,h'}) \geq \det(\Sigma_{k_i,h})$. Thus, we have

$$\prod_{h=1}^H \det(\Sigma_{k_i,h}) \geq 2 \prod_{h=1}^H \det(\Sigma_{k_{i-1},h}). \quad (\text{F.1})$$

Applying the result (F.1) overall episodes in $\{k_1, \dots, k_m\}$, we have

$$\prod_{h=1}^H \det(\Sigma_{k_m,h}) \geq 2^m \prod_{h=1}^H \det(\Sigma_{k_0,h}) = 2^m \prod_{h=1}^H \det(\lambda \mathbf{I}) = 2^m \lambda^{dH}. \quad (\text{F.2})$$

On the other hand, the determinant $\det(\Sigma_{k_m,h})$ is upper bounded by:

$$\prod_{h=1}^H \det(\Sigma_{k_m,h}) \leq \prod_{h=1}^H \det(\Sigma_{K,h}) \leq (\lambda + K)^{dH}, \quad (\text{F.3})$$

where the first inequality holds due to $\Sigma_{K,h} \succeq \Sigma_{k_m,h}$, the last inequality holds due to $\bar{\sigma}_{k,h}^{-1} \leq 1$ and $\|\phi(s, a)\|_2 \leq 1$. Combining the results in (F.2) and (F.3), we have

$$m \leq dH \log(1 + K/\lambda).$$

Thus, we finish the proof of Lemma F.1. \square

F.2 Norm of the Weight Vectors

In this section, we provide the following upper bounds for the norm of the weight vectors.

Lemma F.2. *For all stage $h \in [H]$ and all episode $k \in \mathbb{N}$, the norm of the weight vector $\widehat{\mathbf{w}}_{k,h}$ can be upper bounded by*

$$\|\widehat{\mathbf{w}}_{k,h}\|_2 \leq H\sqrt{dK/\lambda}.$$

Proof of Lemma F.2. According to the definition of weight vector $\widehat{\mathbf{w}}_{k,h}$ in Algorithm 1, we have

$$\begin{aligned} \Sigma_{k,h} &= \lambda \mathbf{I} + \sum_{i=1}^{k-1} \bar{\sigma}_{k,i}^{-2} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^\top, \\ \widehat{\mathbf{w}}_{k,h} &= \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{k,i}^{-2} \phi(s_h^i, a_h^i) V_{k,h+1}(s_{h+1}^i). \end{aligned}$$

Then for the norm $\|\widehat{\mathbf{w}}_{k,h}\|_2$, we have the following inequality

$$\begin{aligned} \|\widehat{\mathbf{w}}_{k,h}\|_2^2 &= \left\| \Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{k,i}^{-2} \phi(s_h^i, a_h^i) V_{k,h+1}(s_{h+1}^i) \right\|_2^2 \\ &\leq k \sum_{i=1}^{k-1} \left\| \Sigma_{k,h}^{-1} \bar{\sigma}_{k,i}^{-2} \phi(s_h^i, a_h^i) V_{k,h+1}(s_{h+1}^i) \right\|_2^2 \\ &\leq kH^2 \sum_{i=1}^{k-1} \bar{\sigma}_{k,i}^{-2} \left\| \Sigma_{k,h}^{-1} \phi(s_h^i, a_h^i) \right\|_2^2 \\ &\leq \frac{kH^2}{\lambda} \sum_{i=1}^{k-1} \bar{\sigma}_{k,i}^{-2} \phi(s_h^i, a_h^i)^\top \Sigma_{k,h}^{-1} \phi(s_h^i, a_h^i) \\ &= \frac{kH^2}{\lambda} \text{tr} \left(\Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{k,i}^{-2} \phi(s_h^i, a_h^i)^\top \phi(s_h^i, a_h^i) \right), \end{aligned} \quad (\text{F.4})$$

where the first inequality holds due to Cauchy-Schwartz inequality, the second inequality holds due to $0 \leq V_{k,h+1}(s,a) \leq H$ and the last inequality holds due to $\Sigma_{k,h} \succeq \lambda \mathbf{I}$. Now, we assume the eigen-decomposition of matrix $\sum_{i=1}^{k-1} \bar{\sigma}_{k,i}^{-2} \phi(s_h^i, a_h^i)^\top \phi(s_h^i, a_h^i)$ is $Q^\top \Lambda Q$ and we have

$$\begin{aligned} \text{tr} \left(\Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{k,i}^{-2} \phi(s_h^i, a_h^i)^\top \phi(s_h^i, a_h^i) \right) &= \text{tr} \left((Q^\top \Lambda Q + \lambda \mathbf{I})^{-1} Q^\top \Lambda Q \right) \\ &= \text{tr} \left((\Lambda + \lambda \mathbf{I})^{-1} \Lambda \right) \\ &= \sum_{i=1}^d \frac{\Lambda_i}{\Lambda_i + \lambda} \\ &\leq d. \end{aligned} \tag{F.5}$$

Substituting (F.5) into (F.4), we have

$$\begin{aligned} \|\widehat{\mathbf{w}}_{k,h}\|_2^2 &\leq \frac{kH^2}{\lambda} \text{tr} \left(\Sigma_{k,h}^{-1} \sum_{i=1}^{k-1} \bar{\sigma}_{k,i}^{-2} \phi(s_h^i, a_h^i)^\top \phi(s_h^i, a_h^i) \right) \\ &\leq \frac{kH^2 d}{\lambda}, \end{aligned} \tag{F.6}$$

where the first inequality holds due to (F.4) and the last inequality holds due to (F.5). Thus, we finish the proof of Lemma F.2 \square

In addition, for the pessimistic weight vector $\check{\mathbf{w}}_{k,h}$ and weight vector $\widetilde{\mathbf{w}}_{k,h}$, we have the following lemmas:

Lemma F.3. For all stage $h \in [H]$ and all episode $k \in \mathbb{N}$, the norm of the weight vector $\check{\mathbf{w}}_{k,h}$ can be upper bounded by

$$\|\check{\mathbf{w}}_{k,h}\|_2 \leq H \sqrt{dK/\lambda}.$$

Proof of Lemma F.3. The proof is almost the same as Lemma F.2 and we only need to replace the optimistic value function $V_{k,h}(s,a)$ by the pessimistic value function $\check{V}_{k,h}(s,a)$. \square

Lemma F.4. For each stage $h \in [H]$ and each episode $k \in \mathbb{N}$, the norm of the weight vector $\widetilde{\mathbf{w}}_{k,h}$ can be upper bounded by

$$\|\widetilde{\mathbf{w}}_{k,h}\|_2 \leq H^2 \sqrt{dK/\lambda}.$$

Proof of Lemma F.4. The proof is almost the same as Lemma F.2 and we only need to replace the optimistic value function $V_{k,h}(s,a)$ with the squared value function $V_{k,h}^2(s,a)$. \square

F.3 Function Class and Covering Number

Combining the update rule (Line 8) with Lemma F.1 and Lemma F.2, for each episodes $k \in [K]$ and $h \in [H]$, the optimistic value function $V_{k,h} = \min_{i \leq k} \max_a Q_{i,h}(s,a)$ belong to the following function class

$$\mathcal{V}_h = \left\{ V \mid V(\cdot) = \max_a \min_{1 \leq i \leq l} \min \left(H, r_h(\cdot, a) + \mathbf{w}_i^\top \phi(\cdot, a) + \beta \sqrt{\phi(\cdot, a)^\top \Sigma_i^{-1} \phi(\cdot, a)} \right), \|\mathbf{w}_i\| \leq L, \Sigma_i \succeq \lambda \mathbf{I} \right\}, \tag{F.7}$$

where $l = dH \log(1 + K/\lambda)$ and $L = H \sqrt{dK/\lambda}$. Similarly, for each episodes $k \in [K]$ and $h \in [H]$, the pessimistic value function $\check{V}_{k,h} = \max_{i \leq k} \max_a \check{Q}_{i,h}(s,a)$ belongs to the following function class

$$\check{\mathcal{V}}_h = \left\{ V \mid V(\cdot) = \max_a \max_{1 \leq i \leq l} \max \left(0, r_h(\cdot, a) + \mathbf{w}_i^\top \phi(\cdot, a) - \beta \sqrt{\phi(\cdot, a)^\top \Sigma_i^{-1} \phi(\cdot, a)} \right), \|\mathbf{w}_i\| \leq L, \Sigma_i \succeq \lambda \mathbf{I} \right\}, \tag{F.8}$$

where $l = dH \log(1 + K/\lambda)$ and $L = H \sqrt{dK/\lambda}$. To compute the covering number of function classes \mathcal{V}_h , \mathcal{V}_h^2 and $\check{\mathcal{V}}_h$, we need the following result on the Euclidean ball.

Lemma F.5 (Lemma D.5, Jin et al. 2020). *For a Euclidean ball with radius R in \mathbb{R}^d , the ϵ -covering number of this ball is upper bounded by $(1 + 2R/\epsilon)^d$.*

With the help of Lemma F.5, the covering number \mathcal{N}_ϵ of optimistic function class \mathcal{V} can be upper bounded by the following lemma:

Lemma F.6. *For optimistic function class \mathcal{V}_h , we define the distance between two function V_1 and V_2 as $V_1, V_2 \in \mathcal{V}_h$ as $\text{dist}(V_1, V_2) = \max_s |V_1(s) - V_2(s)|$. With respect to this distance function, the ϵ -covering number \mathcal{N}_ϵ of the function class \mathcal{V}_h can be upper bounded by*

$$\log \mathcal{N}_\epsilon \leq dl \log(1 + 4L/\epsilon) + d^2 l \log(1 + 8\sqrt{d}\beta^2/\epsilon^2)$$

Proof of Lemma F.6. For any two function $V_1, V_2 \in \mathcal{V}_h$, according to the definition of function class \mathcal{V}_h , we have

$$\begin{aligned} V_1(\cdot) &= \max_a \min_{1 \leq i \leq l} \min \left(H, r_h(\cdot, a) + \mathbf{w}_{1,i}^\top \phi(\cdot, a) + \beta \sqrt{\phi(\cdot, a)^\top \mathbf{\Gamma}_{1,i} \phi(\cdot, a)} \right), \\ V_2(\cdot) &= \max_a \min_{1 \leq i \leq l} \min \left(H, r_h(\cdot, a) + \mathbf{w}_{2,i}^\top \phi(\cdot, a) + \beta \sqrt{\phi(\cdot, a)^\top \mathbf{\Gamma}_{2,i} \phi(\cdot, a)} \right), \end{aligned}$$

where $\|\mathbf{w}_{1,i}\|_2, \|\mathbf{w}_{2,i}\|_2 \leq L$ and $\mathbf{\Gamma}_{1,i}, \mathbf{\Gamma}_{2,i} \preceq \mathbf{I}$. Since all of the functions $\min_{1 \leq i \leq l}, \max_a$ and $\min(H, \cdot)$ are contraction functions, we have

$$\begin{aligned} \text{dist}(V_1, V_2) &= \max_{s \in \mathcal{S}} |V_1(s) - V_2(s)| \\ &\leq \max_{1 \leq i \leq l, s \in \mathcal{S}, a \in \mathcal{A}} \left| \mathbf{w}_{1,i}^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \mathbf{\Gamma}_{1,i} \phi(s, a)} \right. \\ &\quad \left. - \mathbf{w}_{2,i}^\top \phi(s, a) - \beta \sqrt{\phi(s, a)^\top \mathbf{\Gamma}_{2,i} \phi(s, a)} \right| \\ &\leq \beta \max_{1 \leq i \leq l, s \in \mathcal{S}, a \in \mathcal{A}} \left| \sqrt{\phi(s, a)^\top \mathbf{\Gamma}_{1,i} \phi(s, a)} - \sqrt{\phi(s, a)^\top \mathbf{\Gamma}_{2,i} \phi(s, a)} \right| \\ &\quad + \max_{1 \leq i \leq l, s \in \mathcal{S}, a \in \mathcal{A}} |(\mathbf{w}_{1,i} - \mathbf{w}_{2,i})^\top \phi(s, a)| \\ &\leq \beta \max_{1 \leq i \leq l, s \in \mathcal{S}, a \in \mathcal{A}} \left| \sqrt{\phi(s, a)^\top (\mathbf{\Gamma}_{1,i} - \mathbf{\Gamma}_{2,i}) \phi(s, a)} \right| \\ &\quad + \max_{1 \leq i \leq l, s \in \mathcal{S}, a \in \mathcal{A}} |(\mathbf{w}_{1,i} - \mathbf{w}_{2,i})^\top \phi(s, a)| \\ &\leq \beta \max_{1 \leq i \leq l} \sqrt{\|\mathbf{\Gamma}_{1,i} - \mathbf{\Gamma}_{2,i}\|_F} + \max_{1 \leq i \leq l} \|\mathbf{w}_{1,i} - \mathbf{w}_{2,i}\|_2, \end{aligned} \tag{F.9}$$

where the first inequality holds due to the contraction property, the second inequality holds due to the fact that $\max_x |f(x) + g(x)| \leq \max_x |f(x)| + \max_x |g(x)|$, the third inequality holds due to $|\sqrt{x} - \sqrt{y}| \geq |\sqrt{x} - \sqrt{y}|$ and the last inequality holds due to the fact that $\|\phi(s, a)\|_2 \leq 1$. Now, we denote $\mathcal{C}_\mathbf{w}$ as a $\epsilon/2$ -cover of the set $\{\mathbf{w} \in \mathbb{R}^d \mid \|\mathbf{w}\|_2 \leq L\}$ and \mathcal{C}_Γ as a $\epsilon^2/(4\beta^2)$ -cover of the set $\{\mathbf{\Gamma} \in \mathbb{R}^{d \times d} \mid \|\mathbf{\Gamma}\|_F \leq \sqrt{d}\}$ with respect to the Frobenius norms. Thus, according to Lemma F.5, we have following property:

$$|\mathcal{C}_\mathbf{w}| \leq (1 + 4L/\epsilon)^d, |\mathcal{C}_\Gamma| \leq (1 + 8\sqrt{d}\beta^2/\epsilon^2)^{d^2}. \tag{F.10}$$

By the definition of covering number, for any function $V_1 \in \mathcal{V}$ with parameters $\mathbf{w}_{1,i}, \mathbf{\Gamma}_{1,i} (1 \leq i \leq l)$, s other parameters $\mathbf{w}_{2,i}, \mathbf{\Gamma}_{2,i} (1 \leq i \leq l)$ such that $\mathbf{w}_{2,i} \in \mathcal{C}_\mathbf{w}, \mathbf{\Gamma}_{2,i} \in \mathcal{C}_\Gamma$ and $\|\mathbf{w}_{2,i} - \mathbf{w}_{1,i}\|_2 \leq \epsilon/2, \|\mathbf{\Gamma}_{2,i} - \mathbf{\Gamma}_{1,i}\|_F \leq \epsilon^2/(4\beta^2)$. Thus, we have

$$\text{dist}(V_1, V_2) \leq \beta \max_{1 \leq i \leq l} \sqrt{\|\mathbf{\Gamma}_{1,i} - \mathbf{\Gamma}_{2,i}\|_F} + \max_{1 \leq i \leq l} \|\mathbf{w}_{1,i} - \mathbf{w}_{2,i}\|_2 \leq \epsilon,$$

where the inequality holds due to (F.9). Therefore, the ϵ -covering number of optimistic function class \mathcal{V}_h is bounded by $\mathcal{N}_\epsilon \leq |\mathcal{C}_\mathbf{w}|^l \cdot |\mathcal{C}_\Gamma|^l$ and it implies

$$\log \mathcal{N}_\epsilon \leq dl \log(1 + 4L/\epsilon) + d^2 l \log(1 + 8\sqrt{d}\beta^2/\epsilon^2),$$

where the first inequality holds due to (F.10). Thus, we finish the proof of Lemma F.6. \square

With a similar argument as Lemma F.6, the covering number \mathcal{N}_ϵ of pessimistic value function class $\check{\mathcal{V}}_h$ can be upper bounded by the following lemma:

Lemma F.7. *For pessimistic function class $\check{\mathcal{V}}_h$, we define the distance between two function V_1 and V_2 as $V_1, V_2 \in \check{\mathcal{V}}_h$ as $\text{dist}(V_1, V_2) = \max_s |V_1(s) - V_2(s)|$. With respect to this distance function, the ϵ -covering number \mathcal{N}_ϵ of the function class $\check{\mathcal{V}}_h$ can be upper bounded by*

$$\log \mathcal{N}_\epsilon \leq dl \log(1 + 4L/\epsilon) + d^2 l \log(1 + 8\sqrt{d}\beta^2/\epsilon^2)$$

Proof of Lemma F.7. For any two function $V_1, V_2 \in \check{\mathcal{V}}_h$, according to the definition of function class $\check{\mathcal{V}}_h$, we have

$$\begin{aligned} V_1(\cdot) &= \max_a \max_{1 \leq i \leq l} \max \left(0, r_h(\cdot, a) + \mathbf{w}_{1,i}^\top \phi(\cdot, a) - \beta \sqrt{\phi(\cdot, a)^\top \mathbf{\Gamma}_{1,i} \phi(\cdot, a)} \right), \\ V_2(\cdot) &= \max_a \max_{1 \leq i \leq l} \max \left(0, r_h(\cdot, a) + \mathbf{w}_{2,i}^\top \phi(\cdot, a) - \beta \sqrt{\phi(\cdot, a)^\top \mathbf{\Gamma}_{2,i} \phi(\cdot, a)} \right), \end{aligned}$$

where $\|\mathbf{w}_{1,i}\|_2, \|\mathbf{w}_{2,i}\|_2 \leq L$ and $\mathbf{\Gamma}_{1,i}, \mathbf{\Gamma}_{2,i} \preceq \mathbf{I}$. Since all of the functions $\max_{1 \leq i \leq l}, \max_a$ and $\max(0, \cdot)$ are contraction functions, we have

$$\begin{aligned} \text{dist}(V_1, V_2) &= \max_{s \in \mathcal{S}} |V_1(s) - V_2(s)| \\ &\leq \max_{1 \leq i \leq l, s \in \mathcal{S}, a \in \mathcal{A}} \left| \mathbf{w}_{1,i}^\top \phi(s, a) - \beta \sqrt{\phi(s, a)^\top \mathbf{\Gamma}_{1,i} \phi(s, a)} \right. \\ &\quad \left. - \mathbf{w}_{2,i}^\top \phi(s, a) + \beta \sqrt{\phi(s, a)^\top \mathbf{\Gamma}_{2,i} \phi(s, a)} \right| \\ &\leq \beta \max_{1 \leq i \leq l, s \in \mathcal{S}, a \in \mathcal{A}} \left| \sqrt{\phi(s, a)^\top \mathbf{\Gamma}_{1,i} \phi(s, a)} - \sqrt{\phi(s, a)^\top \mathbf{\Gamma}_{2,i} \phi(s, a)} \right| \\ &\quad + \max_{1 \leq i \leq l, s \in \mathcal{S}, a \in \mathcal{A}} |(\mathbf{w}_{1,i} - \mathbf{w}_{2,i})^\top \phi(s, a)| \\ &\leq \beta \max_{1 \leq i \leq l, s \in \mathcal{S}, a \in \mathcal{A}} \left| \sqrt{\phi(s, a)^\top (\mathbf{\Gamma}_{1,i} - \mathbf{\Gamma}_{2,i}) \phi(s, a)} \right| \\ &\quad + \max_{1 \leq i \leq l, s \in \mathcal{S}, a \in \mathcal{A}} |(\mathbf{w}_{1,i} - \mathbf{w}_{2,i})^\top \phi(s, a)| \\ &\leq \beta \max_{1 \leq i \leq l} \sqrt{\|\mathbf{\Gamma}_{1,i} - \mathbf{\Gamma}_{2,i}\|_F} + \max_{1 \leq i \leq l} \|\mathbf{w}_{1,i} - \mathbf{w}_{2,i}\|_2, \end{aligned} \tag{F.11}$$

where the first inequality holds due to the contraction property, the second inequality holds due to the fact that $\max_x |f(x) + g(x)| \leq \max_x |f(x)| + \max_x |g(x)|$, the third inequality holds due to $|\sqrt{x} - \sqrt{y}| \geq |\sqrt{x} - \sqrt{y}|$ and the last inequality holds due to the fact that $\|\phi(s, a)\|_2 \leq 1$. Now, we denote \mathcal{C}_w as a $\epsilon/2$ -cover of the set $\{\mathbf{w} \in \mathbb{R}^d \mid \|\mathbf{w}\|_2 \leq L\}$ and \mathcal{C}_Γ as a $\epsilon^2/(4\beta^2)$ -cover of the set $\{\mathbf{\Gamma} \in \mathbb{R}^{d \times d} \mid \|\mathbf{\Gamma}\|_F \leq \sqrt{d}\}$ with respect to the Frobenius norms. Thus, according to Lemma F.5, we have following property:

$$|\mathcal{C}_w| \leq (1 + 4L/\epsilon)^d, |\mathcal{C}_\Gamma| \leq (1 + 8\sqrt{d}\beta^2/\epsilon^2)^{d^2}. \tag{F.12}$$

By the definition of covering number, for any function $V_1 \in \check{\mathcal{V}}$ with parameters $\mathbf{w}_{1,i}, \mathbf{\Gamma}_{1,i} (1 \leq i \leq l)$, s other parameters $\mathbf{w}_{2,i}, \mathbf{\Gamma}_{2,i} (1 \leq i \leq l)$ such that $\mathbf{w}_{2,i} \in \mathcal{C}_w, \mathbf{\Gamma}_{2,i} \in \mathcal{C}_\Gamma$ and $\|\mathbf{w}_{2,i} - \mathbf{w}_{1,i}\|_2 \leq \epsilon/2, \|\mathbf{\Gamma}_{2,i} - \mathbf{\Gamma}_{1,i}\|_F \leq \epsilon^2/(4\beta^2)$. Thus, we have

$$\text{dist}(V_1, V_2) \leq \beta \max_{1 \leq i \leq l} \sqrt{\|\mathbf{\Gamma}_{1,i} - \mathbf{\Gamma}_{2,i}\|_F} + \max_{1 \leq i \leq l} \|\mathbf{w}_{1,i} - \mathbf{w}_{2,i}\|_2 \leq \epsilon,$$

where the inequality holds due to (F.11). Therefore, the ϵ -covering number of function class $\check{\mathcal{V}}_h$ is bounded by $\mathcal{N}_\epsilon \leq |\mathcal{C}_w|^l \cdot |\mathcal{C}_\Gamma|^l$ and it implies

$$\log \mathcal{N}_\epsilon \leq dl \log(1 + 4L/\epsilon) + d^2 l \log(1 + 8\sqrt{d}\beta^2/\epsilon^2),$$

where the first inequality holds due to (F.12). Thus, we finish the proof of Lemma F.7. \square

In addition, according to the result in Lemma F.6, the covering number \mathcal{N}_ϵ of squared function class \mathcal{V}_h^2 can be upper bounded by:

Lemma F.8. *For squared function class \mathcal{V}_h^2 , we define the distance between two function V_1^2 and V_2^2 as $V_1^2, V_2^2 \in \mathcal{V}_h^2$ as $\text{dist}(V_1^2, V_2^2) = \max_s |V_1^2(s) - V_2^2(s)|$. With respect to this distance function, the ϵ -covering number \mathcal{N}_ϵ of the function class \mathcal{V}_h^2 can be upper bounded by*

$$\log \mathcal{N}_\epsilon \leq dl \log(1 + 8HL/\epsilon) + d^2 l \log(1 + 32\sqrt{d}H^2\beta^2/\epsilon^2)$$

Proof of Lemma F.8. For any function $V_1^2, V_2^2 \in \mathcal{V}_h^2$, the distance can be upper bounded by:

$$\begin{aligned} \text{dist}(V_1^2, V_2^2) &= \max_{s \in \mathcal{S}} |V_1^2(s) - V_2^2(s)| \\ &= \max_{s \in \mathcal{S}} |V_1(s) - V_2(s)| \cdot |V_1(s) + V_2(s)| \\ &\leq 2H \max_{s \in \mathcal{S}} |V_1(s) - V_2(s)| \\ &= 2H \text{dist}(V_1, V_2), \end{aligned} \tag{F.13}$$

where the inequality holds due to the fact that $0 \leq V_1(s), V_2(s) \leq H$. Thus, any $(\epsilon/2H)$ -net for optimistic function class \mathcal{V}_h is also a $(\epsilon/2H)$ -net for the squared function class \mathcal{V}^2 . According to Lemma F.6, the covering number of the squared function class is upper bounded by:

$$\log \mathcal{N}_\epsilon \leq dl \log(1 + 4L/\epsilon) + d^2 l \log(1 + 8\sqrt{d}\beta^2/\epsilon^2).$$

Thus, we finish the proof of Lemma F.8. \square

G Auxiliary Lemmas

Lemma G.1. *For any stage $h \in [h]$ in a linear MDP and any bounded-function $V : \mathcal{S} \rightarrow [0, B]$, there always exists a vector $\mathbf{w} \in \mathbb{R}^d$ such that for all state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$, we have*

$$[\mathbb{P}_h V](s, a) = \mathbf{w}^\top \boldsymbol{\phi}(s, a), \text{ where } \|\mathbf{w}\|_2 \leq B\sqrt{d}.$$

Proof of Lemma G.1. According to the definition of linear MDP (Assumption 3.2), we have

$$\begin{aligned} [\mathbb{P}_h V](s, a) &= \int \mathbb{P}_h(s'|s, a) V(s') ds' \\ &= \int \boldsymbol{\phi}(s, a)^\top V(s') d\boldsymbol{\theta}_h(s') \\ &= \boldsymbol{\phi}(s, a)^\top \int V(s') d\boldsymbol{\theta}_h(s') \\ &= \boldsymbol{\phi}(s, a)^\top \mathbf{w}, \end{aligned}$$

where we set $\mathbf{w} = \int V(s') d\boldsymbol{\theta}_h(s')$. In addition, the norm of \mathbf{w} is upper bounded by:

$$\left\| \int V(s') d\boldsymbol{\theta}_h(s') \right\| \leq \max_{s'} V(s') \cdot \sqrt{d} = B\sqrt{d}.$$

Thus, we finish the proof of Lemma G.1. \square

Lemma G.2 (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*

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

Lemma G.3 (Lemma 11, Abbasi-Yadkori et al. 2011). Let $\{\mathbf{x}_k\}_{k=1}^K$ be a sequence of vectors in \mathbb{R}^d , matrix Σ_0 a $d \times d$ positive definite matrix and define $\Sigma_k = \Sigma_0 + \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^\top$, then we have

$$\sum_{i=1}^k \min \left\{ 1, \mathbf{x}_i^\top \Sigma_{i-1}^{-1} \mathbf{x}_i \right\} \leq 2 \log \left(\frac{\det \Sigma_k}{\det \Sigma_0} \right).$$

In addition, if $\|\mathbf{x}_i\|_2 \leq L$ holds for all $i \in [K]$, then

$$\sum_{i=1}^k \min \left\{ 1, \mathbf{x}_i^\top \Sigma_{i-1}^{-1} \mathbf{x}_i \right\} \leq 2 \log \left(\frac{\det \Sigma_k}{\det \Sigma_0} \right) \leq 2 \left(d \log \left((\text{trace}(\Sigma_0) + kL^2)/d \right) - \log \det \Sigma_0 \right).$$

Lemma G.4 (Lemma 12, Abbasi-Yadkori et al. 2011). Suppose $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times d}$ are two positive definite matrices satisfying that $\mathbf{A} \succeq \mathbf{B}$, then for any $\mathbf{x} \in \mathbb{R}^d$, $\|\mathbf{x}\|_{\mathbf{A}} \leq \|\mathbf{x}\|_{\mathbf{B}} \cdot \sqrt{\det(\mathbf{A})/\det(\mathbf{B})}$.

Lemma G.5 (Confidence Ellipsoid, Theorem 2, Abbasi-Yadkori et al. 2011). Let $\{\mathcal{G}_k\}_{k=1}^\infty$ be a filtration, and $\{\mathbf{x}_k, \eta_k\}_{k \geq 1}$ be a stochastic process such that $\mathbf{x}_k \in \mathbb{R}^d$ is \mathcal{G}_k -measurable and $\eta_k \in \mathbb{R}$ is \mathcal{G}_{k+1} -measurable. Let $L, \sigma, \Sigma, \epsilon > 0$, $\boldsymbol{\mu}^* \in \mathbb{R}^d$. For $k \geq 1$, let $y_k = \langle \boldsymbol{\mu}^*, \mathbf{x}_k \rangle + \eta_k$ and suppose that η_k, \mathbf{x}_k also satisfy

$$\mathbb{E}[\eta_k | \mathcal{G}_k] = 0, |\eta_k| \leq R, \|\mathbf{x}_k\|_2 \leq L. \quad (\text{G.1})$$

For $k \geq 1$, let $\mathbf{Z}_k = \lambda \mathbf{I} + \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^\top$, $\mathbf{b}_k = \sum_{i=1}^k y_i \mathbf{x}_i$, $\boldsymbol{\mu}_k = \mathbf{Z}_k^{-1} \mathbf{b}_k$, and

$$\beta_k = R \sqrt{d \log \left(\frac{1 + kL^2/\lambda}{\delta} \right)}.$$

Then, for any $0 < \delta < 1$, we have with probability at least $1 - \delta$ that,

$$\forall k \geq 1, \left\| \sum_{i=1}^k \mathbf{x}_i \eta_i \right\|_{\mathbf{Z}_k^{-1}} \leq \beta_k, \|\boldsymbol{\mu}_k - \boldsymbol{\mu}^*\|_{\mathbf{Z}_k} \leq \beta_k + \sqrt{\lambda} \|\boldsymbol{\mu}^*\|_2.$$

Lemma G.6 (Lemma 4.4, Zhou and Gu 2022). Let $\{\sigma_k, \hat{\beta}_k\}_{k \geq 1}$ be a sequence of non-negative numbers, $\alpha, \gamma > 0$, $\{\mathbf{a}_k\}_{k \geq 1} \subset \mathbb{R}^d$ and $\|\mathbf{a}_k\|_2 \leq A$. Let $\{\bar{\sigma}_k\}_{k \geq 1}$ and $\{\hat{\Sigma}_k\}_{k \geq 1}$ be (recursively) defined as follows: $\hat{\Sigma}_1 = \lambda \mathbf{I}$,

$$\forall k \geq 1, \bar{\sigma}_k = \max\{\sigma_k, \alpha, \gamma \|\mathbf{a}_k\|_{\hat{\Sigma}_k^{-1}}^{1/2}\}, \hat{\Sigma}_{k+1} = \hat{\Sigma}_k + \mathbf{a}_k \mathbf{a}_k^\top / \bar{\sigma}_k^2.$$

Let $\iota = \log(1 + KA^2/(d\lambda\alpha^2))$. Then we have

$$\sum_{k=1}^K \min \left\{ 1, \|\mathbf{a}_k\|_{\hat{\Sigma}_k^{-1}} \right\} \leq 2d\iota + 2\gamma^2 d\iota + 2\sqrt{d\iota} \sqrt{\sum_{k=1}^K (\sigma_k^2 + \alpha^2)}.$$