

# Doubly Regularized Markov Decision Process for Robust Reinforcement Learning

Yiting He<sup>\*†</sup>      Zhishuai Liu<sup>†‡</sup>      Pan Xu<sup>§</sup>

## Abstract

Empirical successes show that regularization improves the stability and efficiency of reinforcement learning (RL), with applications in robotics and post-training of large language models. Yet, theoretical analyses of regularized Markov decision processes (MDPs) have mostly been confined to the standard RL setting. In this work, we investigate regularized MDPs through the lens of robust RL. We introduce a doubly regularized MDP framework that combines policy and dynamics regularization to enable robust policy learning against reward and dynamics perturbations. Within this framework, we develop an optimism-based online algorithm and provide the first finite-sample regret guarantees in both tabular and rich-observation settings, where the state-action space may be continuous. Our results show that algorithms for doubly regularized MDPs are as sample-efficient as well-studied robust MDP algorithms, while additionally benefiting from the flexibility of soft policies. Finally, experiments on CartPole demonstrate that our approach efficiently and effectively handles function approximation and exploration in large state-action spaces, achieving robust performances.

## 1 INTRODUCTION

Regularizations have been widely studied in reinforcement learning (RL) literature. Practically, existing works leveraging policy regularization have demonstrated tremendous successes in real world applications such as game playing (Mnih et al., 2016), robotics control (Schulman et al., 2015, 2017) and post-training of large language models (Bai et al., 2022; Ouyang et al., 2022). Theoretically, extensive studies (Geist et al., 2019; Yang et al., 2019; Vieillard et al., 2020; Zhao et al., 2025) have investigated the policy regularized MDP in the standard RL setting. Moreover, recent works (Derman et al., 2021; Eysenbach and Levine, 2021; Husain et al., 2021) also point out its relation with reward-robust MDPs. Another line of works (Yang et al., 2023; Zhang et al., 2024; He et al., 2025) propose the regularized robust MDP (RRMDP) framework, where a dynamics regularization is added to the objective function to account for robustness against potential dynamics shift. They have shown that the dynamics regularization leads to more efficient estimation procedures of value functions compared with the conventional distributionally robust MDP framework (Iyengar, 2005; Nilim and El Ghaoui, 2005).

To combine the strength of regularization, in this work, we unify the policy and dynamics regularization in a doubly regularized MDP framework, and investigate it through the lens of robust

---

<sup>\*</sup>University of Science and Technology of China; email: heyiting@mail.ustc.edu.cn

<sup>†</sup>Duke University; email: zhishuai.liu@duke.edu

<sup>‡</sup>Equal contribution

<sup>§</sup>Duke University; email: pan.xu@duke.edu

reinforcement learning. In particular, we first establish the framework in the tabular setting where the state and action spaces are finite. Unlike prior works (Yang et al., 2023; Zhang et al., 2024; Panaganti et al., 2024) assuming access to a fixed dataset or a simulator, we focus on the more realistic online setting, where the agent can only gather data through interaction with the environment. In this setting, the agent faces the unique challenge of exploration, that is, efficiently acquiring information about the environment through strategically selecting the policy. Existing results under this online setting often rely on restrictive structural assumptions, such as the fail-state condition in Lu et al. (2024); Liu and Xu (2024a); Liu et al. (2024), which applies only to uncertainty sets defined by Total Variation distance and does not extend to more general  $f$ -divergences. Recently, He et al. (2025) identified the information deficit issue as the central obstacle in robust RL, arising from the need to generalize experience under distribution shifts to perturbed environments while only interacting with the source environment. Building on this, we adopt the bounded visitation measure ratio assumption proposed by He et al. (2025), extending the analysis to handle the double regularization of our framework. Further, to account for large state and action spaces in practical scenarios, we extend our doubly regularized MDP framework to incorporate function approximation for generalization. We summarize our main contribution as follows:

- We establish a novel doubly regularized MDP framework that incorporates both transition regularization and policy regularization, and show that this framework achieves robustness with respect to both transitions and rewards.
- Under this framework, we propose an online learning algorithm Robust Soft Policy Value Iteration (RSPVI) with general  $f$ -divergence based transition regularization terms and several policy regularization formulations. We establish the first finite sample regret bounds for doubly regularized MDPs.
- We further extend RSPVI to the setting with linear function approximation. With fewer structural assumptions than prior work, we provide the first finite-sample regret guarantees applicable to broader settings defined by general  $f$ -divergences.
- In stark contrast to the conventional robust MDP framework, the optimal policies under doubly regularized MDPs are stochastic. We empirically evaluate RSPVI in the CartPole environment, showing that the stochastic policies learned by RSPVI facilitate exploration and posses robustness performances against environment perturbation.

**Notations** For  $H \in \mathbb{Z}_+$ , we denote  $[H] = \{1, \dots, H\}$ . For any set  $\mathcal{S}$ , define  $\Delta(\mathcal{S})$  as the set of probability distributions over  $\mathcal{S}$ . For a convex function  $f : [0, +\infty) \rightarrow (-\infty, +\infty]$  with  $f(x)$  finite for  $x > 0$ ,  $f(1) = 0$  and  $f(0) = \lim_{t \rightarrow 0^+} f(t)$ . For  $P, Q \in \Delta(\mathcal{S})$  and  $P \ll Q$ , the  $f$ -divergence of  $P$  from  $Q$  is defined as  $D_f(P \| Q) = \int_{\Omega} f(\frac{dP}{dQ}) dQ$ . In this paper, we focus on total variation (TV) distance with  $f(t) = \frac{1}{2}|t - 1|$ , Kullback-Leibler (KL) divergence with  $f(t) = t \ln t$ , and  $\chi^2$ -divergence with  $f(t) = (t - 1)^2$ . We use  $\mathcal{O}(\cdot)$  to hide absolute constant factors and  $\tilde{\mathcal{O}}(\cdot)$  to further hide logarithmic factors. We denote  $\text{proj}_C(x) = \arg \inf_{c \in C} \|c - x\|^2$  as the projection of  $x$  onto the set  $C$ . We denote the support function of a set  $\mathcal{Z} \subseteq \mathbb{R}^n$  as  $\sigma_{\mathcal{Z}}(\mathbf{y}) = \max_{\mathbf{a} \in \mathcal{Z}} \langle \mathbf{a}, \mathbf{y} \rangle$ .

## 2 RELATED WORK

**Robust MDPs** Robust MDPs were first introduced as a control problem with a known nominal transition model (Iyengar, 2005; Nilim and El Ghaoui, 2005; Xu and Mannor, 2006; Wiesemann et al., 2013; Mannor et al., 2016), where robust policies are derived by solving a constrained max-min optimization problem. Subsequent work extended robust MDPs to the learning setting, both with

access to a generative model allowing the agent to sample transitions for arbitrary state-action pair (Zhou et al., 2021; Yang et al., 2022; Panaganti and Kalathil, 2022; Shi et al., 2024) and with pre-collected datasets that provide adequate coverage of the optimal policy under the perturbed environment (Shi and Chi, 2024; Panaganti et al., 2022; Blanchet et al., 2023; Wang et al., 2024a; Liu and Xu, 2024b, 2025). More recently, Lu et al. (2024); Liu and Xu (2024a) studied robust MDPs in the online setting under the fail-state or vanishing minimal-value structure assumption. He et al. (2025) identified the key challenge of online robust learning as the information deficit issue, and addressed this by introducing a maximum visitation ratio assumption, extending the analysis to the general  $f$ -divergence setting. Ghosh et al. (2025) employed a variance-aware proof technique to refine the results of He et al. (2025).

**Regularized Robust MDPs** Unlike the above approaches formulating robust MDPs as a constrained optimization problem, regularized robust MDPs were first introduced by Yang et al. (2023) and Zhang et al. (2024) as an efficient alternative. Yang et al. (2023) focused on the generative model setting, providing a sample complexity analysis and proving the statistical equivalence of the two formulations. Zhang et al. (2024) established the equivalence with risk-sensitive MDPs, analyzed the convergence rate of policy gradient methods, and proposed a value iteration algorithm with general function approximation for the offline setting. Subsequently, Panaganti et al. (2024) extended the framework to general  $f$ -divergences with function approximation in the offline setting, and to the Total Variation distance under a fail-state assumption in the hybrid offline-online setting. More recently, Tang et al. (2025) studied linear function approximation for general  $f$ -divergences using offline datasets.

**Reward Robust MDPs** Most existing works focus solely on transition robustness. For studies addressing reward robustness, notable examples include Zhou et al. (2021), Wang et al. (2023), and Wang et al. (2024b). All of these studies formulate reward robustness through an uncertainty set. Specifically, Zhou et al. (2021) proposed a model-based value iteration algorithm using a pre-collected dataset, while Wang et al. (2023) introduced a model-free  $Q$ -learning algorithm under the assumption of access to a simulator. Building on this, Wang et al. (2024b) incorporated a variance-reduction technique into  $Q$ -learning, also assuming access to a generative model. However, none of these works address the more practical online setting. The work most closely related to ours is Derman et al. (2021), which demonstrates the equivalence between policy regularization and reward robustness. However, it relies on the assumption that the nominal transition is known, which is a fairly strong requirement.

## 3 PRELIMINARIES

### 3.1 Problem Formulations

**Markov Decision Process (MDP)** A finite-horizon Markov Decision Process is denoted as the tuple  $(\mathcal{S}, \mathcal{A}, P, r, H)$ . Here,  $\mathcal{S}$  and  $\mathcal{A}$  represent the state and action spaces. We also denote  $|\mathcal{A}|$  as the volume of  $\mathcal{A}$ , where  $|\mathcal{A}|$  is the number of actions for discrete action space and  $|\mathcal{A}| = \int_{\mathcal{A}} 1 da$  for continuous space. The transition dynamics are given by  $P = \{P_h\}_{h=1}^H$ , where each  $P_h(\cdot|s, a) : \mathcal{S} \rightarrow \Delta(\mathcal{S})$  specifies the probability distribution over the next state when action  $a$  is taken at state  $s$  and step  $h$ . The reward functions are  $r = \{r_h\}_{h=1}^H$ , with each  $r_h : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  assumed to be

deterministic. We define the value function and  $Q$ -function as

$$V_h^\pi(s) = \mathbb{E}_{\pi, P} \left[ \sum_{t=h}^H r_t(s_t, a_t) \mid s_h = s \right],$$

$$Q_h^\pi(s, a) = \mathbb{E}_{\pi, P} \left[ \sum_{t=h}^H r_t(s_t, a_t) \mid s_h = s, a_h = a \right].$$

**Doubly Regularized MDPs** We now introduce the doubly regularized framework. A finite-horizon doubly regularized MDP can be denoted as the tuple  $(\mathcal{S}, \mathcal{A}, P^o, r^o, \beta, D, \Omega^\eta, H)$ . Here,  $D$  is a probability divergence that quantifies the distribution shift (later we will instantiate  $D$  as TV, KL or  $\chi^2$ -divergence), and  $\beta$  is the corresponding regularization parameter.  $\Omega_s^\eta(\pi) \triangleq \Omega^\eta(s, \pi(\cdot|s))$  denotes a penalty function on the chosen policy  $\pi$  at state  $s$ , with the  $\eta$  in  $\Omega^\eta$  controlling the extent of this penalty. We define the doubly regularized value function and  $Q$ -function as

$$V_h^{\pi, \beta, \eta}(s) = \inf_{P_t \in \Delta(\mathcal{S})} \mathbb{E}_{\pi, \{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) + \beta \cdot D(P_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) - \Omega_{s_t}^\eta(\pi_t) \right) \mid s_h = s \right],$$

$$Q_h^{\pi, \beta, \eta}(s, a) = \inf_{P_t \in \Delta(\mathcal{S})} \mathbb{E}_{\pi, \{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) + \beta \cdot D(P_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) - \sum_{t=h+1}^H \Omega_{s_t}^\eta(\pi_t) \mid s_h = s, a_h = a \right].$$

We then define the optimal value function  $V_h^{*, \beta, \eta}(s) = \sup_{\pi \in \Pi} V_h^{\pi, \beta, \eta}(s)$  and the optimal  $Q$ -function  $Q_h^{*, \beta, \eta}(s, a) = \sup_{\pi \in \Pi} Q_h^{\pi, \beta, \eta}(s, a)$  for all  $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$ , where  $\Pi$  is the set of all possible policies. Correspondingly, the optimal policy  $\pi^* = \{\pi_h^*\}_{h=1}^H$  is defined as the policy that achieves the optimal value function for all  $(h, s) \in [H] \times \mathcal{S}$ , that is,  $\pi_h^* = \arg \sup_{\pi \in \Pi} V_h^{\pi, \beta, \eta}(s)$ .

**Dynamic Programming Principles** Under the doubly regularized framework, we establish the dynamic programming principle as follows

**Proposition 3.1.** For doubly regularized tabular MDPs, it holds that for any policy  $\pi$  and any  $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$ ,

$$Q_h^{\pi, \beta, \eta}(s, a) = r_h(s, a) + \inf_{P_h \in \Delta(\mathcal{S})} \left\{ \mathbb{E}_{s' \sim P_h(\cdot|s, a)} [V_{h+1}^{\pi, \beta, \eta}(s')] + \beta \cdot D(P_h(\cdot|s, a) \| P_h^o(\cdot|s, a)) \right\}, \quad (3.1)$$

$$V_h^{\pi, \beta, \eta}(s) = -\Omega_s^\eta(\pi_h) + \mathbb{E}_{a \sim \pi_h(\cdot|s)} [Q_h^{\pi, \beta, \eta}(s, a)]. \quad (3.2)$$

**Double Robustness of the Proposed Framework** The effectiveness of incorporating regularization on transition deviations and accounting for the worst case in transition robustness has been extensively studied in the literature (Yang et al., 2023; Zhang et al., 2024; Panaganti et al., 2024; He et al., 2025; Tang et al., 2025). **Proposition 3.2** further demonstrates that introducing a policy

regularization term is equivalent to considering the worst-case reward within a suitably defined uncertainty set (Derman et al., 2021). Therefore, our proposed doubly regularized framework can achieve double robustness.

**Proposition 3.2.** When setting the policy regularization term as

$$\Omega_s^\eta(\pi) \triangleq \Omega^\eta(s, \pi(\cdot|s)) = \sigma_{\mathcal{R}_s^\eta(\pi)}(-\pi(\cdot|s))$$

for all  $(\pi, s) \in \Pi \times \mathcal{S}$ , where  $\mathcal{R}^\eta : \mathcal{S} \times \mathcal{A} \times \Pi \rightarrow I \subseteq \mathbb{R}$  denotes the uncertainty set within which the reward deviation  $\Delta r = \tilde{r} - r$  may lie, and  $\sigma_{\mathcal{R}_s^\eta(\pi)}(-\pi(\cdot|s)) = \sup_{\Delta r \in \mathcal{R}^\eta(\pi)} \langle -\pi(\cdot|s), \Delta r(s, \cdot) \rangle$ . Then for any state  $s \in \mathcal{S}$  and any function  $V : \mathcal{S} \rightarrow \mathbb{R}$ , we have

$$\begin{aligned} & \inf_{\substack{\tilde{r} \in r + \mathcal{R}^\eta(\pi) \\ P \in \Delta(\mathcal{S})}} \mathbb{E}_{a \sim \pi(\cdot|s)} [\tilde{r}(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} [V(s')] + \beta \cdot D(P(s, a) \| P^o(s, a))] \\ &= \inf_{P \in \Delta(\mathcal{S})} \mathbb{E}_{a \sim \pi(\cdot|s)} [r(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} [V(s')] + \beta \cdot D(P(s, a) \| P^o(s, a)) - \Omega_s^\eta(\pi)]. \end{aligned}$$

**Learning Goal** We consider the online setting, where an agent aims to learn the optimal robust policy by interacting with the environment over  $K$  episodes. At the start of each episode  $k$ , the agent is given an initial state  $s_1^k$ , which we assume to be fixed without loss of generality. For the policy  $\pi^k$  selected by the agent based on the history at episode  $k$ , we measure its sub-optimality by  $V_1^{*, \beta, \eta}(s_1) - V_1^{\pi^k, \beta, \eta}(s_1)$ . The cumulative sub-optimality after  $K$  episodes, which is commonly referred to as  $\text{Regret}(K)$ , is defined as

$$\text{Regret}(K) = \sum_{k=1}^K (V_1^{*, \beta, \eta}(s_1) - V_1^{\pi^k, \beta, \eta}(s_1)).$$

## 4 ALGORITHM

First, we present a meta-algorithm [Algorithm 1](#) for online doubly regularized MDPs, incorporating general  $f$ -divergence based regularization terms and various types of policy regularization terms. We then instantiate [Algorithm 1](#) in each case.

Similar to [Lu et al. \(2024\)](#); [Liu and Xu \(2024a\)](#); [He et al. \(2025\)](#), our algorithm Robust Soft Policy Value Iteration (RSPVI) adopts a value iteration framework. We leverage the robust Bellman optimality equation and incorporate the optimism principle ([Abbasi-Yadkori et al., 2011](#)) to estimate the robust  $Q$ -functions. For the explicit computation expression, we employ the dual formulation corresponding to each  $f$ -divergence setting.

### 4.1 Model Estimation

We use a model based manner to estimate the empirical reward and transition in the tabular setting. In each episode  $k$ , after executing policy  $\pi^k$  and collecting a trajectory  $\tau^k = (s_1^k, a_1^k, r_1^k, \dots, s_H^k, a_H^k, r_H^k)$ , RSPVI updates the empirical estimations as follows

$$n_h^k(s, a) = \sum_{i=1}^k \mathbb{1} \{s_h^i = s, a_h^i = a\},$$

---

**Algorithm 1** Robust Soft Policy Value Iteration (RSPVI) for tabular setting
 

---

**Require:** transition regularizer  $\beta > 0$ , policy regularizer  $\Omega^\eta$ .

```

1: for  $k = 1, \dots, K$  do
2:    $V_{H+1}^{k,\beta,\eta}(\cdot) \leftarrow 0$ .
3:   for  $h = H, \dots, 1$  do
4:     for  $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$  do
5:       Update  $Q$ -function estimation  $Q_h^{k,\beta,\eta}(s, a)$  according to (4.2).
6:     end for
7:     for  $\forall s \in \mathcal{S}$  do
8:       Update policy  $\pi_h^k(\cdot|s)$  according to (4.3).
9:       Update value function estimation  $V_h^{k,\beta,\eta}(s) = \langle Q_h^{k,\beta,\eta}(s, \cdot), \pi_h^k(\cdot|s) \rangle - \Omega_s^\eta(\pi_h^k)$ .
10:    end for
11:  end for
12:  Collect a trajectory  $\tau^k$  by executing  $\pi^k$ .
13:  Update model estimation as (4.1).
14: end for
  
```

---

$$\begin{aligned}
 \hat{r}_h^{k+1}(s, a) &= \frac{\sum_{i=1}^k r_h^i(s, a) \cdot \mathbb{1}\{s_h^i = s, a_h^i = a\}}{n_h^k(s, a) \vee 1}, \\
 \hat{P}_h^{k+1}(s'|s, a) &= \frac{\sum_{i=1}^k \mathbb{1}\{s_h^i = s, a_h^i = a, s_{h+1}^i = s'\}}{n_h^k(s, a) \vee 1}.
 \end{aligned} \tag{4.1}$$

## 4.2 $Q$ -function Estimation

The estimation of optimal robust  $Q$ -function  $Q_h^{k,\beta,\eta}$ ,  $(h, k) \in [H] \times [K]$  is constructed as follows

$$Q_h^{k,\beta,\eta}(s, a) = \min \{ \text{RB}_h^k(s, a) + b_h^k(s, a), H - h + 1 \}, \tag{4.2}$$

which consists of a robust Bellman estimator  $\text{RB}_h^k(s, a)$  and a bonus term  $b_h^k(s, a)$ . In the following, we instantiate this generic form under different  $f$ -divergence settings, providing explicit formulations for robust Bellman estimation as well as bonus design.

**Tabular-TV Setting** According to the dual formulation for regularized TV setting in [Lemma C.3](#), given the estimated value function  $V_{h+1}^{k,\beta,\eta}$  and empirical reward  $\hat{r}_h^k$  and transition  $\hat{P}_h^k$ , we choose the robust Bellman estimator  $\text{RB}_h^k(s, a)$  and the bonus  $b_h^k(s, a)$  as

$$\begin{aligned}
 \text{RB}_h^k(s, a) &= \hat{r}_h^k(s, a) - \mathbb{E}_{\hat{P}_h^k} \left[ \left( \min_{s \in \mathcal{S}} V_{h+1}^{k,\beta,\eta}(s) + \beta - V_{h+1}^{k,\beta,\eta}(s) \right)_+ \right](s, a) + \left( \min_{s \in \mathcal{S}} V_{h+1}^{k,\beta,\eta}(s) + \beta \right), \\
 b_h^k(s, a) &= 2H \sqrt{\frac{2S \ln(2SAHK/\delta)}{n_h^{k-1}(s, a) \vee 1}}.
 \end{aligned}$$

**Tabular-KL Setting** According to the dual formulation for regularized KL setting in [Lemma C.6](#), given the estimated value function  $V_{h+1}^{k,\beta,\eta}$  and empirical reward  $\hat{r}_h^k$  and transition  $\hat{P}_h^k$ , we choose the robust Bellman estimator  $\text{RB}_h^k(s, a)$  and the bonus  $b_h^k(s, a)$  as

$$\begin{aligned}\text{RB}_h^k(s, a) &= \hat{r}_h^k(s, a) - \beta \ln \mathbb{E}_{\hat{P}_h^k} [\exp(-\beta^{-1} V_{h+1}^{k,\beta,\eta})](s, a), \\ b_h^k(s, a) &= (1 + \beta e^{\beta^{-1} H} \sqrt{S}) \sqrt{\frac{2 \ln(2SAHK/\delta)}{n_h^{k-1}(s, a) \vee 1}}.\end{aligned}$$

**Tabular- $\chi^2$  Setting** According to the dual formulation for regularized  $\chi^2$  setting in [Lemma C.9](#), given the estimated value function  $V_{h+1}^{k,\beta,\eta}$  and empirical reward  $\hat{r}_h^k$  and transition  $\hat{P}_h^k$ , we choose the robust Bellman estimator  $\text{RB}_h^k(s, a)$  and the bonus  $b_h^k(s, a)$  as

$$\begin{aligned}\text{RB}_h^k(s, a) &= \hat{r}_h^k(s, a) + \sup_{\lambda \in [0, H]} (\mathbb{E}_{\hat{P}_h^k} [V_{h+1}^{k,\beta,\eta} - \lambda](s, a) - \frac{1}{4\beta} \text{Var}_{\hat{P}_h^k} [V_{h+1}^{k,\beta,\eta} - \lambda](s, a)), \\ b_h^k(s, a) &= \left(2 + \frac{3H}{4\beta}\right) H \sqrt{\frac{2S^2 \ln(48SAH^3K^2/\delta)}{n_h^{k-1}(s, a) \vee 1}} + \frac{1 + 4\beta}{4\beta K}.\end{aligned}$$

### 4.3 Policy Update

For the selection of policy  $\pi_h^k$ ,  $(h, k) \in [H] \times [K]$ , we first present a generic form [\(4.3\)](#), which maximizes the estimated value function  $V_h^k(s)$ .

$$\begin{aligned}\pi_h^k(\cdot|s) &= \arg\max_{\pi \in \Pi} \{ \mathbb{E}_{\pi(\cdot|s)} [Q_h^{k,\beta,\eta}(s, \cdot)] - \Omega_s^\eta(\pi) \} \\ &= \arg\max_{\pi \in \Pi} \{ \mathbb{E}_{\pi(\cdot|s)} [Q_h^{k,\beta,\eta}(s, \cdot)] - \sigma_{\mathcal{R}_s^\eta(\pi)}(-\pi(\cdot|s)) \},\end{aligned}\tag{4.3}$$

where we specify the choice of policy regularization according to [Proposition 3.2](#) in the second equality.

Next, for different types of the policy regularization, we specify the corresponding policy update formulation. The exact form of the reward-uncertainty interval for the associated support function is established in [Derman et al. \(2021, Section 3\)](#).

**Negative Shannon Entropy** Let  $\mathcal{R}_{s,a}^{\text{NS},\eta}(\pi) = [\eta \ln(1/\pi(a|s)), +\infty)$  for all  $(\pi, s, a) \in \Pi \times \mathcal{S} \times \mathcal{A}$ , and the associated support function is

$$\sigma_{\mathcal{R}_s^{\text{NS},\eta}(\pi)}(-\pi(\cdot|s)) = \eta \cdot \sum_{a \in \mathcal{A}} \pi(a|s) \ln(\pi(a|s)).$$

For this setting, we update the policy as

$$\pi_h^k(\cdot|s) \propto \exp(Q_h^{k,\beta,\eta}(s, \cdot)/\eta).\tag{4.4}$$



**Kullback-Leibler Divergence** Fixing a reference policy  $\tilde{\pi}$ , let  $\mathcal{R}_{s,a}^{\text{KL},\eta}(\pi) = \eta \ln(\tilde{\pi}(a|s)) + \mathcal{R}_{s,a}^{\text{NS},\eta}(\pi)$  for all  $(\pi, s, a) \in \Pi \times \mathcal{S} \times \mathcal{A}$ , and the associated support function is

$$\sigma_{\mathcal{R}_s^{\text{KL},\eta}(\pi)}(-\pi(\cdot|s)) = \eta \cdot \sum_{a \in \mathcal{A}} \pi(a|s) \ln(\pi(a|s)/\tilde{\pi}(a|s)).$$

For this setting, we update the policy as

$$\pi_h^k(\cdot|s) \propto \tilde{\pi}_h^k(\cdot|s) \exp(Q_h^{k,\beta,\eta}(s, \cdot)/\eta). \quad (4.5)$$

**Negative Tsallis Entropy** Let  $\mathcal{R}_{s,a}^{\text{NT},\eta}(\pi) = \eta[(1 - \pi(a|s))/2, +\infty)$  for all  $(\pi, s, a) \in \Pi \times \mathcal{S} \times \mathcal{A}$ , and the associated support function is

$$\sigma_{\mathcal{R}_s^{\text{NT},\eta}(\pi)}(-\pi(\cdot|s)) = \eta \cdot (\|\pi(\cdot|s)\|^2 - 1)/2.$$

For this setting, we update the policy as

$$\pi_h^k(\cdot|s) = \text{proj}_{\Delta}(Q_h^{k,\beta,\eta}(s, \cdot)/\eta). \quad (4.6)$$

In [Lu et al. \(2024\)](#); [Liu and Xu \(2024a\)](#); [He et al. \(2025\)](#), the updated policy is chosen as the greedy policy with respect to the estimated  $Q$ -function  $Q_h^{k,\beta,\eta}$ . In contrast, since our formulation incorporates a policy regularization term, the corresponding robust Bellman equations [Proposition 6.2](#) differ from theirs, leading to a different policy update rule. Notably, the formulation in [\(4.3\)](#) allows us to employ a soft policy update.

## 5 THEORETICAL RESULTS

In this section, we present the regret bounds for our proposed algorithm RSPVI. We begin by introducing the assumptions required for online robust learning. Then we discuss our results under different  $f$ -divergence settings.

[He et al. \(2025\)](#) showed that online robust learning can be exponentially hard without appropriate assumptions. Based on this, they introduce the bounded visitation measure ratio assumption and justify its necessity by deriving the lower bounds. Following their approach, we begin by introducing several notations that will be useful throughout the discussion.

**Definition 5.1** (Worst-case transition). In the tabular setting, for any policy  $\pi$ , we define the worst-case transition corresponding to  $\pi$  as

$$P_h^{w,\pi}(\cdot|s, a) = \underset{P_h \in \Delta(\mathcal{S})}{\text{argmin}} \mathbb{E}_{P_h}[V_{h+1}^{\pi,\beta,\eta}(s, a) + \beta \cdot D(P_h(\cdot|s, a) \| P_h^o(\cdot|s, a))].$$

Based on [Definition 5.1](#), we also define the corresponding visitation measure.

**Definition 5.2** (Visitation measure). At timestep  $h \in [H]$ , we denote  $d_h^\pi(\cdot)$  as the visitation measure on  $\mathcal{S}$  induced by policy  $\pi$  under  $P^o$ , and  $q_h^\pi(\cdot)$  as the visitation measure on  $\mathcal{S}$  induced by policy  $\pi$  under  $P^{w,\pi}$ .

We now make the same bounded visitation measure ratio assumption as in [He et al. \(2025\)](#).



**Assumption 5.3** (Bounded visitation measure ratio). Under the definition of [Definition 5.2](#), we define  $C_{vr} := \sup_{\pi, h, s} \frac{q_h^\pi(s)}{d_h^\pi(s)}$  as the supremal ratio between the nominal visitation measure and the worst-case visitation measure. We assume that  $C_{vr}$  is polynomial in  $H$ ,  $S$  and  $A$ .

**Theorem 5.4** (Tabular RRMDP regret bounds). Assume that [Assumption 5.3](#) holds for each  $f$ -divergence-based transition regularization, and consider the policy regularizations in (4.4), (4.5), and (4.6). Then for any  $\delta \in (0, 1/3)$ , with probability at least  $1 - 3\delta$ , the regret of [Algorithm 1](#) satisfies

$$\text{Regret}(K) = \begin{cases} \tilde{\mathcal{O}}(C_{vr}S^{\frac{3}{2}}AH^2 + C_{vr}^{\frac{1}{2}}SA^{\frac{1}{2}}H^2\sqrt{K}) & (\text{TV}) \\ \tilde{\mathcal{O}}((1 + \beta e^{\beta^{-1}H}\sqrt{S})(C_{vr}SAH + C_{vr}^{\frac{1}{2}}S^{\frac{1}{2}}A^{\frac{1}{2}}H\sqrt{K})) & (\text{KL}) \\ \tilde{\mathcal{O}}(C_{vr}S^2AH^3 + C_{vr}^{\frac{1}{2}}S^{\frac{3}{2}}A^{\frac{1}{2}}H^3\sqrt{K}) & (\chi^2) \end{cases}$$

[Theorem 5.4](#) presents the first regret bound for online MDPs that simultaneously accounts for both transition and reward uncertainties, demonstrating that sample-efficient online robust learning is achievable under [Assumption 5.3](#). Unlike transition dynamics, we make no assumptions on reward uncertainty, since reward variations do not affect the difficulty of reaching specific states across the source and target environments. A key distinction of our setting is that the optimal policy is soft, whereas in [He et al. \(2025\)](#) it must be greedy. This advantage arises from policy regularization, and often leads to the learned policy more exploratory and reliable.

The orders in [Theorem 5.4](#) match those in [He et al. \(2025\)](#), showing that incorporating an additional mechanism for reward robustness does not increase the overall regret. This highlights the sample efficiency of RSPVI and the advantages of our proposed doubly regularized framework. The key insight is as follows: in the online setting, we add a bonus to the estimated function to ensure optimism, and the estimation error can be bounded by this bonus term. Thanks to our algorithmic design, the optimism property continues to hold for all estimated value functions and  $Q$ -functions when the bonus terms are chosen in the same way as before. Consequently, by applying [Assumption 5.3](#) in the same manner as [He et al. \(2025\)](#), the order of the regret remains unchanged.

## 6 LINEAR MDP EXTENSION

### 6.1 Preliminaries

**Linear Function Approximation**  $d$ -rectangular linear regularized robust MDP was first proposed in [Tang et al. \(2025\)](#), which changes the uncertainty set into a penalty term while persisting the linear structure of the transitions and reward functions. We make the following assumption the same as [Liu and Xu \(2024a\)](#); [Tang et al. \(2025\)](#).

**Assumption 6.1.** ([Jin et al., 2020](#)) Given a known state-action feature mapping  $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$  satisfying  $\sum_{i=1}^d \phi_i(s, a) = 1, \phi_i(s, a) \geq 0$ . Further more, we assume the reward function  $\{r_h\}_{h=1}^H$  and the nominal transition kernels  $\{P_h^o\}_{h=1}^H$  admit linear structures. Specifically, we have for all  $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$ ,

$$r_h(s, a) = \langle \phi(s, a), \theta_h \rangle, \quad P_h^o(\cdot | s, a) = \langle \phi(s, a), \mu_h^o(\cdot) \rangle,$$

where  $\{\theta_h\}_{h=1}^H$  are unknown vectors with bounded norm  $\|\theta_h\|_2 \leq \sqrt{d}$  and  $\{\mu_h^o\}_{h=1}^H$  are unknown probability measure vectors over  $\mathcal{S}$ , i.e.,  $\mu_h^o = (\mu_{h,1}^o, \mu_{h,2}^o, \dots, \mu_{h,d}^o), \mu_{h,i}^o \in \Delta(\mathcal{S}), \forall i \in [d]$ .

Based on such linear structure, the robust value function and  $Q$ -function are defined as

$$\begin{aligned}
V_h^{\pi, \beta, \eta}(s) &= \inf_{\substack{\boldsymbol{\mu}_t \in \Delta(S)^d \\ P_t = \langle \phi, \boldsymbol{\mu}_t \rangle}} \mathbb{E}_{\pi, \{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) \right. \right. \\
&\quad \left. \left. + \beta \cdot \langle \phi(s_t, a_t), \mathbf{D}(\boldsymbol{\mu}_t \| \boldsymbol{\mu}_t^o) \rangle - \Omega_{s_t}^\eta(\pi_t) \right) \mid s_h = s \right], \\
Q_h^{\pi, \beta, \eta}(s, a) &= \inf_{\substack{\boldsymbol{\mu}_t \in \Delta(S)^d \\ P_t = \langle \phi, \boldsymbol{\mu}_t \rangle}} \mathbb{E}_{\pi, \{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) \right. \right. \\
&\quad \left. \left. + \beta \cdot \langle \phi(s_t, a_t), \mathbf{D}(\boldsymbol{\mu}_t \| \boldsymbol{\mu}_t^o) \rangle \right) - \sum_{t=h+1}^H \Omega_{s_t}^\eta(\pi_t) \mid s_h = s, a_h = a \right],
\end{aligned}$$

where  $\mathbf{D}(\boldsymbol{\mu} \| \boldsymbol{\mu}^o) = (\mathbf{D}(\mu_1 \| \mu_1^o), \dots, \mathbf{D}(\mu_d \| \mu_d^o))^\top$ .

It is worth noting that there has been no prior framework addressing reward robustness in the linear setting, largely due to the constraints of its linear structure. In contrast, under our doubly regularized framework, reward robustness naturally emerges from the policy regularization, making such robustness much easier to achieve.

**Dynamic Programming Principle** Similar to the tabular setting, we have the dynamic programming principle for doubly regularized MDPs as follows

**Proposition 6.2.** For doubly regularized linear MDPs, it holds that for any policy  $\pi$  and any  $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$ ,

$$Q_h^{\pi, \beta, \eta}(s, a) = \inf_{\substack{\boldsymbol{\mu}_h \in \Delta(S)^d \\ P_h = \langle \phi, \boldsymbol{\mu}_h \rangle}} \left\{ \mathbb{E}_{s' \sim P_h(\cdot | s, a)} [V_{h+1}^{\pi, \beta, \eta}(s')] + \beta \cdot \langle \phi(s, a), \mathbf{D}(\boldsymbol{\mu}_h \| \boldsymbol{\mu}_h^o) \rangle \right\} + r_h(s, a), \quad (6.1)$$

$$V_h^{\pi, \beta, \eta}(s) = -\Omega_s^\eta(\pi_h) + \mathbb{E}_{a \sim \pi_h(\cdot | s)} [Q_h^{\pi, \beta, \eta}(s, a)]. \quad (6.2)$$

## 6.2 Algorithm

Then, we present the linear version of [Algorithm 1](#) as [Algorithm 2](#).

**$Q$ -function Estimation** For the estimation of robust  $Q$ -function  $Q_h^{k, \beta, \eta}$ ,  $(h, k) \in [H] \times [K]$ , we first present a generic form [\(6.3\)](#). Since the feature mapping is known, it is sufficient to estimate the weight vectors  $\boldsymbol{\theta}_h^k$  and  $\mathbf{w}_h^k$  to recover the robust  $Q$ -function. Here,  $\boldsymbol{\theta}_h^k$  corresponds to the reward component, and  $\mathbf{w}_h^k$  corresponds to the estimated value function  $V_{h+1}^{k, \beta, \eta}$ . Instead of a model-based method in the tabular setting, we use ridge regression to estimate  $\boldsymbol{\theta}_h^k$  and  $\mathbf{w}_h^k$ .  $\Gamma_h^k(s, a)$  denotes the reward term, where  $c$  needs to be determined according to specific  $f$ -divergence setting. In the following, we instantiate this generic form under different  $f$ -divergence settings and provide explicit formulations for robust Bellman estimation as well as the coefficient  $c$  for the bonus term.

$$\begin{aligned}
Q_h^k(s, a) &= \min \left\{ \langle \phi(s, a), \boldsymbol{\theta}_h^k + \mathbf{w}_h^k \rangle + \Gamma_h^k(s, a), H - h + 1 \right\}, \\
\Gamma_h^k(s, a) &= c \sum_{i=1}^d \|\phi_i(s, a) \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}},
\end{aligned} \quad (6.3)$$

---

**Algorithm 2** RSPVI for linear setting

---

**Require:** transition regularizer  $\beta > 0$ , policy regularizer  $\Omega^\eta$ , ridge regression regularizer  $\lambda > 0$ .

```
1: for  $k = 1, \dots, K$  do
2:    $V_{H+1}^{k,\beta,\eta}(\cdot) \leftarrow 0$ .
3:   for  $h = H, \dots, 1$  do
4:     for  $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$  do
5:       Update  $Q$ -function estimation  $Q_h^{k,\beta,\eta}(s, a)$  according to (6.3).
6:     end for
7:     for  $\forall s \in \mathcal{S}$  do
8:       Update policy  $\pi_h^k(\cdot|s)$  according to (4.3).
9:       Update value function estimation  $V_h^{k,\beta,\eta}(s) = \langle Q_h^{k,\beta,\eta}(s, \cdot), \pi_h^k(\cdot|s) \rangle - \Omega_s^\eta(\pi_h^k)$ .
10:    end for
11:  end for
12:  Collect a trajectory  $\tau^k$  by executing  $\pi^k$ .
13: end for
```

---

where  $\Lambda_h^k = \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top + \lambda I$ ,  $\lambda > 0$  is the regularizer in the ridge regression.

**Linear-TV Setting** According to the linear structure assumption in [Assumption 6.1](#) and the dual formulation for regularized TV setting in [Lemma D.3](#), given the estimated robust value function  $V_{h+1}^{k,\beta,\eta}$ , we estimate the parameters  $\theta_h^k$  and  $w_h^k$  as follows

$$\begin{aligned}\theta_h^k &= \operatorname{argmin}_{\theta \in \mathbb{R}^d} \sum_{\tau=1}^k (r_h^\tau - \phi(s_h^\tau, a_h^\tau)^\top \theta)^2 + \lambda \|\theta\|_2^2, \\ w_h^k &= \operatorname{argmin}_{w \in \mathbb{R}^d} \sum_{\tau=1}^k ([V_{h+1}^{k,\beta,\eta}(s_{h+1}^\tau)]_{\min_{s' \in \mathcal{S}} V_{h+1}^{k,\beta,\eta}(s') + \beta} - \phi(s_h^\tau, a_h^\tau)^\top w)^2 + \lambda \|w\|_2^2.\end{aligned}$$

The bonus coefficient is given by  $c = Hd \cdot \xi_{\text{TV}}$ , where

$$\xi_{\text{TV}} = 720 + 3\sqrt{40 \log(96K^{13/2}H|\mathcal{A}|^3/\delta)}.$$

**Linear-KL Setting** According to the linear structure assumption in [Assumption 6.1](#) and the dual formulation for regularized KL setting in [Lemma D.6](#), we estimate  $w_h^k$  through two steps. First, given the estimated robust value function  $V_{h+1}^{k,\beta,\eta}$ , we estimate  $\widehat{\mathbb{E}}_{s \sim \mu^o}[e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}]$  as

$$\widehat{\mathbb{E}}_{s \sim \mu^o}[e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] = \operatorname{argmin}_{w \in \mathbb{R}^d} \sum_{\tau=1}^k (e^{-V_{h+1}^{k,\beta,\eta}(s_{h+1}^\tau)/\beta} - \phi(s_h^\tau, a_h^\tau)^\top w)^2 + \lambda \|w\|_2^2.$$

And then we estimate the parameters  $\theta_h^k$  and  $w_h^k$  as

$$\begin{aligned}\theta_h^k &= \operatorname{argmin}_{\theta \in \mathbb{R}^d} \sum_{\tau=1}^k (r_h^\tau - \phi(s_h^\tau, a_h^\tau)^\top \theta)^2 + \lambda \|\theta\|_2^2, \\ w_h^k &= -\beta \log \max \{ \widehat{\mathbb{E}}_{s \sim \mu^o}[e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}], e^{-H/\beta} \}.\end{aligned}$$

The truncation here is to make  $V_h^k$  lower bounded and therefore the ridge regression operation is well-defined. And the bonus coefficient is given by  $c = (1 + 2\beta e^{\beta^{-1}H})Hd \cdot \xi_{\text{KL}}$ , where

$$\xi_{\text{KL}} = 80 + \sqrt{40 \log (64\beta K^{11/2} H d^{1/2} (1 + 2\beta e^{\beta^{-1}H})^2 |\mathcal{A}|^3 / \delta)}.$$

**Linear- $\chi^2$  Setting** According to the linear structure assumption in [Assumption 6.1](#) and the dual formulation for regularized  $\chi^2$  setting in [Lemma D.9](#), we estimate  $\mathbf{w}_h^k$  through two steps. First, given the estimated robust value function  $V_{h+1}^{k,\beta,\eta}$  and an arbitrary constant  $\alpha$ , we estimate  $\widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}(s)]_\alpha$  and  $\widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}(s)]_\alpha^2$  as

$$\begin{aligned} \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}(s)]_\alpha &= \left[ \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{\tau=1}^k ([V_{h+1}^{k,\beta,\eta}(s_h^\tau)]_\alpha - \phi(s_h^\tau, a_h^\tau)^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2 \right]_{[0,H]}, \\ \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}(s)]_\alpha^2 &= \left[ \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{\tau=1}^k ([V_{h+1}^{k,\beta,\eta}(s_h^\tau)]_\alpha^2 - \phi(s_h^\tau, a_h^\tau)^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2 \right]_{[0,H]}. \end{aligned}$$

Based on the estimated expectations and the dual variable  $\alpha$ , we estimate the parameters  $\boldsymbol{\theta}_h^k$  and  $\mathbf{w}_h^k$  as

$$\begin{aligned} \boldsymbol{\theta}_h^k &= \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{\tau=1}^k (r_h^\tau - \phi(s_h^\tau, a_h^\tau)^\top \boldsymbol{\theta})^2 + \lambda \|\boldsymbol{\theta}\|_2^2, \\ \mathbf{w}_h^k &= \sup_{\alpha \in [0,H]} \left\{ \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}(s)]_\alpha - \frac{1}{4\beta} \widehat{\operatorname{Var}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}(s)]_\alpha \right\} \\ &= \sup_{\alpha \in [0,H]} \left\{ \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}(s)]_\alpha + \frac{1}{4\beta} (\widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}(s)]_\alpha)^2 - \frac{1}{4\beta} \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}(s)]_\alpha^2 \right\}. \end{aligned}$$

The bonus coefficient is given by  $c = (1 + H/(2\beta))Hd \cdot \xi_{\chi^2}$ , where

$$\xi_{\chi^2} = 720 + 3\sqrt{40 \log (96K^6 H^5 (1 + H/(2\beta))^3 |\mathcal{A}|^3 / \delta)}.$$

### 6.3 Theoretical Results

As in the tabular setting, we first introduce the necessary definitions and assumptions.

**Definition 6.3** (Worst-case transition). In the linear setting, since the transition model  $P_h^o(\cdot|s, a) = \langle \phi(s, a), \boldsymbol{\mu}_h^o(\cdot) \rangle$  is a factor distribution here, for any policy  $\pi$ , we first define  $\boldsymbol{\mu}_h^{w,\pi} = (\mu_{h,1}^{w,\pi}, \dots, \mu_{h,d}^{w,\pi})^\top$ , where for any  $i \in [d]$ ,

$$\mu_{h,i}^{w,\pi}(\cdot) = \underset{\mu_h \in \Delta(S)}{\operatorname{argmin}} \mathbb{E}_{\mu_h} [V_{h+1}^{\pi,\beta,\eta}](s, a) + \beta \cdot D(\mu_h(\cdot) \| \mu_{h,i}^o(\cdot)).$$

Then we can directly obtain the worst-case transition from the linear structure assumption [Assumption 6.1](#), written as

$$P_h^{w,\pi}(\cdot|s, a) = \langle \phi(s, a), \mu_{h,i}^{w,\pi}(\cdot) \rangle.$$

We define the worst-case transition differently in the tabular and linear settings, since in the linear case the transitions additionally admit a linear structure. With such definition in place, we can now naturally define the corresponding visitation measures as follows

**Definition 6.4** (Visitation measure). At timestep  $h \in [H]$ , we denote  $d_h^\pi(\cdot)$  as the visitation measure on  $\mathcal{S}$  induced by policy  $\pi$  under  $P^o$ , and  $q_h^\pi(\cdot)$  as the visitation measure on  $\mathcal{S}$  induced by policy  $\pi$  under  $P^{w,\pi}$ .

We now introduce the linear version of the bounded visitation measure ratio assumption as follows

**Assumption 6.5** (Bounded visitation measure ratio). Under the definition of [Definition 6.4](#), we define  $C_{vr} := \sup_{\pi, h, s} \frac{q_h^\pi(s)}{d_h^\pi(s)}$  as the supremal ratio between the nominal visitation measure and the worst-case visitation measure. We assume that  $C_{vr}$  is polynomial in  $H$ ,  $S$  and  $A$ .

Note that the formulations in [Definition 6.4](#) and [Assumption 6.5](#) are the same as those in [Definition 5.2](#) and [Assumption 5.3](#) in the tabular setting. This showcases that  $C_{vr}$  is an intrinsic quantity characterizing the difficulty of the online robust learning problem, independent of the linear structure assumption. Next, we present our theoretical results in the linear setting. As in [Liu and Xu \(2024a\)](#), we begin by presenting the instance-dependent upper bounds.

**Theorem 6.6.** Assume [Assumption 6.1](#) and [Assumption 6.5](#) holds for each  $f$ -divergence-based transition regularization, and consider the policy regularizations in (4.4), (4.5), and (4.6). Then by choosing  $\lambda = 1$ , for any  $\delta \in (0, 1/3)$ , with probability at least  $1 - 3\delta$ , the regret of [Algorithm 2](#) satisfies

$$\text{Regret}(K) \leq 2C_{vr} \left( \sqrt{2KH^3 \ln(2/\delta)} + c \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^d \|\phi_i(s_h^k, a_h^k) \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \right), \quad (6.4)$$

where  $c$  is the parameter in bonus coefficients defined in [Section 6.2](#) for each  $f$ -divergence setting.

[Theorem 6.6](#) provides the first finite-sample regret bound in both online distributionally robust learning under general  $f$ -divergences with linear function approximation and reward-robust learning. This result demonstrates that when [Assumption 6.5](#) is satisfied, efficient online robust learning is attainable. The formulation of (6.4) is similar to that of [Liu and Xu \(2024a\)](#). However, since we impose the bounded visitation measure assumption rather than the fail-state assumption ([Liu and Xu, 2024a](#), Assumption 4.1), the bound includes an additional dependence on  $C_{vr}$ .

**Theorem 6.7.** Assume [Assumption 6.1](#) holds for each  $f$ -divergence-based transition regularization, and consider the policy regularizations in (4.4), (4.5), and (4.6). Instead of assuming [Assumption 6.5](#), we now assume that

$$\mathbb{E}_\pi [\phi(s_h, a_h) \phi(s_h, a_h)^\top] \geq \alpha I \quad (6.5)$$

for all  $(\pi, h) \in \Pi \times [H]$ , where  $\alpha > 0$  is a constant. Then by choosing  $\lambda = 1$ , for any  $\delta \in (0, 1/3)$ , with probability at least  $1 - 3\delta$ , the regret of [Algorithm 2](#) satisfies

$$\text{Regret}(K) \leq 2cH\sqrt{K} \sqrt{\frac{128}{\alpha^2} \log \left( \frac{dHK}{\delta} \right) + \frac{2}{\alpha} \log(K)},$$

where  $c$  is the parameter in bonus coefficients defined in [Section 6.2](#) for each  $f$ -divergence setting.

As discussed in Liu and Xu (2024a), the  $d$ -rectangular estimation error in (6.4), which arises from the structure of the  $d$ -rectangular uncertainty set, cannot be bounded directly using the elliptical potential lemma (Abbasi-Yadkori et al., 2011). To address this difficulty, Liu and Xu (2024a, Corollary 5.3) introduce an additional (6.5) assumption. When specialized to the tabular setting,  $\phi(s_h, a_h)$  reduces to a one-hot vector. In this case, the diagonal entry of the matrix  $\mathbb{E}_\pi[\phi(s_h, a_h)\phi(s_h, a_h)^\top]$  corresponds to the probability of visiting  $(s, a) \in \mathcal{S} \times \mathcal{A}$  at step  $h$  under policy  $\pi$ . Hence, we have  $d_h^\pi(s, a) \geq \alpha$  for all  $(\pi, h, s, a) \in \Pi \times [H] \times \mathcal{S} \times \mathcal{A}$ , which directly implies Assumption 5.3. In the linear setting, (6.5) can be interpreted as a uniform coverage condition, requiring the nominal environment to be sufficiently exploratory. Theorem 6.7 shows that (6.5) can also serve as a sufficient condition for online robust learning.

Compared with Liu and Xu (2024a, Corollary 5.3), their result additionally relies on a fail-state assumption (Liu and Xu, 2024a, Assumption 4.1) in addition to (6.5). Moreover, Liu and Xu (2024a, Corollary 5.3) is restricted to transition uncertainty characterized by Total Variance, whereas Theorem 6.7 extends to more general  $f$ -divergence-based transition robustness.

## 7 EXPERIMENT

In this section, we validate our theoretical findings through an empirical study. It is worth noting that this is the first attempt in literature to make linear MDP practical in the context of robust RL. We choose the standard CartPole environment as a start, where the objective is to keep a pole balanced by moving a cart left or right. An episode is terminated if the pole falls or is truncated after 200 steps. The state space is continuous and represented by the tuple  $(x, \dot{x}, \theta, \dot{\theta})$ , corresponding to the cart’s position, velocity, the pole’s angle, and angular velocity, respectively. The action space is discrete with two possible actions, 0, 1, indicating the direction of force applied by the agent at each step. As long as the episode continues, the agent receives a reward of  $r = 1$  at each time step.

To implement RSPVI in the linear setting, we initially learn a feature mapping  $\phi$  using a discrete variational autoencoder (VAE) (Rolfe, 2016). The dimension of the feature mapping  $\phi$  is a hyperparameter that we tune from 5 to 100 in our experiment and select the one ( $d = 80$  in this CartPole environment) that effectively models the dynamics. Given the high feature dimension, we observe that computing the bonus in Algorithm 2 incurs significant computational overhead and becomes the primary bottleneck of the algorithm’s speed. Consequently, we opt not to include it in the implementation, albeit at the potential cost of slight performance degradation. During evaluation, we introduce perturbation to the agent by randomly selecting an action with a specified probability at each step. The learned policy is subsequently assessed by averaging the cumulative reward across 50 runs, and all results are further averaged across 10 random seeds.

Experiment results are shown in Figure 1, for which we choose TV and KL as the transition regularization and Negative Shannon Entropy as the policy regularization. As we can see from Figure 1, our algorithm RSPVI is less sensitive to the perturbation compared with the non-robust baseline. We also compare RSPVI with DR-LSVI-UCB (Liu and Xu, 2024a). We find that DR-LSVI-UCB fails to learn the robust policy without the exploration bonus, while the soft policies learned by RSPVI has exploration ability in nature. This further demonstrates the advantages of a soft policy and our doubly regularized framework in applications with large state-action spaces. More experimental results are provided in the appendix.

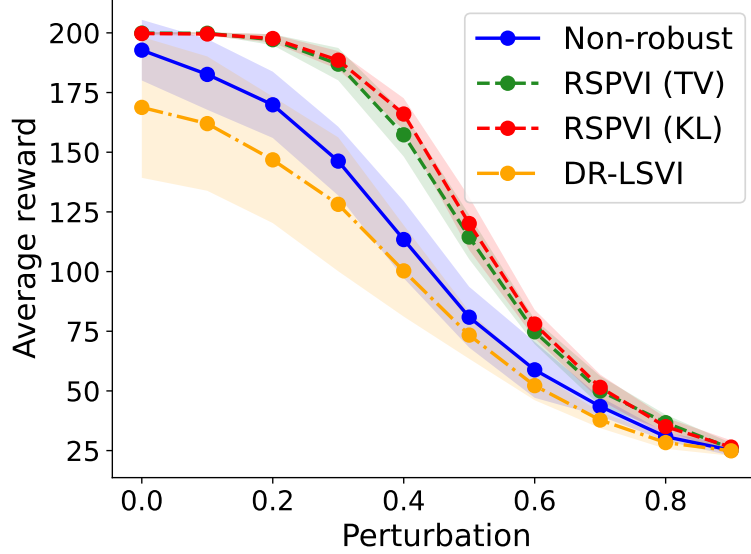


Figure 1: Experiment results for the CartPole environment. RSPVI is equipped with Negative Shannon entropy regularization.  $x$ -axis represents the perturbation level during evaluation.  $y$ -axis represents the average cumulative reward achieved over 50 episodes.

## 8 CONCLUSION

We propose a novel doubly regularized framework that achieves both reward robustness and transition robustness. Within this framework, we introduce an online learning algorithm, RSPVI, and prove that when using general  $f$ -divergence-based transition regularization together with various policy regularization formulations, the sample complexity remains no worse than that of traditional RMDPs. We further extend RSPVI to the linear setting, establishing the first finite-sample regret guarantee under general  $f$ -divergence. Compared with prior work, our results require fewer structural assumptions. Finally, we demonstrate the effectiveness and robustness of our framework through experiments in the CartPole environment.



## A Additional Details on Experiments

In this section, we provide more experimental results and details in addition to [Section 7](#). First, we report the hyper-parameters used in the experiments in [Table 1](#). In practice, we find that the ridge regression coefficient  $\lambda$  plays a crucial role in training and is tuned from  $\{1, 0.1, 0.01, 0.001, 0.0001\}$ .

First, we conduct ablation study to show the impact of  $\eta$  on the algorithm’s performance with negative Shannon entropy policy regularization.

**Negative Shannon Entropy** As shown in [Figures 2\(a\)](#) and [2\(d\)](#), it is clear that selecting an appropriate value for  $\eta$  is crucial for algorithm training. When  $\eta$  is too large, the excessive penalty on policy regularization can drive the policy toward near-random behavior, preventing the agent from gathering sufficient useful information during exploration. Conversely, when  $\eta$  is too small, the policy becomes less stochastic, which reduces exploration and may lead to the agent learning a suboptimal policy.

Since our theoretical framework also covers Kullback-Leibler divergence and negative Tsallis entropy policy regularization, we conduct additional experiments to evaluate the performance of RSPVI under these two settings.

**Kullback-Leibler Divergence** In [Figures 2\(b\)](#) and [2\(e\)](#), we present the results of RSPVI with KL divergence policy regularization for different values of  $\eta$ . The reference policy is learned with NS policy regularization but without transition regularization, making it non-robust to environmental perturbations. As  $\eta$  increases, the stronger policy regularization brings the learned policy closer to the reference policy. However, the overall performance of the learned policy still surpasses that of the reference policy. Further, the KL-divergence defined transition regularization leads to better robustness performance compared to the TV case.

**Negative Tsallis Entropy** In [Figures 2\(c\)](#) and [2\(f\)](#), we present the results of RSPVI with negative Tsallis entropy policy regularization for different values of  $\eta$ . With TV-distance defined transition regularization, similar to the behavior observed with Negative Shannon Entropy, the policy can be too random for large  $\eta$  and becomes too deterministic for small  $\eta$ , while a moderate  $\eta$  lead to the optimal performance. While for the case with KL-divergence defined transition regularization, our algorithm is not sensitive to the choice of  $\eta$ .

In general, this CartPole experiment suggests that our robust algorithm RSPVI outperforms the non-robust baseline.

Table 1: hyper-parameters for experiments

Symbol	Description	Value
$d$	latent space dimension	80
$K$	total episodes	2000
$\lambda$	ridge regression parameter	0.001
$\gamma$	discount factor	0.99
$\beta_{\text{TV}}$	transition regularization parameter	70
$\beta_{\text{KL}}$	transition regularization parameter	120

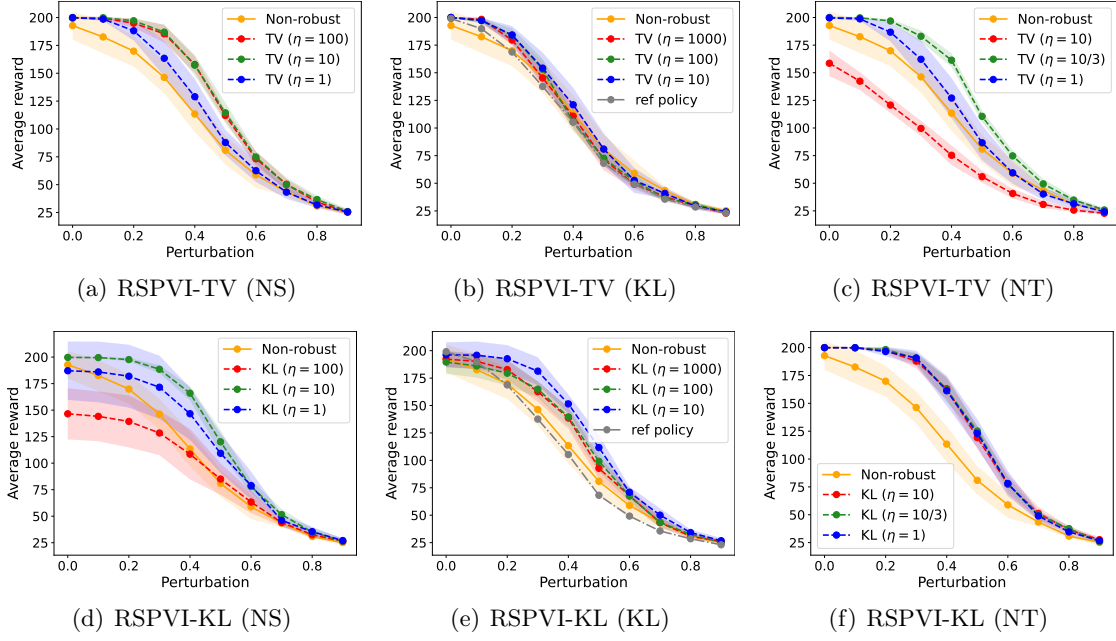


Figure 2: Performance of RSPVI with different types of policy regularization and values of  $\eta$ . Figures 2(a) to 2(c) use TV transition regularization, while Figures 2(d) to 2(f) use KL transition regularization. Additionally, Figures 2(a) and 2(d) apply negative Shannon entropy (NS) policy regularization, Figures 2(b) and 2(e) apply Kullback-Leibler divergence (KL) policy regularization, and Figures 2(c) and 2(f) apply negative Tsallis entropy (NT) policy regularization. The  $x$ -axis represents the perturbation level during evaluation, and the  $y$ -axis represents the average cumulative reward achieved over 50 episodes.

## B Proof of the Policy Update Formulation

In particular, we prove that the policy update formulation in (4.4), (4.5), (4.6) satisfies (4.3).

*Proof of the Kullback-Leibler Divergence case.* Fixing a reference policy  $\tilde{\pi}$ , we need to prove that policy

$$\pi(\cdot|s) \propto \tilde{\pi}_h^k(\cdot|s) \exp(Q_h^{k,\beta,\eta}(s, \cdot)/\eta)$$

maximizes the objective  $J(\pi)$

$$J(\pi) = \langle Q_h^{k,\beta,\eta}(s, \cdot), \pi(\cdot|s) \rangle - \eta \cdot \sum_{a \in \mathcal{A}} \pi(a|s) \ln \left( \frac{\pi(a|s)}{\tilde{\pi}_h^k(a|s)} \right).$$

In order to prove this, we define an auxiliary policy  $\pi'$  as

$$\pi'(a|s) = \frac{\tilde{\pi}_h^k(a|s) \exp(Q_h^{k,\beta,\eta}(s, a)/\eta)}{Z}, \quad \text{where} \quad Z = \sum_a \tilde{\pi}_h^k(a|s) \exp(Q_h^{k,\beta,\eta}(s, a)/\eta).$$

From the non-negativity of KL divergence, we have

$$\begin{aligned} 0 &\leq \sum_a \pi(a|s) \log \left( \frac{\pi(a|s)}{\pi'(a|s)} \right) \\ &= \sum_a \pi(a|s) \log \left( \frac{\pi(a|s)}{\tilde{\pi}_h^k(a|s)} \right) - \frac{1}{\eta} \sum_a \pi(a|s) Q_h^{k,\beta,\eta}(s, a) + \log Z. \end{aligned}$$

Rearrange this and it follows that

$$\sum_a \pi(a|s) Q_h^{k,\beta,\eta}(s, a) - \eta \sum_a \pi(a|s) \log \left( \frac{\pi(a|s)}{\tilde{\pi}_h^k(a|s)} \right) \leq \eta \log Z.$$

The equality is achieved if and only if  $\text{KL}(\pi||\pi') = 0$ , or  $\pi = \pi'$ . Since  $J(\pi)$  is strictly concave with respect to  $\pi$ ,  $\pi'$  is the unique maximizer of  $J(\pi)$ .  $\square$

*Proof of the Negative Shannon Entropy case.* We need to prove that policy

$$\pi(\cdot|s) \propto \exp(Q_h^{k,\beta,\eta}(s, \cdot)/\eta)$$

maximizes the objective  $J(\pi)$

$$J(\pi) = \langle Q_h^{k,\beta,\eta}(s, \cdot), \pi(\cdot|s) \rangle - \eta \cdot \sum_{a \in \mathcal{A}} \pi(a|s) \ln(\pi(a|s)).$$

The result is directly followed by the Kullback-Leibler Divergence case by choosing  $\tilde{\pi}(\cdot|s)$  as the uniform distribution.  $\square$

*Proof of the Negative Tsallis Entropy case.* We need to prove that policy

$$\pi(\cdot|s) = \text{proj}_\Delta(Q_h^{k,\beta,\eta}(s, \cdot)/\eta)$$

maximizes the objective  $J(\pi)$

$$J(\pi) = \langle Q_h^{k,\beta,\eta}(s, \cdot), \pi(\cdot|s) \rangle - \frac{\eta}{2} \|\pi(\cdot|s)\|^2,$$

In order to prove this, we rewrite the objective as

$$\begin{aligned} J(\pi) &= \langle Q_h^{k,\beta,\eta}(s, \cdot), \pi(\cdot|s) \rangle - \frac{\eta}{2} \|\pi(\cdot|s)\|^2 \\ &= -\frac{\eta}{2} \left\| \pi(\cdot|s) - \frac{1}{\eta} Q_h^{k,\beta,\eta}(s, \cdot) \right\|^2 + \frac{1}{2\eta} \|Q_h^{k,\beta,\eta}(s, \cdot)\|^2. \end{aligned}$$

Since  $\|Q_h^{k,\beta,\eta}(s, \cdot)\|^2/(2\eta)$  is independent of  $\pi$ , it is equivalent to minimize  $\|\pi(\cdot|s) - Q_h^{k,\beta,\eta}(s, \cdot)/\eta\|^2$ , subject to  $\pi(\cdot|s)$  being a distribution. As the projection operator  $\text{proj}$  minimizes the distance, its solution is

$$\pi(\cdot|s) = \text{proj}_\Delta(Q_h^{k,\beta,\eta}(s, \cdot)/\eta).$$

This finishes the proof.  $\square$

## C Proofs of the Results in Tabular Setting

### C.1 Proof of the Dynamic Programming Principle

*Proof of Proposition 3.1.* We prove a stronger version of the proposition. Specifically, there exists transition kernels  $\{\tilde{P}_t\}_{t=1}^H$ , such that for all  $(h, s) \in [H] \times \mathcal{S}$ ,

$$V_h^{\pi, \beta, \eta}(s) = \mathbb{E}_{\{\tilde{P}_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \cdot D(P_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) \middle| s_h = s, \pi_{h:H} \right]. \quad (\text{C.1})$$

We prove this statement by induction. For the base case  $h = H$ , the claim holds trivially since no transition kernels are involved. For the inductive step, assume the conclusion holds for step  $h + 1$ , meaning that there exists  $\{\tilde{r}_t\}_{t=1}^H$  and  $\{\tilde{P}_t\}_{t=h+1}^H$  such that

$$V_{h+1}^{\pi, \beta, \eta}(s) = \mathbb{E}_{\{\tilde{P}_t\}_{t=h+1}^H} \left[ \sum_{t=h+1}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \cdot D(P_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) \middle| s_{h+1} = s, \pi_{h+1:H} \right]. \quad (\text{C.2})$$

We now prove (C.1) for the case of  $h$ . For  $Q_h^{\pi, \beta, \eta}(s, a)$ , on the one hand, we upper bound it as

$$\begin{aligned} & Q_h^{\pi, \beta, \eta}(s, a) \\ &= \inf_{P_t \in \Delta(\mathcal{S})} \mathbb{E}_{\{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) + \beta \cdot D(P_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) - \sum_{t=h+1}^H \Omega_{s_t}^\eta(\pi_t) \middle| s_h = s, a_h = a, \pi_{h+1:H} \right] \\ &= r_h(s, a) + \inf_{P_t \in \Delta(\mathcal{S})} \left\{ \beta \cdot D(P_h(\cdot|s_h, a_h) \| P_h^o(\cdot|s_h, a_h)) + \int_{\mathcal{S}} P_h(ds'|s, a) \mathbb{E}_{\{P_t\}_{t=h+1}^H} \left[ \sum_{t=h+1}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \cdot D(P_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) \middle| s_{h+1} = s', \pi_{h+1:H} \right] \right\} \\ &\leq r_h(s, a) + \inf_{P_h \in \Delta(\mathcal{S})} \left\{ \beta \cdot D(P_h(\cdot|s_h, a_h) \| P_h^o(\cdot|s_h, a_h)) + \int_{\mathcal{S}} P_h(ds'|s, a) \mathbb{E}_{\{\tilde{P}_t\}_{t=h+1}^H} \left[ \sum_{t=h+1}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \cdot D(\tilde{P}_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) \middle| s_{h+1} = s', \pi_{h+1:H} \right] \right\} \end{aligned} \quad (\text{C.3})$$

$$= r_h(s, a) + \inf_{P_h \in \Delta(\mathcal{S})} \left\{ \beta \cdot D(P_h(\cdot|s_h, a_h) \| P_h^o(\cdot|s_h, a_h)) + \mathbb{E}_{s' \sim P_h^o(\cdot|s, a)} [V_{h+1}^{\pi, \beta, \eta}(s')] \right\}, \quad (\text{C.4})$$

where (C.3) is because the definition of the infimum operator and (C.4) is from (C.2).

On the other hand, we can lower bound  $Q_h^{\pi, \beta, \eta}(s, a)$  as

$$\begin{aligned} & Q_h^{\pi, \beta, \eta}(s, a) \\ &= \inf_{P_t \in \Delta(\mathcal{S})} \mathbb{E}_{\{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) + \beta \cdot D(P_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) - \sum_{t=h+1}^H \Omega_{s_t}^\eta(\pi_t) \middle| s_h = s, a_h = a, \pi_{h+1:H} \right] \\ &= r_h(s, a) + \inf_{P_t \in \Delta(\mathcal{S})} \left\{ \beta \cdot D(P_h(\cdot|s_h, a_h) \| P_h^o(\cdot|s_h, a_h)) + \int_{\mathcal{S}} P_h(ds'|s, a) \mathbb{E}_{\{P_t\}_{t=h+1}^H} \left[ \sum_{t=h+1}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \cdot D(P_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) \middle| s_{h+1} = s', \pi_{h+1:H} \right] \right\} \end{aligned}$$

$$\geq r_h(s, a) + \inf_{P_h \in \Delta(S)} \left\{ \beta \cdot D(P_h(\cdot|s_h, a_h) \| P_h^o(\cdot|s_h, a_h)) + \int_S P_h(ds'|s, a) \inf_{P_t \in \Delta(S)} \mathbb{E}_{\{P_t\}_{t=h+1}^H} \left[ \sum_{t=h+1}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \cdot D(P_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) \middle| s_{h+1} = s', \pi_{h+1:H} \right] \right\} \quad (\text{C.5})$$

$$= r_h(s, a) + \inf_{P_h \in \Delta(S)} \left\{ \beta \cdot D(P_h(\cdot|s_h, a_h) \| P_h^o(\cdot|s_h, a_h)) + \mathbb{E}_{s' \sim P_h^o(\cdot|s, a)} [V_{h+1}^{\pi, \beta, \eta}(s')] \right\}, \quad (\text{C.6})$$

where (C.5) is because  $\inf \int f \geq \int \inf f$  and (C.6) is from the definition of  $V_{h+1}^{\pi, \beta, \eta}(s')$ .

By combining (C.4) and (C.6), we have

$$Q_h^{\pi, \beta, \eta}(s, a) = r_h(s, a) + \inf_{\substack{\mu_h \in \Delta(S)^d \\ P_h = \langle \phi, \mu_h \rangle}} \left\{ \mathbb{E}_{s' \sim P_h^o(\cdot|s, a)} [V_{h+1}^{\pi, \beta, \eta}(s')] + \beta \cdot D(P_h(\cdot|s_h, a_h) \| P_h^o(\cdot|s_h, a_h)) \right\},$$

which implies (3.1).

Since  $\Delta(S)$  is compact, there exists  $\tilde{P}_h$ , such that

$$Q_h^{\pi, \beta, \eta}(s, a) = r_h(s, a) + \mathbb{E}_{s' \sim \tilde{P}_h(\cdot|s, a)} [V_{h+1}^{\pi, \beta, \eta}(s')] + \beta \cdot D(\tilde{P}_h(\cdot|s, a) \| P_h^o(\cdot|s, a)). \quad (\text{C.7})$$

This finishes the proof of (C.1).

For  $V_h^{\pi, \beta, \eta}(s)$ , on the one hand, we upper bound it as

$$\begin{aligned} & V_h^{\pi, \beta, \eta}(s) \\ &= \inf_{P_t \in \Delta(S)} \mathbb{E}_{\{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \cdot D(P_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) \middle| s_h = s, \pi_{h:H} \right] \\ &= \inf_{P_t \in \Delta(S)} \left\{ -\Omega_s^\eta(\pi_h) + \sum_{a \in \mathcal{A}} \pi_h(a|s) \mathbb{E}_{\{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) + \beta \cdot D(P_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) - \sum_{t=h+1}^H \Omega_{s_t}^\eta(\pi_t) \middle| s_h = s, a_h = a, \pi_{h+1:H} \right] \right\} \\ &\leq -\Omega_s^\eta(\pi_h) + \sum_{a \in \mathcal{A}} \pi_h(a|s) \mathbb{E}_{\{\tilde{P}_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) + \beta \cdot D(\tilde{P}_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) - \sum_{t=h+1}^H \Omega_{s_t}^\eta(\pi_t) \middle| s_h = s, a_h = a, \pi_{h+1:H} \right] \quad (\text{C.8}) \end{aligned}$$

$$= -\Omega_s^\eta(\pi_h) + \sum_{a \in \mathcal{A}} \pi_h(a|s) Q_h^{\pi, \beta, \eta}(s, a), \quad (\text{C.9})$$

where (C.8) is because the definition of the infimum operator and (C.9) is from (C.7) and (C.2).

On the other hand, we can lower bound  $V_h^{\pi, \beta, \eta}(s)$  as

$$\begin{aligned} & V_h^{\pi, \beta, \eta}(s) \\ &= \inf_{P_t \in \Delta(S)} \mathbb{E}_{\{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \cdot D(P_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) \middle| s_h = s, \pi_{h:H} \right] \\ &= \inf_{P_t \in \Delta(S)} \left\{ -\Omega_s^\eta(\pi_h) + \sum_{a \in \mathcal{A}} \pi_h(a|s) \mathbb{E}_{\{P_t\}_{t=h}^H} \left[ \right. \right. \end{aligned}$$

$$\begin{aligned}
& \left. \sum_{t=h}^H \left( r_t(s_t, a_t) + \beta \cdot D(P_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) - \sum_{t=h+1}^H \Omega_{s_t}^\eta(\pi_t) \right|_{s_h = s, a_h = a, \pi_{h+1:H}} \Bigg\} \\
& \geq -\Omega_s^\eta(\pi_h) + \sum_{a \in \mathcal{A}} \pi_h(a|s) \inf_{P_t \in \Delta(S)} \mathbb{E}_{\{P_t\}_{t=h}^H} \left[ \right. \\
& \quad \left. \sum_{t=h}^H \left( r_t(s_t, a_t) + \beta \cdot D(P_t(\cdot|s_t, a_t) \| P_t^o(\cdot|s_t, a_t)) \right) - \sum_{t=h+1}^H \Omega_{s_t}^\eta(\pi_t) \right]_{s_h = s, a_h = a, \pi_{h+1:H}} \quad (\text{C.10})
\end{aligned}$$

$$= -\Omega_s^\eta(\pi_h) + \sum_{a \in \mathcal{A}} \pi_h(a|s) Q_h^{\pi, \beta, \eta}(s, a), \quad (\text{C.11})$$

where (C.10) is because  $\inf \sum f \geq \sum \inf f$  and (C.11) is from the definition of  $Q_h^{\pi, \beta, \eta}(s, a)$ .

By combining (C.9) and (C.11), we have

$$V_h^{\pi, \beta, \eta}(s) = -\Omega_s^\eta(\pi_h) + \sum_{a \in \mathcal{A}} \pi_h(a|s) Q_h^{\pi, \beta, \eta}(s, a),$$

which implies (3.2).  $\square$

## C.2 Proof of the Double Robustness

*Proof of Proposition 3.2.* We have

$$\begin{aligned}
& \inf_{\substack{\tilde{r} \in \mathcal{R}^\eta(\pi) \\ P \in \Delta(S)}} \mathbb{E}_{a \sim \pi(\cdot|s)} [\tilde{r}(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} [V(s')] + \beta \cdot D(P(s, a) \| P^o(s, a))] \\
& = \inf_{\substack{\tilde{r} \in \mathcal{R}^\eta(\pi) \\ P \in \Delta(S)}} \left( \mathbb{E}_{a \sim \pi(\cdot|s)} [\tilde{r}(s, a)] + \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ \mathbb{E}_{s' \sim P(\cdot|s, a)} [V(s')] + \beta \cdot D(P(s, a) \| P^o(s, a)) \right] \right) \\
& = \inf_{\tilde{r} \in \mathcal{R}^\eta(\pi)} \sum_{a \in \mathcal{A}} \pi(a|s) \tilde{r}(s, a) + \inf_{P \in \Delta(S)} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ \mathbb{E}_{s' \sim P(\cdot|s, a)} [V(s')] + \beta \cdot D(P(s, a) \| P^o(s, a)) \right] \\
& = \inf_{\Delta r \in \mathcal{R}^\eta(\pi)} \sum_{a \in \mathcal{A}} \pi(a|s) \Delta r(s, a) + \inf_{P \in \Delta(S)} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ r(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} [V(s')] + \beta \cdot D(P(s, a) \| P^o(s, a)) \right] \\
& = \inf_{\Delta r \in \mathcal{R}^\eta(\pi)} \langle \pi(\cdot|s), \Delta r(s, \cdot) \rangle + \inf_{P \in \Delta(S)} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ r(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} [V(s')] + \beta \cdot D(P(s, a) \| P^o(s, a)) \right] \\
& = - \sup_{\Delta r \in \mathcal{R}^\eta(\pi)} \langle -\pi(\cdot|s), \Delta r(s, \cdot) \rangle + \inf_{P \in \Delta(S)} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ r(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} [V(s')] + \beta \cdot D(P(s, a) \| P^o(s, a)) \right] \\
& = -\sigma_{\mathcal{R}_s^\eta(\pi)}(-\pi(\cdot|s)) + \inf_{P \in \Delta(S)} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ r(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} [V(s')] + \beta \cdot D(P(s, a) \| P^o(s, a)) \right] \quad (\text{C.12})
\end{aligned}$$

$$\begin{aligned}
& = \inf_{P \in \Delta(S)} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ r(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} [V(s')] + \beta \cdot D(P(s, a) \| P^o(s, a)) \right] - \Omega_s^\eta(\pi) \quad (\text{C.13}) \\
& = \inf_{P \in \Delta(S)} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ r(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} [V(s')] + \beta \cdot D(P(s, a) \| P^o(s, a)) - \Omega_s^\eta(\pi) \right],
\end{aligned}$$

where (C.12) is from the definition that  $\sigma_{\mathcal{Z}}(\mathbf{y}) = \max_{a \in \mathcal{Z}} \langle \mathbf{a}, \mathbf{y} \rangle$ , (C.13) is because we choose  $\Omega_s^\eta(\pi) = \sigma_{\mathcal{R}_s^\eta(\pi)}(-\pi(\cdot|s))$ . This finishes the proof.  $\square$

## C.3 Proofs of the Main Results

To establish Theorem 5.4, we first introduce Lemma C.1, which bounds the regret by the sum of bonus terms under the expectation defined by the worst-case transition.

**Lemma C.1.** Under [Assumption 5.3](#), for each  $f$ -divergence setting and any  $\delta \in (0, 1/2)$ , with probability at least  $1 - 2\delta$ , the regret of [Algorithm 1](#) satisfies

$$\text{Regret}(K) \leq 2 \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{P^{w,k}, \pi^k} [b_h^k(s, a)], \quad (\text{C.14})$$

where  $b_h^k(s, a)$  is the bonus term defined in [Section 4.2](#) for each  $f$ -divergence setting.

To upper bound the sum of the expectations of  $\sqrt{\frac{1}{n_h^{k-1} \vee 1}}$  under the worst-case environment, [He et al. \(2025\)](#) proved the following lemma.

**Lemma C.2.** ([He et al., 2025](#), Lemma C.4, Lemma C.9) If the algorithm RSPVI selects  $\pi^k$  for the  $k$ -th episode, with probability at least  $1 - \delta$ , it holds that

$$\sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{P^{w,k}, \pi^k} \left[ \sqrt{\frac{1}{n_h^{k-1}(s, a) \vee 1}} \right] = \tilde{\mathcal{O}}(\sqrt{C_{vr}SAH^2K} + C_{vr}SAH).$$

*Proof of Theorem 5.4.* The result follows immediately from combining [Lemma C.1](#) and [Lemma C.2](#), and substituting the specific bonus formulation defined in [Section 4.2](#).  $\square$

## C.4 Proofs of the Technical Lemmas

Here, we present the proof of [Lemma C.1](#) is various  $f$ -divergence settings. For convenience, we also write  $P^{w,k} := P^{w, \pi^k}$ ,  $d^k := d^{\pi^k}$  and  $q^k := q^{\pi^k}$ .

### C.4.1 Proof of [Lemma C.1](#) in TV Setting

**Lemma C.3** (Dual formulation). ([He et al., 2025](#), Lemma D.1) For the optimization problem  $Q_h(s, a) = r_h(s, a) + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}] + \beta \text{TV}(P \| P_h^o))(s, a)$ , we have its dual formulation as follows

$$Q_h = r_h - \mathbb{E}_{P_h^o} \left[ \left( \min_{s \in S} V_{h+1}(s) + \beta - V_{h+1}(s) \right)_+ \right] + \left( \min_{s \in S} V_{h+1}(s) + \beta \right). \quad (\text{C.15})$$

**Lemma C.4** (Optimism). If we set the bonus term as follows

$$\text{bonus}_h^k(s, a) = 2H \sqrt{\frac{2S \ln(2SAHK/\delta)}{n_h^{k-1}(s, a) \vee 1}}, \quad (\text{C.16})$$

then for any policy  $\pi$  and any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ , with probability at least  $1 - 2\delta$ , we have  $Q_h^{k, \beta, \eta}(s, a) \geq Q_h^{\pi, \beta, \eta}(s, a)$ . Specially, by setting  $\pi = \pi^*$ , we have  $Q_h^{k, \beta, \eta}(s, a) \geq Q_h^{\pi^*, \beta, \eta}(s, a)$ .

*Proof.* We prove this by induction. First, when  $h = H + 1$ ,  $Q_{H+1}^{k, \beta, \eta}(s, a) = 0 = Q_{H+1}^{\pi, \beta, \eta}(s, a)$  holds trivially.

Assume  $Q_{h+1}^{k, \beta, \eta}(s, a) \geq Q_{h+1}^{\pi, \beta, \eta}(s, a)$  holds, due to the choice of  $\pi_{h+1}^k$  in [\(4.3\)](#), we have

$$\begin{aligned} V_{h+1}^{k, \beta, \eta}(s) &= \langle Q_{h+1}^{k, \beta, \eta}(s, \cdot), \pi_{h+1}^k(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}^k) \\ &\geq \langle Q_{h+1}^{k, \beta, \eta}(s, \cdot), \pi_{h+1}(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}) \end{aligned}$$



$$\geq \langle Q_{h+1}^{\pi, \beta, \eta}(s, \cdot), \pi_{h+1}(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}) = V_{h+1}^{\pi, \beta, \eta}(s).$$

Recall that we denote  $Q_h^{k, \beta, \eta}$  as the optimistic estimation in  $k$ -th episode, that is,

$$Q_h^{k, \beta, \eta}(s, a) = \min \left\{ \text{bonus}_h^k(s, a) + \hat{r}_h^k(s, a) + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{TV}(P \| \hat{P}_h^k))(s, a), H - h + 1 \right\}.$$

If  $Q_h^{k, \beta, \eta}(s, a) = H - h + 1$ , then it follows immediately that

$$Q_h^{k, \beta, \eta}(s, a) = H - h + 1 \geq Q_h^{\pi, \beta, \eta}(s, a)$$

by the definition of  $Q_h^{\pi, \beta, \eta}(s, a)$ . Otherwise, we can infer that

$$\begin{aligned} & Q_h^{k, \beta, \eta} - Q_h^{\pi, \beta, \eta} \\ &= \text{bonus}_h^k + \hat{r}_h^k + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{TV}(P \| \hat{P}_h^k)) - r_h - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{\pi, \beta, \eta}] + \beta \text{TV}(P \| P_h^o)) \\ &= \text{bonus}_h^k + \hat{r}_h^k - r_h + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{TV}(P \| \hat{P}_h^k)) - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{TV}(P \| P_h^o)) \\ &\quad + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{TV}(P \| P_h^o)) - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{\pi, \beta, \eta}] + \beta \text{TV}(P \| P_h^o)) \\ &\geq \text{bonus}_h^k + \hat{r}_h^k - r_h + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] - \mathbb{E}_P[V_{h+1}^{\pi, \beta, \eta}]) \\ &\quad + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{TV}(P \| \hat{P}_h^k)) - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{TV}(P \| P_h^o)) \end{aligned} \quad (\text{C.17})$$

$$\geq \text{bonus}_h^k + \hat{r}_h^k - r_h + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{TV}(P \| \hat{P}_h^k)) - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{TV}(P \| P_h^o)) \quad (\text{C.18})$$

$$\begin{aligned} &= \text{bonus}_h^k + \hat{r}_h^k - r_h - \mathbb{E}_{P_h^o} \left[ \left( \min_{s \in S} V_{h+1}^{k, \beta, \eta}(s) + \beta - V_{h+1}^{k, \beta, \eta}(s) \right)_+ \right] + \left( \min_{s \in S} V_{h+1}^{k, \beta, \eta}(s) + \beta \right) \\ &\quad + \mathbb{E}_{\hat{P}_h^k} \left[ \left( \min_{s \in S} V_{h+1}^{k, \beta, \eta}(s) + \beta - V_{h+1}^{k, \beta, \eta}(s) \right)_+ \right] - \left( \min_{s \in S} V_{h+1}^{k, \beta, \eta}(s) + \beta \right) \end{aligned} \quad (\text{C.19})$$

$$\begin{aligned} &= \text{bonus}_h^k + \hat{r}_h^k - r_h \\ &\quad - \mathbb{E}_{P_h^o} \left[ \left( \min_{s \in S} V_{h+1}^{k, \beta, \eta}(s) + \beta - V_{h+1}^{k, \beta, \eta}(s) \right)_+ \right] + \mathbb{E}_{\hat{P}_h^k} \left[ \left( \min_{s \in S} V_{h+1}^{k, \beta, \eta}(s) + \beta - V_{h+1}^{k, \beta, \eta}(s) \right)_+ \right] \\ &\geq \text{bonus}_h^k - \underbrace{|\hat{r}_h^k - r_h|}_{(i)} \\ &\quad - \underbrace{\left| \mathbb{E}_{P_h^o} \left[ \left( \min_{s \in S} V_{h+1}^{k, \beta, \eta}(s) + \beta - V_{h+1}^{k, \beta, \eta}(s) \right)_+ \right] - \mathbb{E}_{\hat{P}_h^k} \left[ \left( \min_{s \in S} V_{h+1}^{k, \beta, \eta}(s) + \beta - V_{h+1}^{k, \beta, \eta}(s) \right)_+ \right] \right|}_{(ii)}, \end{aligned} \quad (\text{C.20})$$

where (C.17) is from  $\inf f(x) - \inf g(x) \geq \inf(f - g)(x)$ , (C.18) is from the induction assumption, we plug in the dual formulation (C.15) in (C.19).

For term (i) in (C.20), from Lemma E.1 and a union bound, with probability at least  $1 - \delta$ , we have

$$|\hat{r}_h^k(s, a) - r_h(s, a)| \leq \sqrt{\frac{\ln(2SAHK/\delta)}{2n_h^{k-1}(s, a) \vee 1}} \quad (\text{C.21})$$

for any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ .

For any fixed  $V$ , we apply [Lemma E.3](#) and have

$$|\mathbb{E}_{P_h^o}[V] - \mathbb{E}_{\hat{P}_h^k}[V]| \leq \|P_h^o - \hat{P}_h^k\|_1 \cdot \|V\|_\infty \leq H \sqrt{\frac{2S \ln(2/\delta)}{n_h^{k-1} \vee 1}}, \quad (\text{C.22})$$

with probability at least  $1 - \delta$ .

For term (ii) in [\(C.20\)](#), by applying [\(C.22\)](#) and a union bound, with probability at least  $1 - \delta$ , we have

$$\begin{aligned} & \left| \mathbb{E}_{P_h^o} \left[ \left( \min_{s \in \mathcal{S}} V_{h+1}^{k, \beta, \eta}(s) + \beta - V_{h+1}^{k, \beta, \eta}(s) \right)_+ \right] - \mathbb{E}_{\hat{P}_h^k} \left[ \left( \min_{s \in \mathcal{S}} V_{h+1}^{k, \beta, \eta}(s) + \beta - V_{h+1}^{k, \beta, \eta}(s) \right)_+ \right] \right| \\ & \leq H \sqrt{\frac{2S \ln(2SAHK/\delta)}{n_h^{k-1} \vee 1}} \end{aligned} \quad (\text{C.23})$$

for any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ .

Apply the union bound again and combine [\(C.20\)](#) with [\(C.21\)](#), [\(C.23\)](#), the definition of bonus and induction assumption. With probability at least  $1 - 2\delta$ , we have  $Q_h^{k, \rho}(s, a) \geq Q_h^{\pi, \rho}(s, a)$  for any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ . This completes the proof.  $\square$

**Lemma C.5.** If we set the bonus term to be the same as in [Lemma C.4](#), then for any  $(k, s, a) \in [K] \times \mathcal{S} \times \mathcal{A}$ , with probability at least  $1 - \delta$ , the sum of estimation errors can be bounded as

$$Q_1^{k, \beta, \eta}(s, a) - Q_1^{\pi^k, \beta, \eta}(s, a) \leq 2 \cdot \mathbb{E}_{P^{w, k}, \pi^k} [\text{bonus}_h^k],$$

*Proof.* From the proof of [Lemma C.4](#), we see that with probability at least  $1 - \delta$ , for any  $(k, s, a) \in [K] \times \mathcal{S} \times \mathcal{A}$ ,

$$\begin{aligned} & |\hat{r}_h^k(s, a) - r_h(s, a)| + \left| \inf_{P \in \Delta(\mathcal{S})} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{TV}(P \| \hat{P}_h^k))(s, a) \right. \\ & \quad \left. - \inf_{P \in \Delta(\mathcal{S})} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{TV}(P \| P_h^o))(s, a) \right| \leq \text{bonus}_h^k(s, a). \end{aligned} \quad (\text{C.24})$$

Recall that we define  $P_h^{w, k} = \underset{P \in \Delta(\mathcal{S})}{\text{argmin}} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{TV}(P \| P_h^o))$  as the worst-case transition in

[Definition 5.1](#), we have

$$\begin{aligned} & Q_h^{k, \beta, \eta} - Q_h^{\pi^k, \beta, \eta} \\ & \leq \text{bonus}_h^k + \hat{r}_h^k + \inf_{P \in \Delta(\mathcal{S})} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{TV}(P \| \hat{P}_h^k)) - r_h - \inf_{P \in \Delta(\mathcal{S})} (\mathbb{E}_P[V_{h+1}^{\pi^k, \beta, \eta}] + \beta \text{TV}(P \| P_h^o)) \\ & \leq 2\text{bonus}_h^k + \inf_{P \in \Delta(\mathcal{S})} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{TV}(P \| P_h^o)) - \inf_{P \in \Delta(\mathcal{S})} (\mathbb{E}_P[V_{h+1}^{\pi^k, \beta, \eta}] + \beta \text{TV}(P \| P_h^o)) \end{aligned} \quad (\text{C.25})$$

$$= 2\text{bonus}_h^k + \mathbb{E}_{P_h^{w, k}, \pi_h^k} [Q_{h+1}^{k, \beta, \eta} - Q_{h+1}^{\pi^k, \beta, \eta}]. \quad (\text{C.26})$$

where [\(C.25\)](#) uses [\(C.24\)](#). Apply [\(C.26\)](#) recursively, we can obtain the result.  $\square$

*Proof of Lemma C.1 in TV Setting.* The result directly follows from combining [Lemma C.4](#) and [Lemma C.5](#), along with applying a union bound.  $\square$

#### C.4.2 Proof of Lemma C.1 in KL Setting

**Lemma C.6** (Dual formulation). (He et al., 2025, Lemma D.5) For the optimization problem  $Q_h(s, a) = r_h(s, a) + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}] + \beta \text{KL}(P \| P_h^o))(s, a)$ , we have its dual formulation as follows

$$Q_h = r_h - \beta \ln \mathbb{E}_{P_h^o} [e^{-\beta^{-1} V_{h+1}}]. \quad (\text{C.27})$$

**Lemma C.7** (Optimism). If we set the bonus term as follows

$$\text{bonus}_h^k(s, a) = (1 + \beta e^{\beta^{-1} H} \sqrt{S}) \sqrt{\frac{2 \ln(2SAHK/\delta)}{n_h^{k-1}(s, a) \vee 1}}, \quad (\text{C.28})$$

then for any policy  $\pi$  and any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ , with probability at least  $1 - 2\delta$ , we have  $Q_h^{k, \beta, \eta}(s, a) \geq Q_h^{\pi, \beta, \eta}(s, a)$ . Specially, by setting  $\pi = \pi^*$ , we have  $Q_h^{k, \beta, \eta}(s, a) \geq Q_h^{\pi^*, \beta, \eta}(s, a)$ .

*Proof.* We prove this by induction. First, when  $h = H + 1$ ,  $Q_{H+1}^{k, \beta, \eta}(s, a) = 0 = Q_{H+1}^{\pi, \beta, \eta}(s, a)$  holds trivially.

Assume  $Q_{h+1}^{k, \beta, \eta}(s, a) \geq Q_{h+1}^{\pi, \beta, \eta}(s, a)$  holds, due to the choice of  $\pi_{h+1}^k$  in (4.3), we have

$$\begin{aligned} V_{h+1}^{k, \beta, \eta}(s) &= \langle Q_{h+1}^{k, \beta, \eta}(s, \cdot), \pi_{h+1}^k(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}^k) \\ &\geq \langle Q_{h+1}^{k, \beta, \eta}(s, \cdot), \pi_{h+1}(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}) \\ &\geq \langle Q_{h+1}^{\pi, \beta, \eta}(s, \cdot), \pi_{h+1}(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}) = V_{h+1}^{\pi, \beta, \eta}(s). \end{aligned}$$

Recall that we denote  $Q_h^{k, \beta, \eta}$  as the optimistic estimation in  $k$ -th episode, that is,

$$Q_h^{k, \beta, \eta}(s, a) = \min \left\{ \text{bonus}_h^k(s, a) + \hat{r}_h^k(s, a) + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{KL}(P \| \hat{P}_h^k))(s, a), H - h + 1 \right\}.$$

If  $Q_h^{k, \beta, \eta}(s, a) = H - h + 1$ , then it follows immediately that

$$Q_h^{k, \beta, \eta}(s, a) = H - h + 1 \geq Q_h^{\pi, \beta, \eta}(s, a)$$

by the definition of  $Q_h^{\pi, \beta, \eta}(s, a)$ . Otherwise, we can infer that

$$\begin{aligned} &Q_h^{k, \beta, \eta} - Q_h^{\pi, \beta, \eta} \\ &= \text{bonus}_h^k + \hat{r}_h^k + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{KL}(P \| \hat{P}_h^k)) - r_h - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{\pi, \beta, \eta}] + \beta \text{KL}(P \| P_h^o)) \\ &= \text{bonus}_h^k + \hat{r}_h^k - r_h + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{KL}(P \| \hat{P}_h^k)) - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{KL}(P \| P_h^o)) \\ &\quad + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{KL}(P \| P_h^o)) - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{\pi, \beta, \eta}] + \beta \text{KL}(P \| P_h^o)) \\ &\geq \text{bonus}_h^k + \hat{r}_h^k - r_h + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] - \mathbb{E}_P[V_{h+1}^{\pi, \beta, \eta}]) \\ &\quad + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{KL}(P \| \hat{P}_h^k)) - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{KL}(P \| P_h^o)) \end{aligned} \quad (\text{C.29})$$

$$\geq \text{bonus}_h^k + \hat{r}_h^k - r_h + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{KL}(P \| \hat{P}_h^k)) - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \text{KL}(P \| P_h^o)) \quad (\text{C.30})$$

$$= \text{bonus}_h^k + \hat{r}_h^k - r_h + \beta \ln \mathbb{E}_{P_h^o} [e^{-\beta^{-1} V_{h+1}^{k,\beta,\eta}}] - \beta \ln \mathbb{E}_{\hat{P}_h^k} [e^{-\beta^{-1} V_{h+1}^{k,\beta,\eta}}] \quad (\text{C.31})$$

$$\geq \text{bonus}_h^k - \underbrace{|\hat{r}_h^k - r_h|}_{(i)} - \underbrace{\beta |\ln \mathbb{E}_{\hat{P}_h^k} [e^{-\beta^{-1} V_{h+1}^{k,\beta,\eta}}] - \ln \mathbb{E}_{P_h^o} [e^{-\beta^{-1} V_{h+1}^{k,\beta,\eta}}]}_{(ii)}, \quad (\text{C.32})$$

where (C.29) is from  $\inf f(x) - \inf g(x) \geq \inf(f - g)(x)$ , (C.30) is from the induction assumption, we plug in the dual formulation (C.27) in (C.31).

For term (i) in (C.32), from Lemma E.1 and a union bound, with probability at least  $1 - \delta$ , we have

$$|\hat{r}_h^k(s, a) - r_h(s, a)| \leq \sqrt{\frac{\ln(2SAHK/\delta)}{2n_h^{k-1}(s, a) \vee 1}} \quad (\text{C.33})$$

for any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ .

For term (ii) in (C.32), by applying (C.22) and a union bound, with probability at least  $1 - \delta$ , we have

$$\left| \ln \mathbb{E}_{\hat{P}_h^k} [e^{-\beta^{-1} V_{h+1}^{k,\beta,\eta}}] - \ln \mathbb{E}_{P_h^o} [e^{-\beta^{-1} V_{h+1}^{k,\beta,\eta}}] \right| \leq e^{\beta^{-1} H} \sqrt{\frac{2S \ln(2SAHK/\delta)}{n_h^{k-1} \vee 1}} \quad (\text{C.34})$$

for any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ .

Apply the union bound again and combine (C.32) with (C.33), (C.34), the definition of bonus and induction assumption. With probability at least  $1 - 2\delta$ , we have  $Q_h^{k,\rho}(s, a) \geq Q_h^{\pi^k,\rho}(s, a)$  for any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ . This completes the proof.  $\square$

**Lemma C.8.** If we set the bonus term to be the same as in Lemma C.7, then for any  $(k, s, a) \in [K] \times \mathcal{S} \times \mathcal{A}$ , with probability at least  $1 - \delta$ , the sum of estimation errors can be bounded as

$$Q_1^{k,\beta,\eta}(s, a) - Q_1^{\pi^k,\beta,\eta}(s, a) \leq 2 \cdot \mathbb{E}_{P^{w,k}, \pi^k} [\text{bonus}_h^k],$$

*Proof.* From the proof of Lemma C.7, we see that with probability at least  $1 - \delta$ , for any  $(k, s, a) \in [K] \times \mathcal{S} \times \mathcal{A}$ ,

$$\begin{aligned} & \left| \hat{r}_h^k(s, a) - r_h(s, a) \right| + \left| \inf_{P \in \Delta(\mathcal{S})} (\mathbb{E}_P [V_{h+1}^{k,\beta,\eta}] + \beta \text{KL}(P \| \hat{P}_h^k))(s, a) \right. \\ & \quad \left. - \inf_{P \in \Delta(\mathcal{S})} (\mathbb{E}_P [V_{h+1}^{k,\beta,\eta}] + \beta \text{KL}(P \| P_h^o))(s, a) \right| \leq \text{bonus}_h^k(s, a). \end{aligned} \quad (\text{C.35})$$

Recall that we define  $P_h^{w,k} = \argmin_{P \in \Delta(\mathcal{S})} (\mathbb{E}_P [V_{h+1}^{\pi^k,\beta,\eta}] + \beta \text{KL}(P \| P_h^o))$  as the worst-case transition in

Definition 5.1, we have

$$\begin{aligned} & Q_h^{k,\beta,\eta} - Q_h^{\pi^k,\beta,\eta} \\ & \leq \text{bonus}_h^k + \hat{r}_h^k + \inf_{P \in \Delta(\mathcal{S})} (\mathbb{E}_P [V_{h+1}^{k,\beta,\eta}] + \beta \text{KL}(P \| \hat{P}_h^k)) - r_h - \inf_{P \in \Delta(\mathcal{S})} (\mathbb{E}_P [V_{h+1}^{\pi^k,\beta,\eta}] + \beta \text{KL}(P \| P_h^o)) \\ & \leq 2\text{bonus}_h^k + \inf_{P \in \Delta(\mathcal{S})} (\mathbb{E}_P [V_{h+1}^{k,\beta,\eta}] + \beta \text{KL}(P \| P_h^o)) - \inf_{P \in \Delta(\mathcal{S})} (\mathbb{E}_P [V_{h+1}^{\pi^k,\beta,\eta}] + \beta \text{KL}(P \| P_h^o)) \end{aligned} \quad (\text{C.36})$$

$$= 2\text{bonus}_h^k + \mathbb{E}_{P_h^{w,k}, \pi^k} [Q_{h+1}^{k,\beta,\eta} - Q_{h+1}^{\pi^k,\beta,\eta}]. \quad (\text{C.37})$$

where (C.36) uses (C.35). Apply (C.37) recursively, we can obtain the result.  $\square$

*Proof of Lemma C.1 in KL Setting.* The result directly follows from combining Lemma C.7 and Lemma C.8, along with applying a union bound.  $\square$

### C.4.3 Proof of Lemma C.1 in $\chi^2$ Setting

**Lemma C.9** (Dual formulation). (He et al., 2025, Lemma D.9) For the optimization problem  $Q_h(s, a) = r_h(s, a) + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}] + \beta \chi^2(P \| P_h^o))(s, a)$ , we have its dual formulation as follows

$$Q_h = r_h + \sup_{\lambda \in [0, H]} (\mathbb{E}_{P_h^o}[V_{h+1} - \lambda] - \frac{1}{4\beta} \text{Var}_{P_h^o}(V_{h+1} - \lambda)). \quad (\text{C.38})$$

**Lemma C.10** (Optimism). If we set the bonus term as follows

$$\text{bonus}_h^k(s, a) = \left(2 + \frac{3H}{4\beta}\right) H \sqrt{\frac{2S^2 \ln(48SAH^3K^2/\delta)}{n_h^{k-1}(s, a) \vee 1}} + \left(1 + \frac{1}{4\beta}\right) \frac{1}{K}, \quad (\text{C.39})$$

then for any policy  $\pi$  and any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ , with probability at least  $1 - 3\delta$ , we have  $Q_h^{k, \beta, \eta}(s, a) \geq Q_h^{\pi, \beta, \eta}(s, a)$ . Specially, by setting  $\pi = \pi^*$ , we have  $Q_h^{k, \beta, \eta}(s, a) \geq Q_h^{\pi^*, \beta, \eta}(s, a)$ .

*Proof.* We prove this by induction. First, when  $h = H + 1$ ,  $Q_{H+1}^{k, \beta, \eta}(s, a) = 0 = Q_{H+1}^{\pi, \beta, \eta}(s, a)$  holds trivially.

Assume  $Q_{h+1}^{k, \beta, \eta}(s, a) \geq Q_{h+1}^{\pi, \beta, \eta}(s, a)$  holds, due to the choice of  $\pi_{h+1}^k$  in (4.3), we have

$$\begin{aligned} V_{h+1}^{k, \beta, \eta}(s) &= \langle Q_{h+1}^{k, \beta, \eta}(s, \cdot), \pi_{h+1}^k(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}^k) \\ &\geq \langle Q_{h+1}^{k, \beta, \eta}(s, \cdot), \pi_{h+1}(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}) \\ &\geq \langle Q_{h+1}^{\pi, \beta, \eta}(s, \cdot), \pi_{h+1}(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}) = V_{h+1}^{\pi, \beta, \eta}(s). \end{aligned}$$

Recall that we denote  $Q_h^{k, \beta, \eta}$  as the optimistic estimation in  $k$ -th episode, that is,

$$Q_h^{k, \beta, \eta}(s, a) = \min \left\{ \text{bonus}_h^k(s, a) + \hat{r}_h^k(s, a) + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \chi^2(P \| \hat{P}_h^k))(s, a), H - h + 1 \right\}.$$

If  $Q_h^{k, \beta, \eta}(s, a) = H - h + 1$ , then it follows immediately that

$$Q_h^{k, \beta, \eta}(s, a) = H - h + 1 \geq Q_h^{\pi, \beta, \eta}(s, a)$$

by the definition of  $Q_h^{\pi, \beta, \eta}(s, a)$ . Otherwise, we can infer that

$$\begin{aligned} &Q_h^{k, \beta, \eta} - Q_h^{\pi, \beta, \eta} \\ &= \text{bonus}_h^k + \hat{r}_h^k + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \chi^2(P \| \hat{P}_h^k)) - r_h - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{\pi, \beta, \eta}] + \beta \chi^2(P \| P_h^o)) \\ &= \text{bonus}_h^k + \hat{r}_h^k - r_h + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \chi^2(P \| \hat{P}_h^k)) - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \chi^2(P \| P_h^o)) \\ &\quad + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta \chi^2(P \| P_h^o)) - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{\pi, \beta, \eta}] + \beta \chi^2(P \| P_h^o)) \\ &\geq \text{bonus}_h^k + \hat{r}_h^k - r_h + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] - \mathbb{E}_P[V_{h+1}^{\pi, \beta, \eta}]) \end{aligned}$$

$$+ \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k,\beta,\eta}] + \beta\chi^2(P\|\hat{P}_h^k)) - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k,\beta,\eta}] + \beta\chi^2(P\|P_h^o)) \quad (\text{C.40})$$

$$\geq \text{bonus}_h^k + \hat{r}_h^k - r_h + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k,\beta,\eta}] + \beta\chi^2(P\|\hat{P}_h^k)) - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k,\beta,\eta}] + \beta\chi^2(P\|P_h^o)) \quad (\text{C.41})$$

$$= \text{bonus}_h^k + \hat{r}_h^k - r_h + \sup_{\lambda \in [0,H]} (\mathbb{E}_{\hat{P}_h^k}[V_{h+1}^{k,\beta,\eta} - \lambda] - \frac{1}{4\beta} \text{Var}_{\hat{P}_h^k}(V_{h+1}^{k,\beta,\eta} - \lambda)) \\ - \sup_{\lambda \in [0,H]} (\mathbb{E}_{P_h^o}[V_{h+1}^{k,\beta,\eta} - \lambda] - \frac{1}{4\beta} \text{Var}_{P_h^o}(V_{h+1}^{k,\beta,\eta} - \lambda)) \quad (\text{C.42})$$

$$\geq \text{bonus}_h^k + \hat{r}_h^k - r_h + \inf_{\lambda \in [0,H]} \left\{ (\mathbb{E}_{\hat{P}_h^k}[V_{h+1}^{k,\beta,\eta} - \lambda] - \frac{1}{4\beta} \text{Var}_{\hat{P}_h^k}(V_{h+1}^{k,\beta,\eta} - \lambda)) \right. \\ \left. - (\mathbb{E}_{P_h^o}[V_{h+1}^{k,\beta,\eta} - \lambda] - \frac{1}{4\beta} \text{Var}_{P_h^o}(V_{h+1}^{k,\beta,\eta} - \lambda)) \right\} \\ \geq \text{bonus}_h^k - |\hat{r}_h^k - r_h| - \sup_{\lambda \in [0,H]} \left| (\mathbb{E}_{\hat{P}_h^k}[V_{h+1}^{k,\beta,\eta} - \lambda] - \frac{1}{4\beta} \text{Var}_{\hat{P}_h^k}(V_{h+1}^{k,\beta,\eta} - \lambda)) \right. \\ \left. - (\mathbb{E}_{P_h^o}[V_{h+1}^{k,\beta,\eta} - \lambda] - \frac{1}{4\beta} \text{Var}_{P_h^o}(V_{h+1}^{k,\beta,\eta} - \lambda)) \right| \\ \geq \text{bonus}_h^k - \underbrace{|\hat{r}_h^k - r_h|}_{(i)} - \underbrace{\sup_{\lambda \in [0,H]} |\mathbb{E}_{\hat{P}_h^k}[V_{h+1}^{k,\beta,\eta} - \lambda] - \mathbb{E}_{P_h^o}[V_{h+1}^{k,\beta,\eta} - \lambda]|}_{(ii)} \\ - \underbrace{\frac{1}{4\beta} \sup_{\lambda \in [0,H]} |\text{Var}_{\hat{P}_h^k}(V_{h+1}^{k,\beta,\eta} - \lambda) - \text{Var}_{P_h^o}(V_{h+1}^{k,\beta,\eta} - \lambda)|}_{(iii)}, \quad (\text{C.43})$$

where (C.40) is from  $\inf f(x) - \inf g(x) \geq \inf(f - g)(x)$ , (C.41) is from the induction assumption, we plug in the dual formulation (C.38) in (C.42).

For term (i) in (C.43), from Lemma E.1 and a union bound, with probability at least  $1 - \delta$ , we have

$$|\hat{r}_h^k(s, a) - r_h(s, a)| \leq \sqrt{\frac{\ln(2SAHK/\delta)}{2n_h^{k-1}(s, a) \vee 1}}, \quad (\text{C.44})$$

for any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ .

We denote  $V(\lambda) = V_{h+1}^{k,\rho} - \lambda \in [-H, H]$  and  $\mathcal{V} = \{V \in \mathbb{R}^S : \|V\|_\infty \leq H\}$ . To bound term (ii) in (C.43), we create a  $\epsilon$ -net  $\mathcal{N}_\mathcal{V}(\epsilon)$  for  $\mathcal{V}$ . From Lemma E.4, it holds that  $\ln |\mathcal{N}_\mathcal{V}(\epsilon)| \leq |S| \cdot \ln(3H/\epsilon)$ .

Therefore, by the definition of  $\mathcal{N}_\mathcal{V}(\epsilon)$ , for any fixed  $V$ , there exists a  $V' \in \mathcal{N}_\mathcal{V}(\epsilon)$  such that  $\|V - V'\|_\infty \leq \epsilon$ , that is

$$|\mathbb{E}_{P_h^o}[V] - \mathbb{E}_{\hat{P}_h^k}[V]| \leq |\mathbb{E}_{P_h^o}[V] - \mathbb{E}_{P_h^o}[V']| + |\mathbb{E}_{P_h^o}[V'] - \mathbb{E}_{\hat{P}_h^k}[V']| + |\mathbb{E}_{\hat{P}_h^k}[V'] - \mathbb{E}_{\hat{P}_h^k}[V]| \\ \leq \|P_h^o\|_1 \|V - V'\|_\infty + |\mathbb{E}_{P_h^o}[V'] - \mathbb{E}_{\hat{P}_h^k}[V']| + \|\hat{P}_h^k\|_1 \|V - V'\|_\infty \\ \leq \sup_{V' \in \mathcal{N}_\mathcal{V}(\epsilon)} |\mathbb{E}_{P_h^o}[V'] - \mathbb{E}_{\hat{P}_h^k}[V']| + 2\epsilon, \quad (\text{C.45})$$

where the second inequality follows from the Holder's inequality.

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

$$\begin{aligned} \sup_{\boldsymbol{\lambda} \in [0, H]} |\mathbb{E}_{\hat{P}_h^k}[V_{h+1}^{k, \rho} - \boldsymbol{\lambda}] - \mathbb{E}_{P_h^o}[V_{h+1}^{k, \rho} - \boldsymbol{\lambda}]| &\leq \sup_{\eta \in [0, H]} |\mathbb{E}_{P_h^o}[V(\boldsymbol{\lambda})] - \mathbb{E}_{\hat{P}_h^k}[V(\boldsymbol{\lambda})]| \\ &\leq \sup_{V \in \mathcal{N}_{\mathcal{V}}(\epsilon)} |\mathbb{E}_{P_h^o}[V] - \mathbb{E}_{\hat{P}_h^k}[V]| + 2\epsilon \end{aligned} \quad (\text{C.46})$$

$$\leq H \sqrt{\frac{2S \ln(2SAHK|\mathcal{N}_{\mathcal{V}}(\epsilon)|/\delta)}{n_h^{k-1} \vee 1}} + 2\epsilon \quad (\text{C.47})$$

$$\begin{aligned} &\leq H \sqrt{\frac{2S^2 \ln(6SAH^2K/\epsilon\delta)}{n_h^{k-1} \vee 1}} + 2\epsilon \\ &= H \sqrt{\frac{2S^2 \ln(12SAH^2K^2/\delta)}{n_h^{k-1} \vee 1}} + \frac{1}{K}, \end{aligned} \quad (\text{C.48})$$

for any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ , where (C.46) follows from (C.45), (C.47) is from (C.22) and a union bound, we set  $\epsilon = 1/2K$  in (C.48).

We denote  $V(\boldsymbol{\lambda}) = V_{h+1}^{k, \rho} - \boldsymbol{\lambda} \in [-H, H]$  and  $\mathcal{V} = \{V \in \mathbb{R}^S : \|V\|_{\infty} \leq H\}$ . To bound term (iii) in (C.43), we create a  $\epsilon$ -net  $\mathcal{N}_{\mathcal{V}}(\epsilon)$  for  $\mathcal{V}$ . From Lemma E.4, it holds that  $\ln |\mathcal{N}_{\mathcal{V}}(\epsilon)| \leq |S| \cdot \ln(3H/\epsilon)$ .

Therefore, by the definition of  $\mathcal{N}_{\mathcal{V}}(\epsilon)$ , for any fixed  $V$ , there exists a  $V' \in \mathcal{N}_{\mathcal{V}}(\epsilon)$  such that  $\|V - V'\|_{\infty} \leq \epsilon$ , that is

$$\begin{aligned} &|\text{Var}_{P_h^o}(V) - \text{Var}_{\hat{P}_h^k}(V)| \\ &\leq |\text{Var}_{P_h^o}(V) - \text{Var}_{P_h^o}(V')| + |\text{Var}_{P_h^o}(V') - \text{Var}_{\hat{P}_h^k}(V')| + |\text{Var}_{\hat{P}_h^k}(V') - \text{Var}_{\hat{P}_h^k}(V)| \\ &\leq |\mathbb{E}_{P_h^o}[V^2 - V'^2]| + |\mathbb{E}_{P_h^o}^2[V] - \mathbb{E}_{P_h^o}^2[V']| + |\mathbb{E}_{\hat{P}_h^k}[V'^2 - V^2]| + |\mathbb{E}_{\hat{P}_h^k}^2[V'] - \mathbb{E}_{\hat{P}_h^k}^2[V]| \\ &\quad + |\text{Var}_{P_h^o}(V') - \text{Var}_{\hat{P}_h^k}(V')| \\ &\leq |\text{Var}_{P_h^o}(V') - \text{Var}_{\hat{P}_h^k}(V')| + 8H\epsilon \\ &\leq \sup_{V' \in \mathcal{N}_{\mathcal{V}}} |\text{Var}_{P_h^o}(V') - \text{Var}_{\hat{P}_h^k}(V')| + 8H\epsilon. \end{aligned} \quad (\text{C.49})$$

For any fixed  $V$ , following the same analysis as (C.22), we have

$$\begin{aligned} |\text{Var}_{P_h^o}(V) - \text{Var}_{\hat{P}_h^k}(V)| &= |(\mathbb{E}_{P_h^o}[V^2] - \mathbb{E}_{P_h^o}^2[V]) - (\mathbb{E}_{\hat{P}_h^k}[V^2] - \mathbb{E}_{\hat{P}_h^k}^2[V])| \\ &\leq |\mathbb{E}_{P_h^o}[V^2] - \mathbb{E}_{\hat{P}_h^k}[V^2]| + |\mathbb{E}_{P_h^o}^2[V] - \mathbb{E}_{\hat{P}_h^k}^2[V]| \\ &\leq H^2 \sqrt{\frac{2S \ln(2/\delta)}{n_h^{k-1} \vee 1}} + (\mathbb{E}_{P_h^o}[V] + \mathbb{E}_{\hat{P}_h^k}[V]) \cdot |\mathbb{E}_{P_h^o}[V] - \mathbb{E}_{\hat{P}_h^k}[V]| \\ &\leq H^2 \sqrt{\frac{2S \ln(2/\delta)}{n_h^{k-1} \vee 1}} + 2H^2 \sqrt{\frac{2S \ln(2/\delta)}{n_h^{k-1} \vee 1}} \\ &\leq 3H^2 \sqrt{\frac{2S \ln(2/\delta)}{n_h^{k-1} \vee 1}} \end{aligned} \quad (\text{C.50})$$

with probability at least  $1 - \delta$ .



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

$$\begin{aligned} \sup_{\lambda \in [0, H]} |\text{Var}_{\hat{P}_h^k}(V_{h+1}^{k, \beta, \eta} - \lambda) - \text{Var}_{P_h^o}(V_{h+1}^{k, \beta, \eta} - \lambda)| &\leq \sup_{\lambda \in [0, H]} |\text{Var}_{P_h^o}(V(\lambda)) - \text{Var}_{\hat{P}_h^k}(V(\lambda))| \\ &\leq \sup_{V \in \mathcal{N}_V(\epsilon)} |\text{Var}_{P_h^o}(V) - \text{Var}_{\hat{P}_h^k}(V)| + 8H\epsilon \quad (\text{C.51}) \end{aligned}$$

$$\leq 3H^2 \sqrt{\frac{2S \ln(2SAHK|\mathcal{N}_V|/\delta)}{n_h^{k-1} \vee 1}} + 8H\epsilon \quad (\text{C.52})$$

$$\begin{aligned} &\leq 3H^2 \sqrt{\frac{2S^2 \ln(6SAH^2K/\epsilon\delta)}{n_h^{k-1} \vee 1}} + 8H\epsilon \\ &= 3H^2 \sqrt{\frac{2S^2 \ln(48SAH^3K^2/\delta)}{n_h^{k-1} \vee 1}} + \frac{1}{K}, \quad (\text{C.53}) \end{aligned}$$

for any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ , where (C.51) follows from (C.49), (C.52) is from (C.50) and a union bound, we set  $\epsilon = \frac{1}{8HK}$  in (C.53).

Apply the union bound again and combine (C.43) with (C.44), (C.48), (C.53), the definition of bonus and induction assumption. With probability at least  $1 - 3\delta$ , we have  $Q_h^{k, \rho}(s, a) \geq Q_h^{\pi, \rho}(s, a)$  for any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ . This completes the proof.  $\square$

**Lemma C.11.** If we set the bonus term to be the same as in Lemma C.10, then for any  $(k, s, a) \in [K] \times \mathcal{S} \times \mathcal{A}$ , with probability at least  $1 - \delta$ , the sum of estimation errors can be bounded as

$$Q_1^{k, \beta, \eta}(s, a) - Q_1^{\pi^k, \beta, \eta}(s, a) \leq 2 \cdot \mathbb{E}_{P^{w, k, \pi^k}}[\text{bonus}_h^k],$$

*Proof.* From the proof of Lemma C.10, we see that with probability at least  $1 - \delta$ , for any  $(k, s, a) \in [K] \times \mathcal{S} \times \mathcal{A}$ ,

$$\begin{aligned} &|\hat{r}_h^k(s, a) - r_h(s, a)| + \left| \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta\chi^2(P \|\hat{P}_h^k))(s, a) \right. \\ &\quad \left. - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta\chi^2(P \|P_h^o))(s, a) \right| \leq \text{bonus}_h^k(s, a). \quad (\text{C.54}) \end{aligned}$$

Recall that we define  $P_h^{w, k} = \argmin_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{\pi^k, \beta, \eta}] + \beta\chi^2(P \|P_h^o))$  as the worst-case transition in Definition 5.1, we have

$$\begin{aligned} &Q_h^{k, \beta, \eta} - Q_h^{\pi^k, \beta, \eta} \\ &\leq \text{bonus}_h^k + \hat{r}_h^k + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta\chi^2(P \|\hat{P}_h^k)) - r_h - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{\pi^k, \beta, \eta}] + \beta\chi^2(P \|P_h^o)) \\ &\leq 2\text{bonus}_h^k + \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{k, \beta, \eta}] + \beta\chi^2(P \|P_h^o)) - \inf_{P \in \Delta(S)} (\mathbb{E}_P[V_{h+1}^{\pi^k, \beta, \eta}] + \beta\chi^2(P \|P_h^o)) \quad (\text{C.55}) \end{aligned}$$

$$= 2\text{bonus}_h^k + \mathbb{E}_{P_h^{w, k, \pi^k}}[Q_{h+1}^{k, \beta, \eta} - Q_{h+1}^{\pi^k, \beta, \eta}]. \quad (\text{C.56})$$

where (C.55) uses (C.54). Apply (C.56) recursively, we can obtain the result.  $\square$

*Proof of Lemma C.1 in  $\chi^2$  Setting.* The result directly follows from combining Lemma C.10 and Lemma C.11, along with applying a union bound.  $\square$

## D Proofs of the Results in Linear Setting

### D.1 Proof of the Dynamic Programming Principle

*Proof of Proposition 6.2.* We prove a stronger version of the proposition. Specifically, there exists transition weights  $\{\tilde{\mu}_t\}_{t=1}^H$  together with  $\tilde{P}_t = \langle \phi, \tilde{\mu}_t \rangle$ , such that for all  $(h, s) \in [H] \times \mathcal{S}$ ,

$$V_h^{\pi, \beta, \eta}(s) = \mathbb{E}_{\{\tilde{P}_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \langle \phi(s_t, a_t), \mathbf{D}(\tilde{\mu}_t \| \mu_t^o) \rangle \right) \middle| s_h = s, \pi_{h:H} \right]. \quad (\text{D.1})$$

We prove this statement by induction. For the base case  $h = H$ , the claim holds trivially since no transition kernels are involved. For the inductive step, assume the conclusion holds for step  $h + 1$ , meaning that there exists  $\{\tilde{r}_t\}_{t=1}^H$  and  $\{\tilde{P}_t\}_{t=h+1}^H$  such that

$$V_{h+1}^{\pi, \beta, \eta}(s) = \mathbb{E}_{\{\tilde{P}_t\}_{t=h+1}^H} \left[ \sum_{t=h+1}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \langle \phi(s_t, a_t), \mathbf{D}(\tilde{\mu}_t \| \mu_t^o) \rangle \right) \middle| s_{h+1} = s, \pi_{h+1:H} \right]. \quad (\text{D.2})$$

We now prove (D.1) for the case of  $h$ . For  $Q_h^{\pi, \beta, \eta}(s, a)$ , on the one hand, we upper bound it as

$$\begin{aligned} & Q_h^{\pi, \beta, \eta}(s, a) \\ &= \inf_{\substack{\mu_t \in \Delta(S)^d \\ P_t = \langle \phi, \mu_t \rangle}} \mathbb{E}_{\{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) + \beta \langle \phi(s_t, a_t), \mathbf{D}(\mu_t \| \mu_t^o) \rangle \right) - \sum_{t=h+1}^H \Omega_{s_t}^\eta(\pi_t) \middle| s_h = s, a_h = a, \pi_{h+1:H} \right] \\ &= r_h(s, a) + \inf_{\substack{\mu_h \in \Delta(S)^d \\ P_h = \langle \phi, \mu_h \rangle}} \left\{ \beta \langle \phi(s_h, a_h), \mathbf{D}(\mu_h \| \mu_h^o) \rangle + \int_{\mathcal{S}} P_h(ds' | s, a) \mathbb{E}_{\{P_t\}_{t=h+1}^H} \left[ \right. \right. \\ &\quad \left. \left. \sum_{t=h+1}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \langle \phi(s_t, a_t), \mathbf{D}(\mu_t \| \mu_t^o) \rangle \right) \middle| s_{h+1} = s', \pi_{h+1:H} \right] \right\} \\ &\leq r_h(s, a) + \inf_{\substack{\mu_h \in \Delta(S)^d \\ P_h = \langle \phi, \mu_h \rangle}} \left\{ \beta \langle \phi(s_h, a_h), \mathbf{D}(\mu_h \| \mu_h^o) \rangle + \int_{\mathcal{S}} P_h(ds' | s, a) \mathbb{E}_{\{\tilde{P}_t\}_{t=h+1}^H} \left[ \right. \right. \\ &\quad \left. \left. \sum_{t=h+1}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \langle \phi(s_t, a_t), \mathbf{D}(\tilde{\mu}_t \| \mu_t^o) \rangle \right) \middle| s_{h+1} = s', \pi_{h+1:H} \right] \right\} \end{aligned} \quad (\text{D.3})$$

$$= r_h(s, a) + \inf_{\substack{\mu_h \in \Delta(S)^d \\ P_h = \langle \phi, \mu_h \rangle}} \left\{ \beta \langle \phi(s_h, a_h), \mathbf{D}(\mu_h \| \mu_h^o) \rangle + \mathbb{E}_{s' \sim P_h(\cdot | s, a)} [V_{h+1}^{\pi, \beta, \eta}(s')] \right\}, \quad (\text{D.4})$$

where (D.3) is because the definition of the infimum operator and (D.4) is from (D.2).

On the other hand, we can lower bound  $Q_h^{\pi, \beta, \eta}(s, a)$  as

$$\begin{aligned} & Q_h^{\pi, \beta, \eta}(s, a) \\ &= \inf_{\substack{\mu_t \in \Delta(S)^d \\ P_t = \langle \phi, \mu_t \rangle}} \mathbb{E}_{\{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) + \beta \langle \phi(s_t, a_t), \mathbf{D}(\mu_t \| \mu_t^o) \rangle \right) - \sum_{t=h+1}^H \Omega_{s_t}^\eta(\pi_t) \middle| s_h = s, a_h = a, \pi_{h+1:H} \right] \\ &= r_h(s, a) + \inf_{\substack{\mu_h \in \Delta(S)^d \\ P_h = \langle \phi, \mu_h \rangle}} \left\{ \beta \langle \phi(s_h, a_h), \mathbf{D}(\mu_h \| \mu_h^o) \rangle + \int_{\mathcal{S}} P_h(ds' | s, a) \mathbb{E}_{\{P_t\}_{t=h+1}^H} \left[ \right. \right. \end{aligned}$$

$$\begin{aligned}
& \sum_{t=h+1}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \langle \phi(s_t, a_t), \mathbf{D}(\mu_t \| \mu_t^o) \rangle \right) \Big|_{s_{h+1} = s', \pi_{h+1:H}} \Big] \Big\} \\
& \geq r_h(s, a) + \inf_{\substack{\mu_h \in \Delta(S)^d \\ P_h = \langle \phi, \mu_h \rangle}} \left\{ \beta \langle \phi(s_h, a_h), \mathbf{D}(\mu_h \| \mu_h^o) \rangle + \int_S P_h(ds' | s, a) \inf_{\substack{\mu_t \in \Delta(S)^d \\ P_t = \langle \phi, \mu_t \rangle}} \mathbb{E}_{\{P_t\}_{t=h+1}^H} \left[ \right. \right. \\
& \quad \left. \left. \sum_{t=h+1}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \langle \phi(s_t, a_t), \mathbf{D}(\mu_t \| \mu_t^o) \rangle \right) \Big|_{s_{h+1} = s', \pi_{h+1:H}} \right] \right\} \tag{D.5}
\end{aligned}$$

$$= r_h(s, a) + \inf_{\substack{\mu_h \in \Delta(S)^d \\ P_h = \langle \phi, \mu_h \rangle}} \left\{ \beta \langle \phi(s_h, a_h), \mathbf{D}(\mu_h \| \mu_h^o) \rangle + \mathbb{E}_{s' \sim P_h^o(\cdot | s, a)} [V_{h+1}^{\pi, \beta, \eta}(s')] \right\}, \tag{D.6}$$

where (D.5) is because  $\inf \int f \geq \int \inf f$  and (D.6) is from the definition of  $V_{h+1}^{\pi, \beta, \eta}(s')$ .

By combining (D.4) and (D.6), we have

$$Q_h^{\pi, \beta, \eta}(s, a) = r_h(s, a) + \inf_{\substack{\mu_h \in \Delta(S)^d \\ P_h = \langle \phi, \mu_h \rangle}} \left\{ \mathbb{E}_{s' \sim P_h^o(\cdot | s, a)} [V_{h+1}^{\pi, \beta, \eta}(s')] + \beta \langle \phi(s, a), \mathbf{D}(\mu_h \| \mu_h^o) \rangle \right\},$$

which implies (6.1).

Since  $\Delta(\mathcal{S})$  is compact, there exists  $\tilde{\mu}_h$  together with  $\tilde{P}_h = \langle \phi, \tilde{\mu}_h \rangle$ , such that

$$Q_h^{\pi, \beta, \eta}(s, a) = r_h(s, a) + \mathbb{E}_{s' \sim \tilde{P}_h(\cdot | s, a)} [V_{h+1}^{\pi, \beta, \eta}(s')] + \beta \langle \phi(s, a), \mathbf{D}(\tilde{\mu}_h \| \mu_h^o) \rangle. \tag{D.7}$$

This finishes the proof of (D.1).

For  $V_h^{\pi, \beta, \eta}(s)$ , on the one hand, we upper bound it as

$$\begin{aligned}
& V_h^{\pi, \beta, \eta}(s) \\
& = \inf_{\substack{\mu_t \in \Delta(S)^d \\ P_t = \langle \phi, \mu_t \rangle}} \mathbb{E}_{\{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \langle \phi(s_t, a_t), \mathbf{D}(\mu_t \| \mu_t^o) \rangle \right) \Big|_{s_h = s, \pi_{h:H}} \right] \\
& = \inf_{\substack{\mu_t \in \Delta(S)^d \\ P_t = \langle \phi, \mu_t \rangle}} \left\{ -\Omega_s^\eta(\pi_h) + \sum_{a \in \mathcal{A}} \pi_h(a | s) \mathbb{E}_{\{P_t\}_{t=h}^H} \left[ \right. \right. \\
& \quad \left. \left. \sum_{t=h}^H \left( r_t(s_t, a_t) + \beta \langle \phi(s_t, a_t), \mathbf{D}(\mu_t \| \mu_t^o) \rangle \right) - \sum_{t=h+1}^H \Omega_{s_t}^\eta(\pi_t) \Big|_{s_h = s, a_h = a, \pi_{h+1:H}} \right] \right\} \\
& \leq -\Omega_s^\eta(\pi_h) + \sum_{a \in \mathcal{A}} \pi_h(a | s) \mathbb{E}_{\{\tilde{P}_t\}_{t=h}^H} \left[ \right. \\
& \quad \left. \sum_{t=h}^H \left( r_t(s_t, a_t) + \beta \langle \phi(s_t, a_t), \mathbf{D}(\tilde{\mu}_t \| \mu_t^o) \rangle \right) - \sum_{t=h+1}^H \Omega_{s_t}^\eta(\pi_t) \Big|_{s_h = s, a_h = a, \pi_{h+1:H}} \right] \tag{D.8}
\end{aligned}$$

$$= -\Omega_s^\eta(\pi_h) + \sum_{a \in \mathcal{A}} \pi_h(a | s) Q_h^{\pi, \beta, \eta}(s, a), \tag{D.9}$$

where (D.8) is because the definition of the infimum operator and (D.9) is from (D.7) and (D.2).

On the other hand, we can lower bound  $V_h^{\pi, \beta, \eta}(s)$  as

$$V_h^{\pi, \beta, \eta}(s)$$

$$\begin{aligned}
&= \inf_{\substack{\boldsymbol{\mu}_t \in \Delta(S)^d \\ P_t = \langle \boldsymbol{\phi}, \boldsymbol{\mu}_t \rangle}} \mathbb{E}_{\{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) - \Omega_{s_t}^\eta(\pi_t) + \beta \langle \boldsymbol{\phi}(s_t, a_t), \mathbf{D}(\boldsymbol{\mu}_t \| \boldsymbol{\mu}_t^o) \rangle \right) \middle| s_h = s, \pi_{h:H} \right] \\
&= \inf_{\substack{\boldsymbol{\mu}_t \in \Delta(S)^d \\ P_t = \langle \boldsymbol{\phi}, \boldsymbol{\mu}_t \rangle}} \left\{ -\Omega_s^\eta(\pi_h) + \sum_{a \in \mathcal{A}} \pi_h(a|s) \mathbb{E}_{\{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) + \beta \langle \boldsymbol{\phi}(s_t, a_t), \mathbf{D}(\boldsymbol{\mu}_t \| \boldsymbol{\mu}_t^o) \rangle \right) - \sum_{t=h+1}^H \Omega_{s_t}^\eta(\pi_t) \middle| s_h = s, a_h = a, \pi_{h+1:H} \right] \right\} \\
&\geq -\Omega_s^\eta(\pi_h) + \sum_{a \in \mathcal{A}} \pi_h(a|s) \inf_{\substack{\boldsymbol{\mu}_t \in \Delta(S)^d \\ P_t = \langle \boldsymbol{\phi}, \boldsymbol{\mu}_t \rangle}} \mathbb{E}_{\{P_t\}_{t=h}^H} \left[ \sum_{t=h}^H \left( r_t(s_t, a_t) + \beta \langle \boldsymbol{\phi}(s_t, a_t), \mathbf{D}(\boldsymbol{\mu}_t \| \boldsymbol{\mu}_t^o) \rangle \right) - \sum_{t=h+1}^H \Omega_{s_t}^\eta(\pi_t) \middle| s_h = s, a_h = a, \pi_{h+1:H} \right] \quad (\text{D.10}) \\
&= -\Omega_s^\eta(\pi_h) + \sum_{a \in \mathcal{A}} \pi_h(a|s) Q_h^{\pi, \beta, \eta}(s, a), \quad (\text{D.11})
\end{aligned}$$

where (D.10) is because  $\inf \sum f \geq \sum \inf f$  and (D.11) is from the definition of  $Q_h^{\pi, \beta, \eta}(s, a)$ .

By combining (D.9) and (D.11), we have

$$V_h^{\pi, \beta, \eta}(s) = -\Omega_s^\eta(\pi_h) + \sum_{a \in \mathcal{A}} \pi_h(a|s) Q_h^{\pi, \beta, \eta}(s, a),$$

which implies (6.2).  $\square$

## D.2 Proofs of the Main Results

To establish Theorems 6.6 and 6.7, we first introduce Lemma D.1, which bounds the regret by the sum of bonus terms under the expectation defined by the worst-case transition.

**Lemma D.1.** Under Assumption 6.5 and Assumption 6.1, for each  $f$ -divergence setting and any  $\delta \in (0, 1/2)$ , if we choose  $\lambda = 1$  and set  $c$  according to Theorem 6.6 in Algorithm 2, then with probability at least  $1 - 2\delta$ , the regret of Algorithm 2 satisfies

$$\text{Regret}(K) \leq 2c \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^d \mathbb{E}_{P^{w, k}, \pi^k} \left[ \|\phi_i(s, a) \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \right]. \quad (\text{D.12})$$

*Proof of Theorem 6.6.* Recall that

$$\Gamma_h^k(s, a) = c \sum_{i=1}^d \|\phi_i(s, a) \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}}$$

is the bonus assigned to the state-action pair  $(s, a)$  at episode  $k$  and step  $h$ . Due to the truncation in (6.3), we may assume without loss of generality that  $\Gamma_h^k(s, a) \leq H$  for all  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ . Since otherwise we can simply truncate  $\Gamma_h^k(s, a)$  at  $H$ , which does not affect the results.

Therefore, by Azuma-Hoeffding inequality, with probability at least  $1 - \delta$ , we have

$$\left| \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{P^o, \pi^k} [\Gamma_h^k(s, a)] - \sum_{k=1}^K \sum_{h=1}^H \Gamma_h^k(s_h^k, a_h^k) \right| \leq \sqrt{2KH^3 \ln(2/\delta)}.$$

Now we further analyze (D.12) as

$$\begin{aligned}
\text{Regret}(K) &\leq 2 \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{P^{w,k}, \pi^k} [\Gamma_h^k(s, a)] \\
&\leq 2C_{vr} \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{P^o, \pi^k} [\Gamma_h^k(s, a)] \\
&\leq 2C_{vr} \left( \sqrt{2KH^3 \ln(2/\delta)} + \sum_{k=1}^K \sum_{h=1}^H \Gamma_h^k(s_h^k, a_h^k) \right) \\
&= 2C_{vr} \left( \sqrt{2KH^3 \ln(2/\delta)} + c \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^d \|\phi_i(s_h^k, a_h^k) \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \right),
\end{aligned} \tag{D.13}$$

where (D.13) follows from [Assumption 6.5](#). This finishes the proof.  $\square$

*Proof of Theorem 6.7.* Now we further analyze (D.12) as

$$\begin{aligned}
\text{Regret}(K) &\leq 2c \sum_{k=1}^K \sum_{h=1}^H \sum_{i=1}^d \mathbb{E}_{P^{w,k}, \pi^k} \left[ \|\phi_i(s, a) \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \right] \\
&\leq 2c \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{P^{w,k}, \pi^k} \left[ \sum_{i=1}^d \phi_i(s, a) \sqrt{\lambda_{\max}((\Lambda_h^k)^{-1})} \right] \\
&= 2c \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{P^{w,k}, \pi^k} \left[ \sqrt{\lambda_{\max}((\Lambda_h^k)^{-1})} \right]
\end{aligned} \tag{D.14}$$

$$\begin{aligned}
&= 2c \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{P^{w,k}, \pi^k} \left[ \sqrt{\lambda_{\max}((\Lambda_h^k)^{-1})} \right] \\
&\leq 2c \sum_{k=1}^K \sum_{h=1}^H \sqrt{\mathbb{E}_{P^{w,k}, \pi^k} [\lambda_{\max}((\Lambda_h^k)^{-1})]}
\end{aligned} \tag{D.15}$$

$$\leq 2c\sqrt{K} \sum_{h=1}^H \sqrt{\sum_{k=1}^K \mathbb{E}_{P^{w,k}, \pi^k} [\lambda_{\max}((\Lambda_h^k)^{-1})]}, \tag{D.16}$$

where (D.14) is because  $\|x\|_A \leq \sqrt{\lambda_{\max}(A)} \|x\|_2$ , (D.15) is from Jensen's inequality, (D.16) is from Cauchy-Schwarz inequality.

Using [Lemma E.8](#), with probability at least  $1 - \delta$ , we have

$$\begin{aligned}
&\sum_{k=1}^K \mathbb{E}_{P^{w,k}, \pi^k} [\lambda_{\max}((\Lambda_h^k)^{-1})] \\
&= \sum_{k=1}^K \mathbb{E}_{P^{w,k}, \pi^k} \left[ \frac{1}{\lambda_{\min}(\Lambda_h^k)} \right] \\
&\leq \sum_{k=1}^K \mathbb{E}_{P^{w,k}, \pi^k} \left[ \frac{1}{\max \{ \alpha(k-1) + \lambda - \sqrt{32k \log(dKH/\delta)}, \lambda \}} \right] \\
&= \sum_{k=1}^K \frac{1}{\max \{ \alpha(k-1) + \lambda - \sqrt{32k \log(dKH/\delta)}, \lambda \}}
\end{aligned}$$

$$\leq \frac{128}{\alpha^2} \log \left( \frac{dHK}{\delta} \right) + \frac{2}{\alpha} \log(K). \quad (\text{D.17})$$

Combining (D.16) and (D.17), we finish the proof.  $\square$

### D.3 Proofs of the Technical Lemmas

Before proving [Lemma D.1](#) in different settings, we state a technical lemma about the covering number.

**Lemma D.2.** For any  $h \in [H]$ , let  $\mathcal{V}_h$  denote a class of functions mapping from  $\mathcal{S}$  to  $\mathbb{R}$  with the following form

$$V_h(s; w, \beta, \Lambda) = \left[ \max_{\pi \in \Pi} \left\{ \left\langle \pi(s, a), \phi(s, a)^\top \mathbf{w} + \beta \sum_{i=1}^d \|\phi_i(s, a) \mathbf{1}_i\|_{\Lambda^{-1}} \right\rangle_{\mathcal{A}} - \sigma_{\mathcal{R}_h}(-\pi_h(\cdot|s)) \right\} \right]_{[0, H-h+1]},$$

where the parameters  $(w, \beta, \Lambda)$  satisfy  $\|w\| \leq L$ ,  $\beta \in [0, B]$ ,  $\lambda_{\min}(\Lambda) \geq \lambda$ . Let  $\mathcal{N}_h(\epsilon; \mathcal{V}_h)$  be the  $\epsilon$ -covering number of  $\mathcal{V}_h$  with respect to the distance  $\text{dist}(V_1, V_2) = \sup_{s \in \mathcal{S}} |V_1(s) - V_2(s)|$ . Then

$$\log \mathcal{N}_h(\epsilon) \leq d \log(1 + 4L|\mathcal{A}|/\epsilon) + d^2 \log(1 + 8d^{1/2}B^2|\mathcal{A}|^2/(\lambda\epsilon^2)).$$

*Proof.* The argument is similar to that of [Liu and Xu \(2024a, Lemma D.3\)](#). Denoting  $A = \beta^2 \Lambda^{-1}$ , we have

$$V_h(s; w, \beta, \Lambda) = \left[ \max_{\pi \in \Pi} \left\{ \left\langle \pi(s, a), \phi(s, a)^\top \mathbf{w} + \beta \sum_{i=1}^d \|\phi_i(s, a) \mathbf{1}_i\|_{\Lambda^{-1}} \right\rangle_{\mathcal{A}} - \sigma_{\mathcal{R}_h}(-\pi_h(\cdot|s)) \right\} \right]_{[0, H-h+1]},$$

for  $\|w\| \leq L$  and  $\|A\| \leq B^2 \lambda^{-1}$ . For any two functions  $V_1, V_2 \in \mathcal{V}_h$ , let them take the form above with parameters  $(\mathbf{w}_1, A_1)$  and  $(\mathbf{w}_2, A_2)$  respectively. Since  $\min\{\cdot, H-h+1\}$  is a concentration map, we have

$$\begin{aligned} & \text{dist}(V_1, V_2) \\ & \leq \sup_{s, \pi} \left| \left\{ \left\langle \pi(s, a), \phi(s, a)^\top \mathbf{w}_1 + \sum_{i=1}^d \|\phi_i(s, a) \mathbf{1}_i\|_{A_1} \right\rangle_{\mathcal{A}} - \sigma_{\mathcal{R}_h}(-\pi_h(\cdot|s)) \right\} \right. \\ & \quad \left. - \left\{ \left\langle \pi(s, a), \phi(s, a)^\top \mathbf{w}_2 + \sum_{i=1}^d \|\phi_i(s, a) \mathbf{1}_i\|_{A_2} \right\rangle_{\mathcal{A}} - \sigma_{\mathcal{R}_h}(-\pi_h(\cdot|s)) \right\} \right| \\ & \leq \sup_{s, \pi} \left| \left\langle \pi(s, a), \left( \phi(s, a)^\top \mathbf{w}_1 + \sum_{i=1}^d \|\phi_i(s, a) \mathbf{1}_i\|_{A_1} \right) - \left( \phi(s, a)^\top \mathbf{w}_2 + \sum_{i=1}^d \|\phi_i(s, a) \mathbf{1}_i\|_{A_2} \right) \right\rangle_{\mathcal{A}} \right| \\ & \leq \sup_{s, \pi} \left| \langle \pi(s, a), \phi(s, a)^\top (\mathbf{w}_1 - \mathbf{w}_2) \rangle \right| + \sup_{s, \pi} \left| \left\langle \pi(s, a), \sum_{i=1}^d \|\phi_i(s, a) \mathbf{1}_i\|_{A_1} - \sum_{i=1}^d \|\phi_i(s, a) \mathbf{1}_i\|_{A_2} \right\rangle_{\mathcal{A}} \right| \\ & \leq \sup_{s, \pi} \left| \langle \pi(s, a), \phi(s, a)^\top (\mathbf{w}_1 - \mathbf{w}_2) \rangle \right| + \sup_{s, \pi} \left| \left\langle \pi(s, a), \sum_{i=1}^d \|\phi_i(s, a) \mathbf{1}_i\|_{A_1 - A_2} \right\rangle_{\mathcal{A}} \right| \end{aligned} \quad (\text{D.18})$$

$$\leq |\mathcal{A}| \left( \|\mathbf{w}_1 - \mathbf{w}_2\| + \sqrt{\|A_1 - A_2\|_F} \right),$$

where (D.18) is because  $|\sqrt{x} - \sqrt{y}| \leq \sqrt{|x - y|}$  for  $x, y \geq 0$ . For matrices,  $\|\cdot\|$  and  $\|\cdot\|_F$  denote the matrix operator norm and Frobenius norm respectively.

Let  $\mathcal{C}_{\mathbf{w}}$  be an  $\epsilon/(2|\mathcal{A}|)$ -cover of  $\{\mathbf{w} \in \mathbb{R}^d \mid \|\mathbf{w}\|_2 \leq L\}$  with respect to the 2-norm and  $\mathcal{C}_A$  be an  $\epsilon^2/(4|\mathcal{A}|^2)$ -cover of  $\{A \in \mathbb{R}^{d \times d} \mid \|A\|_F \leq d^{1/2} B^2 \lambda^{-1}\}$  with respect to the Frobenius norm. By Lemma E.6, we have

$$\log \mathcal{N}_\epsilon \leq \log |\mathcal{C}_{\mathbf{w}}| + \log |\mathcal{C}_A| \leq d \log(1 + 4L|\mathcal{A}|/\epsilon) + d^2 \log(1 + 8d^{1/2} B^2 |\mathcal{A}|^2 / (\lambda \epsilon^2)).$$

This finishes the proof.  $\square$

Next, we present the proof of Lemma D.1 in various  $f$ -divergence settings. For convenience, we also write  $P^{w,k} := P^{w, \pi^k}$ ,  $d^k := d^{\pi^k}$  and  $q^k := q^{\pi^k}$ .

### D.3.1 Proof of Lemma D.1 in TV Setting

**Lemma D.3** (Dual formulation). (Tang et al., 2025, Proposition 4.2) For the optimization problem  $\inf_{\mu \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu} [V_{h+1}] + \beta \text{TV}(\mu \| \mu^o))$ , we have its dual formulation as follows

$$\inf_{\mu \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu} [V_{h+1}] + \beta \text{TV}(\mu \| \mu^o)) = \mathbb{E}_{s \sim \mu^o} [V_{h+1}(s)] \min_{s' \in \mathcal{S}} V_{h+1}(s') + \beta.$$

**Lemma D.4** (Optimism). If we set the bonus term as follows

$$\Gamma_h^k(s, a) = c_{\text{TV}} \sum_{i=1}^d \|\phi_i(s, a) \mathbb{1}_i\|_{(\Lambda_h^k)^{-1}},$$

$$\text{where } c_{\text{TV}} = Hd \cdot \xi_{\text{TV}},$$

$$\text{where } \xi_{\text{TV}} = 720 + 3\sqrt{40 \log(96K^{13/2} H |\mathcal{A}|^3 / \delta)},$$

then for any policy  $\pi$  and any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ , with probability at least  $1 - \delta$ , we have  $Q_h^{k, \beta, \eta}(s, a) \geq Q_h^{\pi, \beta, \eta}(s, a)$ . Specially, by setting  $\pi = \pi^*$ , we have  $Q_h^{k, \beta, \eta}(s, a) \geq Q_h^{\pi^*, \beta, \eta}(s, a)$ .

*Proof.* We prove this by induction. First,  $h = H + 1$  holds trivially since  $Q_{H+1}^{k, \beta, \eta}(s, a) = 0 = Q_{H+1}^{\pi, \beta, \eta}(s, a)$ .

Assume  $Q_{h+1}^{k, \beta, \eta}(s, a) \geq Q_{h+1}^{\pi, \beta, \eta}(s, a)$  holds, due to the choice of  $\pi_{h+1}^k$  in (4.3), we have

$$\begin{aligned} V_{h+1}^{k, \beta, \eta}(s) &= \langle Q_{h+1}^{k, \beta, \eta}(s, \cdot), \pi_{h+1}^k(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}^k) \\ &\geq \langle Q_{h+1}^{k, \beta, \eta}(s, \cdot), \pi_{h+1}(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}) \\ &\geq \langle Q_{h+1}^{\pi, \beta, \eta}(s, \cdot), \pi_{h+1}(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}) = V_{h+1}^{\pi, \beta, \eta}(s). \end{aligned}$$

Recall that we denote  $Q_h^{k, \beta, \eta}$  as the optimistic estimation in  $k$ -th episode, that is,

$$Q_h^{k, \beta, \eta}(s, a) \leftarrow \min \{ \phi(s, a)^\top (\theta_h^k + \mathbf{w}_h^k) + \Gamma_h^k(s, a), H - h + 1 \}. \quad (\text{D.19})$$



If  $Q_h^{k,\beta,\eta}(s, a) = H - h + 1$ , then it follows immediately that

$$Q_h^{k,\beta,\eta}(s, a) = H - h + 1 \geq Q_h^{\pi,\beta,\eta}(s, a)$$

by the definition of  $Q_h^{\pi,\beta,\eta}(s, a)$ . Otherwise, we can infer that

$$\begin{aligned} & Q_h^{k,\beta,\eta} - Q_h^{\pi,\beta,\eta} \\ &= \Gamma_h^k + \phi^\top (\theta_h^k + w_h^k) - \phi^\top \left[ \theta_h + \inf_{\mu \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu} [V_{h+1}^{\pi,\beta,\eta}] + \beta \text{TV}(\mu \| \mu^o)) \right] \\ &\geq \Gamma_h^k + \phi^\top (\theta_h^k + w_h^k) - \phi^\top \left[ \theta_h + \inf_{\mu \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu} [V_{h+1}^{k,\beta,\eta}] + \beta \text{TV}(\mu \| \mu^o)) \right] \end{aligned} \quad (\text{D.20})$$

$$= \Gamma_h^k + \left\langle \phi, \theta_h^k + \widehat{\mathbb{E}}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s') + \beta} - \theta_h - \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s') + \beta} \right\rangle \quad (\text{D.21})$$

$$\begin{aligned} &\geq \Gamma_h^k - \left\langle \phi, |\theta_h^k - \theta_h| + \left| \widehat{\mathbb{E}}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s') + \beta} - \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s') + \beta} \right| \right\rangle \\ &= \Gamma_h^k - \underbrace{\sum_{i=1}^d \phi_i \mathbb{1}_i^\top |\theta_h^k - \theta_h|}_{(i)} - \underbrace{\sum_{i=1}^d \phi_i \mathbb{1}_i^\top \left| \widehat{\mathbb{E}}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s') + \beta} - \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s') + \beta} \right|}_{(ii)}, \end{aligned} \quad (\text{D.22})$$

where (D.20) follows from the induction assumption, and (D.21) is obtained from the algorithm update formulation and the dual formulation.

We introduce the following notation for the sake of convenience

$$\text{err}_h^\tau(f) = \mathbb{E}_{s' \sim \mathbb{P}_h^o(\cdot | s_h^\tau, a_h^\tau)} [f(s')] - f(s_{h+1}^\tau). \quad (\text{D.23})$$

For term (i) in (D.22), we have

$$\begin{aligned} & |\phi_i \mathbb{1}_i^\top (\theta_h^k - \theta_h)| \\ &= \left| \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot r_h^\tau \right) \right. \\ &\quad \left. - \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top + \lambda \cdot I \right) \cdot \theta_h \right| \\ &= \lambda |\phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \cdot \theta_h| \end{aligned} \quad (\text{D.24})$$

$$\leq \lambda \|\phi_i \mathbb{1}_i\|_{(\Lambda_h^k)^{-1}} \cdot \|\theta_h\|_{(\Lambda_h^k)^{-1}} \quad (\text{D.25})$$

$$\leq \lambda \|\phi_i \mathbb{1}_i\|_{(\Lambda_h^k)^{-1}} \cdot \lambda_{\max}((\Lambda_h^k)^{-1})^{\frac{1}{2}} \|\theta_h\|_2 \quad (\text{D.26})$$

$$\leq \lambda^{\frac{1}{2}} \sqrt{d} \cdot \|\phi_i \mathbb{1}_i\|_{(\Lambda_h^k)^{-1}},$$

where (D.24) is obtained from the algorithm update formulation and the definition of  $\Lambda_h^k$ , (D.25) follows from the Cauchy-Schwartz inequality, (D.26) is because  $\|x\|_A \leq \sqrt{\lambda_{\max}(A)} \|x\|_2$ .

For term (ii) in (D.22), we decompose it as follows.

$$\left| \phi_i \mathbb{1}_i^\top \left( \widehat{\mathbb{E}}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s') + \beta} - \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s') + \beta} \right) \right|$$

$$\begin{aligned}
&= \left| \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot [V_{h+1}^{k,\beta,\eta}(s_{h+1}^\tau)]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s')+\beta} \right) \right. \\
&\quad \left. - \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top + \lambda \cdot I \right) \left( \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}(s)]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s')+\beta} \right) \right| \quad (\text{D.27})
\end{aligned}$$

$$\begin{aligned}
&= \left| \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau \left( [V_{h+1}^{k,\beta,\eta}(s)]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s')+\beta} \right) \right) \right. \\
&\quad \left. + \lambda \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \left( \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}(s)]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s')+\beta} \right) \right| \quad (\text{D.28})
\end{aligned}$$

$$\begin{aligned}
&\leq \left| \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau \left( [V_{h+1}^{k,\beta,\eta}(s)]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s')+\beta} \right) \right) \right| \\
&\quad + \lambda \left| \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \left( \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}(s)]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s')+\beta} \right) \right| \\
&\leq \underbrace{\|\phi_i \mathbb{1}_i\|_{(\Lambda_h^k)^{-1}} \cdot \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau \left( [V_{h+1}^{k,\beta,\eta}(s)]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s')+\beta} \right) \right\|_{(\Lambda_h^k)^{-1}}}_{(\text{iii})} \\
&\quad + \underbrace{\lambda \|\phi_i \mathbb{1}_i\|_{(\Lambda_h^k)^{-1}} \cdot \left\| \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}(s)]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s')+\beta} \right\|_{(\Lambda_h^k)^{-1}}}_{(\text{iv})}, \quad (\text{D.29})
\end{aligned}$$

where (D.27) is obtained from the algorithm update formulation and the definition of  $\Lambda_h^k$ , (D.28) is from the definition in (D.23), (D.29) follows from the Cauchy-Schwartz inequality.

For term (iv) in (D.29), since  $V_{h+1}^{k,\beta,\eta} \leq H$ , we have

$$\begin{aligned}
&\left\| \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s')+\beta} \right\|_{(\Lambda_h^k)^{-1}} \\
&\leq \lambda_{\max}((\Lambda_h^k)^{-1})^{\frac{1}{2}} \left\| \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s')+\beta} \right\|_2 \\
&\leq \lambda^{-\frac{1}{2}} H \sqrt{d}, \quad (\text{D.30})
\end{aligned}$$

where (D.30) is because  $\|x\|_A \leq \sqrt{\lambda_{\max}(A)} \|x\|_2$ .

Before we continue to bound term (iii), we prove a auxiliary result first. That is,  $\|\mathbf{w}_h^k\|_2 \leq H \sqrt{kd/\lambda}$ . From the algorithm update formulation, we have

$$\begin{aligned}
\|\mathbf{w}_h^k\|_2 &= \left\| (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot [V_{h+1}^{k,\beta,\eta}(s_{h+1}^\tau)]_{\min_{s'} V_{h+1}^{k,\beta,\eta}(s')+\beta} \right) \right\|_2 \\
&\leq H \sum_{\tau=1}^{k-1} \left\| (\Lambda_h^k)^{-1} \phi(s_h^\tau, a_h^\tau) \right\|_2 \\
&= H \sum_{\tau=1}^{k-1} \sqrt{\phi(s_h^\tau, a_h^\tau)^\top (\Lambda_h^k)^{-1/2} (\Lambda_h^k)^{-1} (\Lambda_h^k)^{-1/2} \phi(s_h^\tau, a_h^\tau)}
\end{aligned}$$

$$\begin{aligned}
&\leq \frac{H}{\sqrt{\lambda}} \sum_{\tau=1}^{k-1} \sqrt{\phi(s_h^\tau, a_h^\tau)^\top (\Lambda_h^k)^{-1} \phi(s_h^\tau, a_h^\tau)} \\
&\leq \frac{H\sqrt{k}}{\sqrt{\lambda}} \sqrt{\sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau)^\top (\Lambda_h^k)^{-1} \phi(s_h^\tau, a_h^\tau)} \\
&= \frac{H\sqrt{k}}{\sqrt{\lambda}} \sqrt{\text{tr}((\Lambda_h^k)^{-1} (\Lambda_h^k - \lambda I))} \\
&\leq \frac{H\sqrt{k}}{\sqrt{\lambda}} \sqrt{\text{tr}(I)} = H\sqrt{kd/\lambda}.
\end{aligned}$$

Now we are ready to bound term (iii) in (D.29). Let  $\mathcal{V}_h$  denote a class of functions mapping from  $\mathcal{S}$  to  $\mathbb{R}$  with the following form

$$V(\cdot) = \min \left\{ \max_{a \in \mathcal{A}} \left\{ \phi(\cdot, a)^\top \mathbf{w} + \beta \sum_{i=1}^d \|\phi_i(\cdot, a)\|_{\Lambda^{-1}} \right\}, H - h + 1 \right\},$$

where the parameters  $(w, \beta, \Lambda)$  satisfy  $\|w\| \leq L$ ,  $\beta \in [0, B]$ ,  $\lambda_{\min} \geq \lambda$ . Here we set  $L = H\sqrt{kd/\lambda}$  and  $B = Hd \cdot \xi_{\text{TV}}$ , where  $\xi_{\text{TV}} = 720 + 3\sqrt{40 \log(96K^{13/2}H|\mathcal{A}|^3/\delta)}$ .

Let  $\mathcal{N}_h(\epsilon; \mathcal{V}_h)$  be the minimum  $\epsilon$ -cover of  $\mathcal{V}_h$ ,  $\mathcal{N}_h(\epsilon; [0, H])$  be the minimum  $\epsilon$ -cover of  $[0, H]$ . Therefore, for any  $V \in \mathcal{V}_h$  and  $\alpha \in [0, H]$ , there exists a function  $V_\epsilon \in \mathcal{N}_h(\epsilon; \mathcal{V}_h)$  and  $\alpha_\epsilon \in \mathcal{N}_h(\epsilon; [0, H])$  such that

$$\sup_{s \in \mathcal{S}} |V(s) - V_\epsilon(s)| \leq \epsilon, \quad \left| \min_{s'} V_{h+1}^{k, \beta, \eta}(s') + \beta - \alpha_\epsilon \right| \leq \epsilon.$$

Then we have

$$\begin{aligned}
(\text{iii})^2 &= \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau \left( [V_{h+1}^{k, \beta, \eta}]_{\min_{s'} V_{h+1}^{k, \beta, \eta}(s') + \beta} \right) \right\|_{(\Lambda_h^k)^{-1}}^2 \\
&\leq 2 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau ([V_{h+1}^{k, \beta, \eta}]_{\alpha_\epsilon}) \right\|_{(\Lambda_h^k)^{-1}}^2 \\
&\quad + 2 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau \left( [V_{h+1}^{k, \beta, \eta}]_{\min_{s'} V_{h+1}^{k, \beta, \eta}(s') + \beta} - [V_{h+1}^{k, \beta, \eta}]_{\alpha_\epsilon} \right) \right\|_{(\Lambda_h^k)^{-1}}^2 \tag{D.31}
\end{aligned}$$

$$\begin{aligned}
&\leq 4 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau ([V_\epsilon]_{\alpha_\epsilon}) \right\|_{(\Lambda_h^k)^{-1}}^2 + 4 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau ([V_{h+1}^{k, \beta, \eta}]_{\alpha_\epsilon} - [V_\epsilon]_{\alpha_\epsilon}) \right\|_{(\Lambda_h^k)^{-1}}^2 \\
&\quad + 2(2\epsilon)^2 k^2 / \lambda \tag{D.32}
\end{aligned}$$

$$\begin{aligned}
&\leq 4 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau ([V_\epsilon]_{\alpha_\epsilon}) \right\|_{(\Lambda_h^k)^{-1}}^2 + 4(2\epsilon)^2 k^2 / \lambda + 2(2\epsilon)^2 k^2 / \lambda \\
&= 4 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau ([V_\epsilon]_{\alpha_\epsilon}) \right\|_{(\Lambda_h^k)^{-1}}^2 + 24\epsilon^2 k^2 / \lambda
\end{aligned}$$

$$\leq 4H^2 \left( d \log(1 + k/\lambda) + 2 \log \frac{KH |\mathcal{N}_h(\epsilon; \mathcal{V}_h)| \cdot |\mathcal{N}_h(\epsilon; [0, H])|}{\delta} \right) + 24\epsilon^2 k^2 / \lambda \quad (\text{D.33})$$

$$\leq 4H^2 \left( d \log(1 + k/\lambda) + 2d \log(1 + 4L|\mathcal{A}|/\epsilon) + 2d^2 \log(1 + 8d^{1/2}B^2|\mathcal{A}|^2/(\lambda\epsilon^2)) \right. \\ \left. + 2 \log(3H/\epsilon) + 2 \log(KH/\delta) \right) + 24\epsilon^2 k^2 / \lambda \quad (\text{D.34})$$

$$\leq 4H^2 \left( 2d \log(k) + 4d \log(4k^{3/2}|\mathcal{A}|/(d^{1/2})) + 4d^2 \log(8k^2B^2|\mathcal{A}|^2/(H^2d^{3/2})) \right. \\ \left. + 2 \log(3k/d) + 2 \log(KH/\delta) \right) + 24H^2d^2 \quad (\text{D.35}) \\ \leq 16H^2d^2 \log(96K^{13/2}B^2|\mathcal{A}|^3/(Hd^3\delta)) + 24H^2d^2,$$

where (D.31) and (D.32) is because  $\|a + b\|_A^2 \leq 2\|a\|_A^2 + 2\|b\|_A^2$ , (D.33) is from Lemma E.7 together with a union bound over all  $(k, h) \in [K] \times [H]$ , (D.34) makes use of Lemmas D.2 and E.5, we set  $\epsilon = Hd/k$  and  $\lambda = 1$  and apply the inequality  $\log(1 + x) \leq 2 \log x$  for  $x \geq 2$  in (D.35).

Combining everything together, recall that we set  $\lambda = 1$ , we have

$$(i) + (ii) \leq \left( \sqrt{d} + \sqrt{16H^2d^2 \log(96K^{13/2}B^2|\mathcal{A}|^3/(Hd^3\delta)) + 24H^2d^2} + H\sqrt{d} \right) \cdot \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \\ \leq 3Hd \sqrt{40 \log(96K^{13/2}B^2|\mathcal{A}|^3/(Hd^3\delta))} \cdot \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}}.$$

With our choice of  $B$ , and noting that  $x \geq \max\{(2B + \sqrt{A})^2, 1\}$  guarantees  $x \geq A + B \log x$ , we conclude the proof.  $\square$

**Lemma D.5.** If we set the bonus term to be the same as in Lemma D.4, then for any  $(k, s, a) \in [K] \times \mathcal{S} \times \mathcal{A}$ , with probability at least  $1 - \delta$ , the sum of estimation errors can be bounded as

$$Q_1^{k,\beta,\eta}(s, a) - Q_1^{\pi^k,\beta,\eta}(s, a) \leq 2 \cdot \mathbb{E}_{P^{w,k}, \pi^k} \left[ \sum_{h=1}^H \Gamma_h^k(s_h^k, a_h^k) \right],$$

where  $P_h^{w,k}$  is defined by  $P_h^{w,k}(s, a) = \phi(s, a)^\top \cdot \mu_h^{w,k}$ .

*Proof.* From the definition, We have that

$$Q_h^{k,\beta,\eta} - Q_h^{\pi^k,\beta,\eta} \\ = \Gamma_h^k + \phi^\top(\theta_h^k + w_h^k) - \phi^\top \left[ \theta_h + \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{\pi^k,\beta,\eta}] + \beta \text{TV}(\mu_h \| \mu_h^o)) \right] \\ = \Gamma_h^k + \phi^\top(\theta_h^k - \theta_h) + \phi^\top \left[ w_h^k - \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{k,\beta,\eta}] + \beta \text{TV}(\mu_h \| \mu_h^o)) \right] \\ + \phi^\top \left[ \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{k,\beta,\eta}] + \beta \text{TV}(\mu_h \| \mu_h^o)) - \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{\pi^k,\beta,\eta}] + \beta \text{TV}(\mu_h \| \mu_h^o)) \right] \\ \leq 2 \cdot \Gamma_h^k + \phi^\top \left[ \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{k,\beta,\eta}] + \beta \text{TV}(\mu_h \| \mu_h^o)) - \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{\pi^k,\beta,\eta}] + \beta \text{TV}(\mu_h \| \mu_h^o)) \right] \quad (\text{D.36})$$

$$\leq 2 \cdot \Gamma_h^k + \phi^\top \mathbb{E}_{s \sim \mu_h^{w,k}} [V_{h+1}^{k,\beta,\eta} - V_{h+1}^{\pi^k,\beta,\eta}] \quad (\text{D.37}) \\ \leq 2 \cdot \Gamma_h^k + \mathbb{E}_{s \sim P_h^w} [V_{h+1}^{k,\beta,\eta} - V_{h+1}^{\pi^k,\beta,\eta}],$$

where (D.36) can be derived from Lemma D.4, and (D.37) uses the definition of the worst-case transition. Apply the above inequality recursively, we have

$$Q_1^{k,\beta,\eta} - Q_1^{\pi^k,\beta,\eta} \leq 2 \cdot \mathbb{E}_{P^{w,k},\pi^k} \left[ \sum_{h=1}^H \Gamma_h^k \right].$$

This finishes the proof.  $\square$

*Proof of Lemma D.1 in TV Setting.* The result directly follows from combining Lemma D.4 and Lemma D.5, along with applying a union bound.  $\square$

### D.3.2 Proof of Lemma D.1 in KL Setting

**Lemma D.6** (Dual formulation). (Tang et al., 2025, Proposition 4.5) For the optimization problem  $\inf_{\mu \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu} [V_{h+1}] + \beta \text{KL}(\mu \| \mu^o))$ , we have its dual formulation as follows

$$\inf_{\mu \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu} [V_{h+1}] + \beta \text{KL}(\mu \| \mu^o)) = -\beta \log \mathbb{E}_{s \sim \mu^o} [e^{-V_{h+1}(s)/\beta}].$$

**Lemma D.7** (Optimism). If we set the bonus term as follows

$$\begin{aligned} \Gamma_h^k(s, a) &= c_{\text{KL}} \sum_{i=1}^d \|\phi_i(s, a) \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}}, \\ \text{where } c_{\text{KL}} &= (1 + 2\beta e^{\beta^{-1}H}) H d \cdot \xi_{\text{KL}}, \\ \text{where } \xi_{\text{KL}} &= 80 + \sqrt{40 \log (64\beta K^{11/2} H d^{1/2} (1 + 2\beta e^{\beta^{-1}H})^2 |\mathcal{A}|^3 / \delta)}, \end{aligned}$$

then for any policy  $\pi$  and any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ , with probability at least  $1 - \delta$ , we have  $Q_h^{k,\beta,\eta}(s, a) \geq Q_h^{\pi,\beta,\eta}(s, a)$ . Specially, by setting  $\pi = \pi^*$ , we have  $Q_h^{k,\beta,\eta}(s, a) \geq Q_h^{\pi^*,\beta,\eta}(s, a)$ .

*Proof.* We prove this by induction. First,  $h = H + 1$  holds trivially since  $Q_{H+1}^{k,\beta,\eta}(s, a) = 0 = Q_{H+1}^{\pi,\beta,\eta}(s, a)$ .

Assume  $Q_{h+1}^{k,\beta,\eta}(s, a) \geq Q_{h+1}^{\pi,\beta,\eta}(s, a)$  holds, due to the choice of  $\pi_{h+1}^k$  in (4.3), we have

$$\begin{aligned} V_{h+1}^{k,\beta,\eta}(s) &= \langle Q_{h+1}^{k,\beta,\eta}(s, \cdot), \pi_{h+1}^k(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}^k) \\ &\geq \langle Q_{h+1}^{k,\beta,\eta}(s, \cdot), \pi_{h+1}(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}) \\ &\geq \langle Q_{h+1}^{\pi,\beta,\eta}(s, \cdot), \pi_{h+1}(\cdot | s) \rangle - \Omega_s^\eta(\pi_{h+1}) = V_{h+1}^{\pi,\beta,\eta}(s). \end{aligned}$$

Recall that we denote  $Q_h^{k,\beta,\eta}$  as the optimistic estimation in  $k$ -th episode, that is,

$$Q_h^{k,\beta,\eta}(s, a) \leftarrow \min \{ \phi(s, a)^\top (\theta_h^k + \mathbf{w}_h^k) + \Gamma_h^k(s, a), H - h + 1 \}. \quad (\text{D.38})$$

If  $Q_h^{k,\beta,\eta}(s, a) = H - h + 1$ , then it follows immediately that

$$Q_h^{k,\beta,\eta}(s, a) = H - h + 1 \geq Q_h^{\pi,\beta,\eta}(s, a)$$

by the definition of  $Q_h^{\pi,\beta,\eta}(s, a)$ . Otherwise, we can infer that

$$\begin{aligned}
& Q_h^{k,\beta,\eta} - Q_h^{\pi,\beta,\eta} \\
&= \Gamma_h^k + \phi^\top(\theta_h^k + w_h^k) - \phi^\top \left[ \theta_h + \inf_{\mu \in \Delta(S)} (\mathbb{E}_{s \sim \mu} [V_{h+1}^{\pi,\beta,\eta}] + \beta \text{KL}(\mu \| \mu^o)) \right] \\
&\geq \Gamma_h^k + \phi^\top(\theta_h^k + w_h^k) - \phi^\top \left[ \theta_h + \inf_{\mu \in \Delta(S)} (\mathbb{E}_{s \sim \mu} [V_{h+1}^{k,\beta,\eta}] + \beta \text{KL}(\mu \| \mu^o)) \right] \tag{D.39}
\end{aligned}$$

$$\begin{aligned}
&= \Gamma_h^k + \left\langle \phi, \theta_h^k - \beta \log \max \left\{ \widehat{\mathbb{E}}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}], e^{-H/\beta} \right\} - \theta_h + \beta \log \mathbb{E}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] \right\rangle \tag{D.40}
\end{aligned}$$

$$\begin{aligned}
&\geq \Gamma_h^k - \left\langle \phi, |\theta_h^k - \theta_h| + \beta e^{\beta^{-1}H} \left| \max \left\{ \widehat{\mathbb{E}}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}], e^{-H/\beta} \right\} - \mathbb{E}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] \right| \right\rangle \\
&\geq \Gamma_h^k - \left\langle \phi, |\theta_h^k - \theta_h| + \beta e^{\beta^{-1}H} \left| \widehat{\mathbb{E}}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] - \mathbb{E}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] \right| \right\rangle \\
&= \Gamma_h^k - \sum_{i=1}^d \underbrace{\phi_i \mathbf{1}_i^\top |\theta_h^k - \theta_h|}_{(i)} - \beta e^{\beta^{-1}H} \sum_{i=1}^d \underbrace{\phi_i \mathbf{1}_i^\top \left| \widehat{\mathbb{E}}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] - \mathbb{E}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] \right|}_{(ii)}, \tag{D.41}
\end{aligned}$$

where (D.39) follows from the induction assumption, and (D.40) is obtained from the algorithm update formulation and the dual formulation.

For term (i) in (D.41), we have

$$\begin{aligned}
& \left| \phi_i \mathbf{1}_i^\top (\theta_h^k - \theta_h) \right| \\
&= \left| \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot r_h^\tau \right) \right. \\
&\quad \left. - \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top + \lambda \cdot I \right) \cdot \theta_h \right| \\
&= \lambda \left| \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \cdot \theta_h \right| \tag{D.42}
\end{aligned}$$

$$\leq \lambda \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \cdot \|\theta_h\|_{(\Lambda_h^k)^{-1}} \tag{D.43}$$

$$\leq \lambda \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \cdot \lambda_{\max}((\Lambda_h^k)^{-1})^{\frac{1}{2}} \|\theta_h\|_2 \tag{D.44}$$

$$\leq \lambda^{\frac{1}{2}} \sqrt{d} \cdot \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}},$$

where (D.42) is obtained from the algorithm update formulation and the definition of  $\Lambda_h^k$ , (D.43) follows from the Cauchy-Schwartz inequality, (D.44) is because  $\|x\|_A \leq \sqrt{\lambda_{\max}(A)} \|x\|_2$ .

For term (ii) in (D.41), we decompose it as follows.

$$\begin{aligned}
& \left| \phi_i \mathbf{1}_i^\top \left( \widehat{\mathbb{E}}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] - \mathbb{E}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] \right) \right| \\
&= \left| \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot e^{-V_{h+1}^{k,\beta,\eta}(s_h^\tau)/\beta} \right) \right. \\
&\quad \left. - \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top + \lambda \cdot I \right) \mathbb{E}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] \right| \tag{D.45}
\end{aligned}$$

$$\begin{aligned}
&= \left| \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau(e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}) \right) \right. \\
&\quad \left. + \lambda \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \mathbb{E}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] \right| \tag{D.46}
\end{aligned}$$

$$\begin{aligned}
&\leq \left| \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau(e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}) \right) \right| \\
&\quad + \lambda \left| \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \mathbb{E}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] \right| \\
&\leq \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \cdot \underbrace{\left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau(e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}) \right\|_{(\Lambda_h^k)^{-1}}}_{\text{(iii)}} \\
&\quad + \lambda \underbrace{\|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \cdot \left\| \mathbb{E}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] \right\|_{(\Lambda_h^k)^{-1}}}_{\text{(iv)}}, \tag{D.47}
\end{aligned}$$

where (D.45) is obtained from the algorithm update formulation and the definition of  $\Lambda_h^k$ , (D.46) is from the definition in (D.23), (D.47) follows from the Cauchy-Schwartz inequality.

For term (iv) in (D.47), since  $e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta} \leq 1$ , we have

$$\begin{aligned}
&\left\| \mathbb{E}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] \right\|_{(\Lambda_h^k)^{-1}} \\
&\leq \lambda_{\max}((\Lambda_h^k)^{-1})^{\frac{1}{2}} \left\| \mathbb{E}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] \right\|_2 \\
&\leq \lambda^{-\frac{1}{2}} \sqrt{d}, \tag{D.48}
\end{aligned}$$

where (D.48) is because  $\|x\|_A \leq \sqrt{\lambda_{\max}(A)} \|x\|_2$ .

Before we continue to bound term (iii), we prove a auxiliary result first. That is,  $\|\mathbf{w}_h^k\|_2 \leq \sqrt{kd/\lambda}$ . From the algorithm update formulation, we have

$$\begin{aligned}
\left\| \widehat{\mathbb{E}}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}] \right\|_2 &= \left\| (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot e^{-V_{h+1}^{k,\beta,\eta}(s_h^\tau)/\beta} \right) \right\|_2 \\
&\leq \sum_{\tau=1}^{k-1} \left\| (\Lambda_h^k)^{-1} \phi(s_h^\tau, a_h^\tau) \right\|_2 \\
&= \sum_{\tau=1}^{k-1} \sqrt{\phi(s_h^\tau, a_h^\tau)^\top (\Lambda_h^k)^{-1/2} (\Lambda_h^k)^{-1} (\Lambda_h^k)^{-1/2} \phi(s_h^\tau, a_h^\tau)} \\
&\leq \frac{1}{\sqrt{\lambda}} \sum_{\tau=1}^{k-1} \sqrt{\phi(s_h^\tau, a_h^\tau)^\top (\Lambda_h^k)^{-1} \phi(s_h^\tau, a_h^\tau)} \\
&\leq \frac{\sqrt{k}}{\sqrt{\lambda}} \sqrt{\sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau)^\top (\Lambda_h^k)^{-1} \phi(s_h^\tau, a_h^\tau)}
\end{aligned}$$

$$\begin{aligned}
&= \frac{\sqrt{k}}{\sqrt{\lambda}} \sqrt{\text{tr}\left((\Lambda_h^k)^{-1}(\Lambda_h^k - \lambda I)\right)} \\
&\leq \frac{\sqrt{k}}{\sqrt{\lambda}} \sqrt{\text{tr}(I)} = \sqrt{kd/\lambda}.
\end{aligned}$$

Therefore,

$$\begin{aligned}
\|\mathbf{w}_h^k\|_2 &= \|\beta \log \max \left\{ \widehat{\mathbb{E}}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}], e^{-H/\beta} \right\}\|_2 \\
&\leq \beta \sqrt{d} \max \left\{ \log \left( \|\widehat{\mathbb{E}}_{s \sim \mu^o} [e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}]\|_2 \right), H/\beta \right\} \\
&\leq \beta \sqrt{d} \max \left\{ \log (\sqrt{kd/\lambda}), H/\beta \right\} \\
&\leq \beta \sqrt{d} (\sqrt{kd/\lambda} + H/\beta) \\
&\leq 2\beta H d \sqrt{k}.
\end{aligned}$$

Now we are ready to bound term (iii) in (D.47). Let  $\mathcal{V}_h$  denote a class of functions mapping from  $\mathcal{S}$  to  $\mathbb{R}$  with the following form

$$V(\cdot) = \min \left\{ \max_{a \in \mathcal{A}} \left\{ \phi(\cdot, a)^\top \mathbf{w} + \beta \sum_{i=1}^d \|\phi_i(\cdot, a) \mathbb{1}_i\|_{\Lambda^{-1}} \right\}, H - h + 1 \right\},$$

where the parameters  $(w, \beta, \Lambda)$  satisfy  $\|w\| \leq L$ ,  $\beta \in [0, B]$ ,  $\lambda_{\min} \geq \lambda$ . Here we set  $L = 2\beta H d \sqrt{k}$  and  $B = (1 + 2\beta e^{\beta^{-1}H}) H d \cdot \xi_{\text{KL}}$ , where  $\xi_{\text{KL}} = 80 + \sqrt{40 \log(64\beta K^{11/2} H d^{1/2} (1 + 2\beta e^{\beta^{-1}H})^2 |\mathcal{A}|^3 / \delta)}$ .

Let  $\mathcal{N}_h(\epsilon; \mathcal{V}_h)$  be the minimum  $\epsilon$ -cover of  $\mathcal{V}_h$ . Therefore, for any  $V \in \mathcal{V}_h$  and  $\alpha \in [0, H]$ , there exists a function  $V_\epsilon \in \mathcal{N}_h(\epsilon; \mathcal{V}_h)$  and  $\alpha_\epsilon \in \mathcal{N}_h(\epsilon; [0, H])$  such that

$$\sup_{s \in \mathcal{S}} |V(s) - V_\epsilon(s)| \leq \epsilon.$$

Then we have

$$\begin{aligned}
(\text{iii})^2 &= \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau(e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta}) \right\|_{(\Lambda_h^k)^{-1}}^2 \\
&\leq 2 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau(e^{-V_\epsilon(s)/\beta}) \right\|_{(\Lambda_h^k)^{-1}}^2 \\
&\quad + 2 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau(e^{-V_{h+1}^{k,\beta,\eta}(s)/\beta} - e^{-V_\epsilon(s)/\beta}) \right\|_{(\Lambda_h^k)^{-1}}^2 \tag{D.49}
\end{aligned}$$

$$\leq 2 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau(e^{-V_\epsilon(s)/\beta}) \right\|_{(\Lambda_h^k)^{-1}}^2 + 2(4\epsilon)^2 k^2 / \lambda \tag{D.50}$$

$$\leq 4 \left( d \log(1 + k/\lambda) + 2 \log \frac{KH |\mathcal{N}_h(\epsilon; \mathcal{V}_h)|}{\delta} \right) + 32\epsilon^2 k^2 / \lambda \tag{D.51}$$

$$\begin{aligned}
&\leq 4H^2 \left( d \log(1 + k/\lambda) + 2d \log(1 + 4L|\mathcal{A}|/\epsilon) + 2d^2 \log(1 + 8d^{1/2} B^2 |\mathcal{A}|^2 / (\lambda \epsilon^2)) \right. \\
&\quad \left. + 2 \log(KH/\delta) \right) + 32\epsilon^2 k^2 / \lambda \tag{D.52}
\end{aligned}$$



$$\begin{aligned}
&\leq 4H^2 \left( 2d \log(k) + 4d \log(8\beta k^{3/2}|\mathcal{A}|) + 4d^2 \log(8k^2 B^2 |\mathcal{A}|^2 / (H^2 d^{3/2})) \right. \\
&\quad \left. + 2 \log(KH/\delta) \right) + 32H^2 d^2 \\
&\leq 16H^2 d^2 \log(64\beta K^{11/2} B^2 |\mathcal{A}|^3 / (H d^{3/2} \delta)) + 32H^2 d^2,
\end{aligned} \tag{D.53}$$

where (D.49) is because  $\|a + b\|_A^2 \leq 2\|a\|_A^2 + 2\|b\|_A^2$ , (D.50) follows from the inequality  $|e^x - e^y| \leq e^{|x-y|} - 1$  for  $\max\{x, y\} \leq 0$  and  $e^t \leq 1 + 2t$  for  $t \in [0, 1]$ , (D.51) is from Lemma E.7 together with a union bound over all  $(k, h) \in [K] \times [H]$ , (D.52) makes use of Lemma D.2, we set  $\epsilon = Hd/k$  and  $\lambda = 1$  and apply the inequality  $\log(1 + x) \leq 2 \log x$  for  $x \geq 2$  in (D.53).

Combining everything together, recall that we set  $\lambda = 1$ , we have

$$\begin{aligned}
\text{(i)} + \beta e^{\beta^{-1}H} \text{(ii)} &\leq \sqrt{d} \cdot \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \\
&\quad + \beta e^{\beta^{-1}H} \left( \sqrt{16H^2 d^2 \log(64\beta K^{11/2} B^2 |\mathcal{A}|^3 / (H d^{3/2} \delta)) + 32H^2 d^2} + \sqrt{d} \right) \cdot \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \\
&\leq (1 + 2\beta e^{\beta^{-1}H}) Hd \sqrt{40 \log(64\beta K^{11/2} B^2 |\mathcal{A}|^3 / (H d^{3/2} \delta))} \cdot \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}}.
\end{aligned}$$

With our choice of  $B$ , and noting that  $x \geq \max\{(2B + \sqrt{A})^2, 1\}$  guarantees  $x \geq A + B \log x$ , we conclude the proof.  $\square$

**Lemma D.8.** If we set the bonus term to be the same as in Lemma D.7, then for any  $(k, s, a) \in [K] \times \mathcal{S} \times \mathcal{A}$ , with probability at least  $1 - \delta$ , the sum of estimation errors can be bounded as

$$Q_1^{k,\beta,\eta}(s, a) - Q_1^{\pi^k,\beta,\eta}(s, a) \leq 2 \cdot \mathbb{E}_{P^{w,k}, \pi^k} \left[ \sum_{h=1}^H \Gamma_h^k(s_h^k, a_h^k) \right],$$

where  $P_h^{w,k}$  is defined by  $P_h^{w,k}(s, a) = \phi(s, a)^\top \cdot \mu_h^{w,k}$ .

*Proof.* From the definition, We have that

$$\begin{aligned}
&Q_h^{k,\beta,\eta} - Q_h^{\pi^k,\beta,\eta} \\
&= \Gamma_h^k + \phi^\top(\theta_h^k + w_h^k) - \phi^\top \left[ \theta_h + \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{\pi^k,\beta,\eta}] + \beta \text{KL}(\mu_h \|\mu_h^o)) \right] \\
&= \Gamma_h^k + \phi^\top(\theta_h^k - \theta_h) + \phi^\top \left[ w_h^k - \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{k,\beta,\eta}] + \beta \text{KL}(\mu_h \|\mu_h^o)) \right] \\
&\quad + \phi^\top \left[ \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{k,\beta,\eta}] + \beta \text{KL}(\mu_h \|\mu_h^o)) - \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{\pi^k,\beta,\eta}] + \beta \text{KL}(\mu_h \|\mu_h^o)) \right] \\
&\leq 2 \cdot \Gamma_h^k + \phi^\top \left[ \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{k,\beta,\eta}] + \beta \text{KL}(\mu_h \|\mu_h^o)) - \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{\pi^k,\beta,\eta}] + \beta \text{KL}(\mu_h \|\mu_h^o)) \right]
\end{aligned} \tag{D.54}$$

$$\begin{aligned}
&\leq 2 \cdot \Gamma_h^k + \phi^\top \mathbb{E}_{s \sim \mu_h^{w,k}} [V_{h+1}^{k,\beta,\eta} - V_{h+1}^{\pi^k,\beta,\eta}] \\
&\leq 2 \cdot \Gamma_h^k + \mathbb{E}_{s \sim P_h^w} [V_{h+1}^{k,\beta,\eta} - V_{h+1}^{\pi^k,\beta,\eta}],
\end{aligned} \tag{D.55}$$

where (D.54) can be derived from Lemma D.7, and (D.55) uses the definition of the worst-case transition. Apply the above inequality recursively, we have

$$Q_1^{k,\beta,\eta} - Q_1^{\pi^k,\beta,\eta} \leq 2 \cdot \mathbb{E}_{P^{w,k}, \pi^k} \left[ \sum_{h=1}^H \Gamma_h^k \right].$$

This finishes the proof.  $\square$

*Proof of Lemma D.1 in KL Setting.* The result directly follows from combining Lemma D.7 and Lemma D.8, along with applying a union bound.  $\square$

### D.3.3 Proof of Lemma D.1 in $\chi^2$ Setting

**Lemma D.9** (Dual formulation). (Tang et al., 2025, Proposition 4.7) For the optimization problem  $\inf_{\mu \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu} [V_{h+1}] + \beta \chi^2(\mu \| \mu^o))$ , we have its dual formulation as follows

$$\inf_{\mu \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu} [V_{h+1}] + \beta \chi^2(\mu \| \mu^o)) = \sup_{\alpha \in [V_{\min}, V_{\max}]} \left\{ \mathbb{E}_{s \sim \mu^o} [V_{h+1}(s)]_{\alpha} - \frac{1}{4\beta} \text{Var}_{s \sim \mu^o} [V_{h+1}(s)]_{\alpha} \right\}.$$

**Lemma D.10** (Optimism). If we set the bonus term as follows

$$\begin{aligned} \Gamma_h^k(s, a) &= c_{\chi^2} \sum_{i=1}^d \|\phi_i(s, a) \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}}, \\ \text{where } c_{\chi^2} &= (1 + H/(2\beta)) H d \cdot \xi_{\chi^2}, \\ \text{where } \xi_{\chi^2} &= 720 + 3 \sqrt{40 \log(96 K^6 H^5 (1 + H/(2\beta))^3 |\mathcal{A}|^3 / \delta)}, \end{aligned}$$

then for any policy  $\pi$  and any  $(k, h, s, a) \in [K] \times [H] \times \mathcal{S} \times \mathcal{A}$ , with probability at least  $1 - \delta$ , we have  $Q_h^{k, \beta, \eta}(s, a) \geq Q_h^{\pi, \beta, \eta}(s, a)$ . Specially, by setting  $\pi = \pi^*$ , we have  $Q_h^{k, \beta, \eta}(s, a) \geq Q_h^{\pi^*, \beta, \eta}(s, a)$ .

*Proof.* We prove this by induction. First,  $h = H + 1$  holds trivially since  $Q_{H+1}^{k, \beta, \eta}(s, a) = 0 = Q_{H+1}^{\pi, \beta, \eta}(s, a)$ .

Assume  $Q_{h+1}^{k, \beta, \eta}(s, a) \geq Q_{h+1}^{\pi, \beta, \eta}(s, a)$  holds, due to the choice of  $\pi_{h+1}^k$  in (4.3), we have

$$\begin{aligned} V_{h+1}^{k, \beta, \eta}(s) &= \langle Q_{h+1}^{k, \beta, \eta}(s, \cdot), \pi_{h+1}^k(\cdot | s) \rangle - \Omega_s^{\eta}(\pi_{h+1}^k) \\ &\geq \langle Q_{h+1}^{k, \beta, \eta}(s, \cdot), \pi_{h+1}(\cdot | s) \rangle - \Omega_s^{\eta}(\pi_{h+1}) \\ &\geq \langle Q_{h+1}^{\pi, \beta, \eta}(s, \cdot), \pi_{h+1}(\cdot | s) \rangle - \Omega_s^{\eta}(\pi_{h+1}) = V_{h+1}^{\pi, \beta, \eta}(s). \end{aligned}$$

Recall that we denote  $Q_h^{k, \beta, \eta}$  as the optimistic estimation in  $k$ -th episode, that is,

$$Q_h^{k, \beta, \eta}(s, a) \leftarrow \min \{ \phi(s, a)^{\top} (\theta_h^k + \mathbf{w}_h^k) + \Gamma_h^k(s, a), H - h + 1 \}. \quad (\text{D.56})$$

If  $Q_h^{k, \beta, \eta}(s, a) = H - h + 1$ , then it follows immediately that

$$Q_h^{k, \beta, \eta}(s, a) = H - h + 1 \geq Q_h^{\pi, \beta, \eta}(s, a)$$

by the definition of  $Q_h^{\pi, \beta, \eta}(s, a)$ . Otherwise, we introduce a few auxiliary notations

$$\widetilde{\mathbb{E}}_{s \sim \mu^o} [V_{h+1}^{k, \beta, \eta}]_{\alpha_i} = \underset{\mathbf{w} \in \mathbb{R}^d}{\text{argmin}} \sum_{\tau=1}^k ([V_{h+1}^{k, \beta, \eta}(s_{h+1}^{\tau})]_{\alpha_i} - \phi(s_h^{\tau}, a_h^{\tau})^{\top} \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2,$$

$$\tilde{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \sum_{\tau=1}^k ([V_{h+1}^{k,\beta,\eta}(s_{h+1}^\tau)]_{\alpha_i}^2 - \phi(s_h^\tau, a_h^\tau)^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2,$$

and denote

$$\alpha_i = \operatorname{argmax}_{\alpha \in [0, H]} \left\{ \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_\alpha + \frac{1}{4\beta} (\mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_\alpha)^2 - \frac{1}{4\beta} \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_\alpha^2 \right\}.$$

Then

$$\begin{aligned} & Q_h^{k,\beta,\eta} - Q_h^{\pi,\beta,\eta} \\ &= \Gamma_h^k + \phi^\top(\boldsymbol{\theta}_h^k + \mathbf{w}_h^k) - \phi^\top \left[ \boldsymbol{\theta}_h + \inf_{\mu \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu}[V_{h+1}^{\pi,\beta,\eta}] + \beta \chi^2(\mu \| \mu^o)) \right] \\ &\geq \Gamma_h^k + \phi^\top(\boldsymbol{\theta}_h^k + \mathbf{w}_h^k) - \phi^\top \left[ \boldsymbol{\theta}_h + \inf_{\mu \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu}[V_{h+1}^{k,\beta,\eta}] + \beta \chi^2(\mu \| \mu^o)) \right] \end{aligned} \quad (\text{D.57})$$

$$\begin{aligned} &= \Gamma_h^k + \left\langle \phi, \boldsymbol{\theta}_h^k + \sup_{\alpha \in [0, H]} \left\{ \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_\alpha + \frac{1}{4\beta} (\widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_\alpha)^2 - \frac{1}{4\beta} \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_\alpha^2 \right\} \right. \\ &\quad \left. - \boldsymbol{\theta}_h - \sup_{\alpha \in [0, H]} \left\{ \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_\alpha + \frac{1}{4\beta} (\mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_\alpha)^2 - \frac{1}{4\beta} \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_\alpha^2 \right\} \right\rangle \end{aligned} \quad (\text{D.58})$$

$$\begin{aligned} &\geq \Gamma_h^k + \left\langle \phi, \boldsymbol{\theta}_h^k + \left\{ \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i} + \frac{1}{4\beta} (\widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i})^2 - \frac{1}{4\beta} \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 \right\} \right. \\ &\quad \left. - \boldsymbol{\theta}_h - \left\{ \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i} + \frac{1}{4\beta} (\mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i})^2 - \frac{1}{4\beta} \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 \right\} \right\rangle \end{aligned} \quad (\text{D.59})$$

$$\begin{aligned} &\geq \Gamma_h^k - \left\langle \phi, |\boldsymbol{\theta}_h^k - \boldsymbol{\theta}_h| + \left| \frac{1}{4\beta} \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 - \frac{1}{4\beta} \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 \right| \right. \\ &\quad \left. + \left| \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i} + \frac{1}{4\beta} (\widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i})^2 - \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i} - \frac{1}{4\beta} (\mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i})^2 \right| \right\rangle \\ &= \Gamma_h^k - \sum_{i=1}^d \phi_i \mathbb{1}_i^\top |\boldsymbol{\theta}_h^k - \boldsymbol{\theta}_h| - \frac{1}{4\beta} \sum_{i=1}^d \phi_i \mathbb{1}_i^\top \left| \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 - \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 \right| \\ &\quad - \sum_{i=1}^d \phi_i \mathbb{1}_i^\top \left| \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i} - \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i} \right| \cdot \left| \frac{1}{4\beta} \left( \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i} + \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i} \right) + 1 \right| \\ &= \Gamma_h^k - \sum_{i=1}^d \phi_i \mathbb{1}_i^\top |\boldsymbol{\theta}_h^k - \boldsymbol{\theta}_h| - \frac{1}{4\beta} \sum_{i=1}^d \phi_i \mathbb{1}_i^\top \left| \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 - \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 \right| \\ &\quad \times \underbrace{\left| \frac{\widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 - \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2}{\widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 - \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2} \right|}_{(i)} - \sum_{i=1}^d \phi_i \mathbb{1}_i^\top \left| \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i} - \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i} \right| \\ &\quad \times \underbrace{\left[ \frac{1}{4\beta} \left( \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i} + \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i} \right) + 1 \right] \cdot \frac{\widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i} - \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}}{\widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i} - \mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}}}_{(ii)} \end{aligned} \quad (\text{D.60})$$

$$\begin{aligned}
&\geq \Gamma_h^k - \sum_{i=1}^d \underbrace{\phi_i \mathbf{1}_i^\top |\boldsymbol{\theta}_h^k - \boldsymbol{\theta}_h|}_{\text{(iii)}} - \frac{1}{4\beta} \sum_{i=1}^d \underbrace{\phi_i \mathbf{1}_i^\top \left| \tilde{\mathbb{E}}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 - \mathbb{E}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 \right|}_{\text{(v)}} \\
&\quad - \left(1 + \frac{H}{2\beta}\right) \sum_{i=1}^d \underbrace{\phi_i \mathbf{1}_i^\top \left| \tilde{\mathbb{E}}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i} - \mathbb{E}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i} \right|}_{\text{(iv)}}, \tag{D.61}
\end{aligned}$$

where (D.57) follows from the induction assumption, (D.58) is obtained from the algorithm update formulation and the dual formulation, and (D.59) comes from our choice of  $\alpha_i$ . To establish in (D.60) that (i)  $\leq 1$  and (ii)  $\leq 1 + H/(2\beta)$  and hence that (D.61) holds, we analyze the value of  $\tilde{\mathbb{E}}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i}$  or  $\tilde{\mathbb{E}}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2$  under three cases:

$$\begin{cases} \text{(i)} \leq 1 & \tilde{\mathbb{E}}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 < 0 \\ \text{(i)} \leq 1 & \tilde{\mathbb{E}}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 \in [0, H] \\ \text{(i)} \leq 1 & \tilde{\mathbb{E}}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 > H, \end{cases} \quad \begin{cases} \text{(ii)} \leq 1 + H/(4\beta) & \tilde{\mathbb{E}}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i} < 0 \\ \text{(ii)} \leq 1 + H/(2\beta) & \tilde{\mathbb{E}}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i} \in [0, H] \\ \text{(ii)} \leq 1 + H/(2\beta) & \tilde{\mathbb{E}}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i} > H. \end{cases}$$

For term (iii) in (D.61), we have

$$\begin{aligned}
&\left| \phi_i \mathbf{1}_i^\top (\boldsymbol{\theta}_h^k - \boldsymbol{\theta}_h) \right| \\
&= \left| \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \boldsymbol{\phi}(s_h^\tau, a_h^\tau) \cdot r_h^\tau \right) \right. \\
&\quad \left. - \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \boldsymbol{\phi}(s_h^\tau, a_h^\tau) \boldsymbol{\phi}(s_h^\tau, a_h^\tau)^\top + \lambda \cdot I \right) \cdot \boldsymbol{\theta}_h \right| \\
&= \lambda \left| \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \cdot \boldsymbol{\theta}_h \right| \tag{D.62}
\end{aligned}$$

$$\leq \lambda \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \cdot \|\boldsymbol{\theta}_h\|_{(\Lambda_h^k)^{-1}} \tag{D.63}$$

$$\leq \lambda \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \cdot \lambda_{\max}((\Lambda_h^k)^{-1})^{\frac{1}{2}} \|\boldsymbol{\theta}_h\|_2 \tag{D.64}$$

$$\leq \lambda^{\frac{1}{2}} \sqrt{d} \cdot \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}},$$

where (D.62) is obtained from the algorithm update formulation and the definition of  $\Lambda_h^k$ , (D.63) follows from the Cauchy-Schwartz inequality, (D.64) is because  $\|x\|_A \leq \sqrt{\lambda_{\max}(A)} \|x\|_2$ .

For term (iv) in (D.61), we decompose it as follows.

$$\begin{aligned}
&\left| \phi_i \mathbf{1}_i^\top \left( \tilde{\mathbb{E}}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i} - \mathbb{E}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i} \right) \right| \\
&= \left| \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \boldsymbol{\phi}(s_h^\tau, a_h^\tau) \cdot [V_{h+1}^{k,\beta,\eta}(s_{h+1}^\tau)]_{\alpha_i} \right) \right. \\
&\quad \left. - \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \boldsymbol{\phi}(s_h^\tau, a_h^\tau) \boldsymbol{\phi}(s_h^\tau, a_h^\tau)^\top + \lambda \cdot I \right) \cdot \mathbb{E}_{s \sim \boldsymbol{\mu}^o} [V_{h+1}^{k,\beta,\eta}(s)]_{\alpha_i} \right| \tag{D.65} \\
&= \left| \phi_i \mathbf{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \boldsymbol{\phi}(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_{h+1}^{k,\beta,\eta}(s)]_{\alpha_i}) \right) \right|
\end{aligned}$$

$$+ \lambda \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}(s)]_{\alpha_i} \Big| \quad (\text{D.66})$$

$$\begin{aligned} &\leq \left| \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_{h+1}^{k,\beta,\eta}(s)]_{\alpha_i}) \right) \right| \\ &\quad + \lambda \left| \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}(s)]_{\alpha_i} \right| \\ &\leq \underbrace{\|\phi_i \mathbb{1}_i\|_{(\Lambda_h^k)^{-1}} \cdot \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_{h+1}^{k,\beta,\eta}(s)]_{\alpha_i}) \right\|_{(\Lambda_h^k)^{-1}}}_{(\text{vi})} \\ &\quad + \underbrace{\lambda \|\phi_i \mathbb{1}_i\|_{(\Lambda_h^k)^{-1}} \cdot \left\| \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}(s)]_{\alpha_i} \right\|_{(\Lambda_h^k)^{-1}}}_{(\text{vii})}, \end{aligned} \quad (\text{D.67})$$

where (D.65) is obtained from the algorithm update formulation and the definition of  $\Lambda_h^k$ , (D.66) is from the definition in (D.23), (D.67) follows from the Cauchy-Schwartz inequality.

For term (v) in (D.61), we decompose it as follows.

$$\begin{aligned} &\left| \phi_i \mathbb{1}_i^\top \left( \widetilde{\mathbb{E}}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 - \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 \right) \right| \\ &= \left| \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot [V_{h+1}^{k,\beta,\eta}(s_h^\tau)]_{\alpha_i}^2 \right) \right. \\ &\quad \left. - \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top + \lambda \cdot I \right) \cdot \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}(s)]_{\alpha_i}^2 \right| \end{aligned} \quad (\text{D.68})$$

$$\begin{aligned} &= \left| \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_{h+1}^{k,\beta,\eta}(s)]_{\alpha_i}^2) \right) \right. \\ &\quad \left. + \lambda \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}(s)]_{\alpha_i}^2 \right| \end{aligned} \quad (\text{D.69})$$

$$\begin{aligned} &\leq \left| \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \left( \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_{h+1}^{k,\beta,\eta}(s)]_{\alpha_i}^2) \right) \right| \\ &\quad + \lambda \left| \phi_i \mathbb{1}_i^\top (\Lambda_h^k)^{-1} \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}(s)]_{\alpha_i}^2 \right| \\ &\leq \underbrace{\|\phi_i \mathbb{1}_i\|_{(\Lambda_h^k)^{-1}} \cdot \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_{h+1}^{k,\beta,\eta}(s)]_{\alpha_i}^2) \right\|_{(\Lambda_h^k)^{-1}}}_{(\text{viii})} \\ &\quad + \underbrace{\lambda \|\phi_i \mathbb{1}_i\|_{(\Lambda_h^k)^{-1}} \cdot \left\| \mathbb{E}_{s \sim \mu^o} [V_{h+1}^{k,\beta,\eta}(s)]_{\alpha_i}^2 \right\|_{(\Lambda_h^k)^{-1}}}_{(\text{ix})}, \end{aligned} \quad (\text{D.70})$$

where (D.68) is obtained from the algorithm update formulation and the definition of  $\Lambda_h^k$ , (D.69) is from the definition in (D.23), (D.70) follows from the Cauchy-Schwartz inequality.

For term (vii) in (D.67) and For term (ix) in (D.70), since  $V_{h+1}^{k,\beta,\eta} \leq H$ , we have

$$\begin{aligned}\left\|\mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}\right\|_{(\Lambda_h^k)^{-1}} &\leq \lambda_{\max}((\Lambda_h^k)^{-1})^{\frac{1}{2}} \left\|\mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}\right\|_2 \leq \lambda^{-\frac{1}{2}} H \sqrt{d}, \\ \left\|\mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2\right\|_{(\Lambda_h^k)^{-1}} &\leq \lambda_{\max}((\Lambda_h^k)^{-1})^{\frac{1}{2}} \left\|\mathbb{E}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2\right\|_2 \leq \lambda^{-\frac{1}{2}} H^2 \sqrt{d}.\end{aligned}$$

Before we continue to bound term (vi) and (viii), we prove a auxiliary result first. That is,  $\|\mathbf{w}_h^k\|_2 \leq (H + H^2/(2\beta))\sqrt{d}$ . From the algorithm update formulation, we have

$$\begin{aligned}\|\mathbf{w}_h^k\|_2 &= \left\| \sup_{\alpha \in [0, H]} \left\{ \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha} + \frac{1}{4\beta} (\widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha})^2 - \frac{1}{4\beta} \widehat{\mathbb{E}}_{s \sim \mu^o}[V_{h+1}^{k,\beta,\eta}]_{\alpha}^2 \right\} \right\|_2 \\ &\leq \left\| \mathbf{1} \left( H + \frac{H^2}{2\beta} \right) \right\|_2 = \left( H + \frac{H^2}{2\beta} \right) \sqrt{d}.\end{aligned}$$

Now we are ready to bound term (vi) and (viii). Let  $\mathcal{V}_h$  denote a class of functions mapping from  $\mathcal{S}$  to  $\mathbb{R}$  with the following from

$$V(\cdot) = \min \left\{ \max_{a \in \mathcal{A}} \left\{ \phi(\cdot, a)^\top \mathbf{w} + \beta \sum_{i=1}^d \|\phi_i(\cdot, a) \mathbb{1}_i\|_{\Lambda^{-1}} \right\}, H - h + 1 \right\},$$

where the parameters  $(w, \beta, \Lambda)$  satisfy  $\|\mathbf{w}\| \leq L$ ,  $\beta \in [0, B]$ ,  $\lambda_{\min} \geq \lambda$ . Here we set  $L = (H + H^2/(2\beta))\sqrt{d}$  and  $B = (1 + H/(2\beta))Hd \cdot \xi_{\chi^2}$ , where  $\xi_{\chi^2} = 720 + 3\sqrt{40 \log(96K^6 H^5 (1 + H/(2\beta))^3 |\mathcal{A}|^3 / \delta)}$ .

Let  $\mathcal{N}_h(\epsilon; \mathcal{V}_h)$  be the minimum  $\epsilon$ -cover of  $\mathcal{V}_h$ ,  $\mathcal{N}_h(\epsilon; [0, H])$  be the minimum  $\epsilon$ -cover of  $[0, H]$ . Therefore, for any  $V \in \mathcal{V}_h$  and  $\alpha \in [0, H]$ , there exists a function  $V_\epsilon \in \mathcal{N}_h(\epsilon; \mathcal{V}_h)$  and  $\alpha_\epsilon \in \mathcal{N}_h(\epsilon; [0, H])$  such that

$$\sup_{s \in \mathcal{S}} |V(s) - V_\epsilon(s)| \leq \epsilon, \quad \left| \min_{s'} V_{h+1}^{k,\beta,\eta}(s') + \beta - \alpha_\epsilon \right| \leq \epsilon.$$

Then we have

$$\begin{aligned}(\text{vi})^2 &= \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_{h+1}^{k,\beta,\eta}]_{\alpha_i}) \right\|_{(\Lambda_h^k)^{-1}}^2 \\ &\leq 2 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_{h+1}^{k,\beta,\eta}]_{\alpha_\epsilon}) \right\|_{(\Lambda_h^k)^{-1}}^2 \\ &\quad + 2 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_{h+1}^{k,\beta,\eta}]_{\alpha_i} - [V_{h+1}^{k,\beta,\eta}]_{\alpha_\epsilon}) \right\|_{(\Lambda_h^k)^{-1}}^2 \tag{D.71}\end{aligned}$$

$$\begin{aligned}&\leq 4 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_\epsilon]_{\alpha_\epsilon}) \right\|_{(\Lambda_h^k)^{-1}}^2 + 4 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_{h+1}^{k,\beta,\eta}]_{\alpha_\epsilon} - [V_\epsilon]_{\alpha_\epsilon}) \right\|_{(\Lambda_h^k)^{-1}}^2 \\ &\quad + 2(2\epsilon)^2 k^2 / \lambda \tag{D.72} \\ &\leq 4 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_\epsilon]_{\alpha_\epsilon}) \right\|_{(\Lambda_h^k)^{-1}}^2 + 4(2\epsilon)^2 k^2 / \lambda + 2(2\epsilon)^2 k^2 / \lambda\end{aligned}$$

$$\begin{aligned}
&\leq 4 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_\epsilon]_{\alpha_\epsilon}) \right\|_{(\Lambda_h^k)^{-1}}^2 + 24\epsilon^2 k^2 / \lambda \\
&\leq 4H^2 \left( d \log(1 + k/\lambda) + 2 \log \frac{KH |\mathcal{N}_h(\epsilon; \mathcal{V}_h)| \cdot |\mathcal{N}_h(\epsilon; [0, H])|}{\delta} \right) + 24\epsilon^2 k^2 / \lambda
\end{aligned} \tag{D.73}$$

$$\begin{aligned}
&\leq 4H^2 \left( d \log(1 + k/\lambda) + 2d \log(1 + 4L|\mathcal{A}|/\epsilon) + 2d^2 \log(1 + 8d^{1/2}B^2|\mathcal{A}|^2/(\lambda\epsilon^2)) \right. \\
&\quad \left. + 2 \log(3H/\epsilon) + 2 \log(KH/\delta) \right) + 24\epsilon^2 k^2 / \lambda
\end{aligned} \tag{D.74}$$

$$\begin{aligned}
&\leq 4H^2 \left( 2d \log(k) + 4d \log(4k(1 + H/(2\beta))|\mathcal{A}|/(d^{1/2})) + 4d^2 \log(8k^2 B^2 |\mathcal{A}|^2 / (H^2 d^{3/2})) \right. \\
&\quad \left. + 2 \log(3k/d) + 2 \log(KH/\delta) \right) + 24H^2 d^2 \\
&\leq 16H^2 d^2 \log(96K^6 B^2 (1 + H/(2\beta))|\mathcal{A}|^3 / (H d^3 \delta)) + 24H^2 d^2,
\end{aligned} \tag{D.75}$$

where (D.71) and (D.72) is because  $\|a + b\|_A^2 \leq 2\|a\|_A^2 + 2\|b\|_A^2$ , (D.73) is from Lemma E.7 together with a union bound over all  $(k, h) \in [K] \times [H]$ , (D.74) makes use of Lemmas D.2 and E.5, we set  $\epsilon = Hd/k$  and  $\lambda = 1$  and apply the inequality  $\log(1 + x) \leq 2 \log x$  for  $x \geq 2$  in (D.75).

And

$$\begin{aligned}
(\text{viii})^2 &= \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2) \right\|_{(\Lambda_h^k)^{-1}}^2 \\
&\leq 2 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_{h+1}^{k,\beta,\eta}]_{\alpha_\epsilon}^2) \right\|_{(\Lambda_h^k)^{-1}}^2 \\
&\quad + 2 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_{h+1}^{k,\beta,\eta}]_{\alpha_i}^2 - [V_{h+1}^{k,\beta,\eta}]_{\alpha_\epsilon}^2) \right\|_{(\Lambda_h^k)^{-1}}^2
\end{aligned} \tag{D.76}$$

$$\begin{aligned}
&\leq 4 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_\epsilon]_{\alpha_\epsilon}^2) \right\|_{(\Lambda_h^k)^{-1}}^2 + 4 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_{h+1}^{k,\beta,\eta}]_{\alpha_\epsilon}^2 - [V_\epsilon]_{\alpha_\epsilon}^2) \right\|_{(\Lambda_h^k)^{-1}}^2 \\
&\quad + 2(2H\epsilon)^2 k^2 / \lambda
\end{aligned} \tag{D.77}$$

$$\begin{aligned}
&\leq 4 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_\epsilon]_{\alpha_\epsilon}^2) \right\|_{(\Lambda_h^k)^{-1}}^2 + 4(2H\epsilon)^2 k^2 / \lambda + 2(2H\epsilon)^2 k^2 / \lambda \\
&\leq 4 \left\| \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \cdot \text{err}_h^\tau([V_\epsilon]_{\alpha_\epsilon}^2) \right\|_{(\Lambda_h^k)^{-1}}^2 + 24H^2 \epsilon^2 k^2 / \lambda \\
&\leq 4H^2 \left( d \log(1 + k/\lambda) + 2 \log \frac{KH |\mathcal{N}_h(\epsilon; \mathcal{V}_h)| \cdot |\mathcal{N}_h(\epsilon; [0, H])|}{\delta} \right) + 24H^2 \epsilon^2 k^2 / \lambda
\end{aligned} \tag{D.78}$$

$$\begin{aligned}
&\leq 4H^2 \left( d \log(1 + k/\lambda) + 2d \log(1 + 4L|\mathcal{A}|/\epsilon) + 2d^2 \log(1 + 8d^{1/2}B^2|\mathcal{A}|^2/(\lambda\epsilon^2)) \right. \\
&\quad \left. + 2 \log(3H/\epsilon) + 2 \log(KH/\delta) \right) + 24H^2 \epsilon^2 k^2 / \lambda
\end{aligned} \tag{D.79}$$

$$\begin{aligned}
&\leq 4H^2 \left( 2d \log(k) + 4d \log(4kH(1 + H/(2\beta))|\mathcal{A}|/(d^{1/2})) + 4d^2 \log(8k^2 B^2 |\mathcal{A}|^2 / (d^3/2)) \right. \\
&\quad \left. + 2 \log(3kH/d) + 2 \log(KH/\delta) \right) + 24H^2 d^2
\end{aligned} \tag{D.80}$$

$$\leq 16H^2d^2 \log(96K^6H^3B^2(1+H/(2\beta))|\mathcal{A}|^3/(d^3\delta)) + 24H^2d^2,$$

where (D.76) and (D.77) is because  $\|a+b\|_A^2 \leq 2\|a\|_A^2 + 2\|b\|_A^2$ , (D.78) is from Lemma E.7 together with a union bound over all  $(k, h) \in [K] \times [H]$ , (D.79) makes use of Lemmas D.2 and E.5, we set  $\epsilon = d/k$  and  $\lambda = 1$  and apply the inequality  $\log(1+x) \leq 2\log x$  for  $x \geq 2$  in (D.80).

Combining everything together, recall that we set  $\lambda = 1$ , we have

$$\begin{aligned} & \text{(iii)} + \frac{1}{4\beta} \text{(v)} + \left(1 + \frac{H}{2\beta}\right) \text{(iv)} \\ & \leq \sqrt{d} \cdot \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \\ & \quad + \frac{1}{4\beta} \left( \sqrt{16H^2d^2 \log(96K^6H^3B^2(1+H/(2\beta))|\mathcal{A}|^3/(d^3\delta)) + 24H^2d^2 + H^2\sqrt{d}} \right) \cdot \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \\ & \quad + \left(1 + \frac{H}{2\beta}\right) \left( \sqrt{16H^2d^2 \log(96K^6B^2(1+H/(2\beta))|\mathcal{A}|^3/(Hd^3\delta)) + 24H^2d^2 + H\sqrt{d}} \right) \cdot \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}} \\ & \leq 3 \left(1 + \frac{H}{2\beta}\right) Hd \sqrt{40 \log(96K^6H^3B^2(1+H/(2\beta))|\mathcal{A}|^3/(d^3\delta))} \cdot \|\phi_i \mathbf{1}_i\|_{(\Lambda_h^k)^{-1}}. \end{aligned}$$

With our choice of  $B$ , and noting that  $x \geq \max\{(2B + \sqrt{A})^2, 1\}$  guarantees  $x \geq A + B \log x$ , we conclude the proof.  $\square$

**Lemma D.11.** If we set the bonus term to be the same as in Lemma D.10, then for any  $(k, s, a) \in [K] \times \mathcal{S} \times \mathcal{A}$ , with probability at least  $1 - \delta$ , the sum of estimation errors can be bounded as

$$Q_1^{k,\beta,\eta}(s, a) - Q_1^{\pi^k,\beta,\eta}(s, a) \leq 2 \cdot \mathbb{E}_{P^{w,k}, \pi^k} \left[ \sum_{h=1}^H \Gamma_h^k(s_h^k, a_h^k) \right],$$

where  $P_h^{w,k}$  is defined by  $P_h^{w,k}(s, a) = \phi(s, a)^\top \cdot \mu_h^{w,k}$ .

*Proof.* From the definition, We have that

$$\begin{aligned} & Q_h^{k,\beta,\eta} - Q_h^{\pi^k,\beta,\eta} \\ & = \Gamma_h^k + \phi^\top(\theta_h^k + w_h^k) - \phi^\top \left[ \theta_h + \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{\pi^k,\beta,\eta}] + \beta \chi^2(\mu_h \| \mu_h^o)) \right] \\ & = \Gamma_h^k + \phi^\top(\theta_h^k - \theta_h) + \phi^\top \left[ w_h^k - \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{k,\beta,\eta}] + \beta \chi^2(\mu_h \| \mu_h^o)) \right] \\ & \quad + \phi^\top \left[ \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{k,\beta,\eta}] + \beta \chi^2(\mu_h \| \mu_h^o)) - \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{\pi^k,\beta,\eta}] + \beta \chi^2(\mu_h \| \mu_h^o)) \right] \\ & \leq 2 \cdot \Gamma_h^k + \phi^\top \left[ \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{k,\beta,\eta}] + \beta \chi^2(\mu_h \| \mu_h^o)) - \inf_{\mu_h \in \Delta(\mathcal{S})} (\mathbb{E}_{s \sim \mu_h} [V_{h+1}^{\pi^k,\beta,\eta}] + \beta \chi^2(\mu_h \| \mu_h^o)) \right] \end{aligned} \tag{D.81}$$

$$\leq 2 \cdot \Gamma_h^k + \phi^\top \mathbb{E}_{s \sim \mu_h^{w,k}} [V_{h+1}^{k,\beta,\eta} - V_{h+1}^{\pi^k,\beta,\eta}] \tag{D.82}$$

$$\leq 2 \cdot \Gamma_h^k + \mathbb{E}_{s \sim P_h^w} [V_{h+1}^{k,\beta,\eta} - V_{h+1}^{\pi^k,\beta,\eta}],$$



where (D.81) can be derived from Lemma D.10, and (D.82) uses the definition of the worst-case transition. Apply the above inequality recursively, we have

$$Q_1^{k,\beta,\eta} - Q_1^{\pi^k,\beta,\eta} \leq 2 \cdot \mathbb{E}_{P^{w,k},\pi^k} \left[ \sum_{h=1}^H \Gamma_h^k \right].$$

This finishes the proof.  $\square$

*Proof of Lemma D.1 in  $\chi^2$  Setting.* The result directly follows from combining Lemma D.10 and Lemma D.11, along with applying a union bound.  $\square$

## E Auxiliary Lemmas

Here, we present some auxiliary lemmas which are useful in the proof.

**Lemma E.1** (Hoeffding's inequality). (Vershynin, 2018, Theorem 2.2.6) Let  $X_1, \dots, X_T$  be independent random variables. Assume that  $X_t \in [0, M]$  for every  $t$  with  $M > 0$ . Let  $S_T = \frac{1}{T} \sum_{t=1}^T X_t$ , then for any  $\epsilon > 0$ , we have

$$\mathbb{P}(|S_T - \mathbb{E}[S_T]| \geq \epsilon) \leq 2 \exp \left( - \frac{2T\epsilon^2}{M^2} \right).$$

**Lemma E.2** (Self-bounding variance inequality). (Maurer and Pontil, 2009, Theorem 10) Let  $X_1, \dots, X_T$  be independent and identically distributed random variables with finite variance. Assume that  $X_t \in [0, M]$  for every  $t$  with  $M > 0$ . Let  $S_T^2 = \frac{1}{T} \sum_{t=1}^T X_t^2 - (\frac{1}{T} \sum_{t=1}^T X_t)^2$ , then for any  $\epsilon > 0$ , we have

$$\mathbb{P}(|S_T - \sqrt{\text{Var}(X_1)}| \geq \epsilon) \leq 2 \exp \left( - \frac{T\epsilon^2}{2M^2} \right).$$

**Lemma E.3.** (Weissman et al., 2003, Theorem 2.1) Let  $P$  be a probability distribution over  $\mathcal{S} = \{s_1, \dots, s_S\}$ ,  $X_1, \dots, X_T$  be independent and identically distributed random variables distributed according to  $P$ . Let  $\hat{P}(s) = \frac{1}{T} \sum_{t=1}^T \mathbb{1}\{X_t = s\}$ , then for any  $\epsilon > 0$ , we have

$$\mathbb{P}(\|P - \hat{P}\|_1 \geq \epsilon) \leq 2^S \exp \left( - \frac{T\epsilon^2}{2} \right).$$

**Lemma E.4.** (Panaganti and Kalathil, 2022, Lemma 7) We define  $\mathcal{V} = \{V \in \mathbb{R}^S : \|V\|_\infty \leq V_{\max}\}$ . Let  $\mathcal{N}_{\mathcal{V}}(\epsilon)$  be a minimal  $\epsilon$ -cover of  $\mathcal{V}$  with respect to the distance metric  $d(V, V') = \|V - V'\|_\infty$  for some fixed  $\epsilon \in (0, 1)$ . Then we have

$$\log |\mathcal{N}_{\mathcal{V}}(\epsilon)| \leq |\mathcal{S}| \cdot \log \left( \frac{3V_{\max}}{\epsilon} \right).$$

**Lemma E.5.** (Van Handel, 2014, Lemma 5.13) Denote the  $\epsilon$ -covering number of the closed interval  $[a, b]$  for some real number  $b > a$  with respect to the distance metric  $\text{dist}(\alpha_1, \alpha_2) = |\alpha_1 - \alpha_2|$  as  $\mathcal{N}_h(\epsilon; [a, b])$ , then we have  $\mathcal{N}_h(\epsilon; [a, b]) \leq 3(b - a)/\epsilon$ .

**Lemma E.6.** (Jin et al., 2020, Lemma D.5) For any  $\epsilon > 0$ , the  $\epsilon$ -covering number of the Euclidean ball in  $\mathbb{R}^d$  with radius  $R > 0$  is upper bounded by  $(1 + 2R/\epsilon)^d$ .

**Lemma E.7.** (Abbasi-Yadkori et al., 2011, Theorem 1) Let  $\{\epsilon_t\}_{t=1}^\infty$  be a real-valued stochastic process with corresponding filtration  $\{\mathcal{F}_t\}_{t=0}^\infty$ . Let  $\epsilon_t|\mathcal{F}_{t-1}$  be mean-zero and  $\sigma$ -sub-Gaussian. Let  $\{\phi_t\}_{t=1}^\infty$  be a  $\mathbb{R}^d$ -valued stochastic process where  $\phi_t$  is  $\mathcal{F}_{t-1}$  measurable. Assume  $\Lambda_0$  is a  $d \times d$  positive definite matrix, and let  $\Lambda_t = \Lambda_0 + \sum_{s=1}^t \phi_s \phi_s^\top$ . Then for any  $\delta > 0$ , with probability at least  $1 - \delta$ , we have for all  $t \geq 0$ ,

$$\left\| \sum_{s=1}^t \phi_s \epsilon_s \right\|_{\Lambda_t^{-1}}^2 \leq 2\sigma^2 \log \left[ \frac{\det(\Lambda_t)^{\frac{1}{2}} \det(\Lambda_0)^{-\frac{1}{2}}}{\delta} \right].$$

**Lemma E.8.** (Liu and Xu, 2024a, Corollary 5.3) For all  $(\pi, h) \in \Pi \times [H]$ , assume that

$$\mathbb{E}_\pi[\phi(s_h, a_h)\phi(s_h, a_h)^\top] \geq \alpha I,$$

where  $\alpha > 0$ . Then with probability at least  $1 - \delta$ , we have for all  $(k, h) \in [K] \times [H]$ ,

$$\lambda_{\min}(\Lambda_h^k) \geq \max \{ \alpha(k-1) + \lambda - \sqrt{32k \log(dKH/\delta)}, \lambda \}.$$

## References

- ABBASI-YADKORI, Y., PÁL, D. and SZEPESVÁRI, C. (2011). Improved algorithms for linear stochastic bandits. *Advances in Neural Information Processing Systems* **24**. 5, 14, 54
- BAI, Y., JONES, A., NDOUSSE, K., ASKELL, A., CHEN, A., DASARMA, N., DRAIN, D., FORT, S., GANGULI, D., HENIGHAN, T. ET AL. (2022). Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862* . 1
- BLANCHET, J., LU, M., ZHANG, T. and ZHONG, H. (2023). Double pessimism is provably efficient for distributionally robust offline reinforcement learning: Generic algorithm and robust partial coverage. *Advances in Neural Information Processing Systems* **36** 66845–66859. 3
- DERMAN, E., GEIST, M. and MANNOR, S. (2021). Twice regularized mdps and the equivalence between robustness and regularization. *Advances in Neural Information Processing Systems* **34** 22274–22287. 1, 3, 5, 7
- EYSENBACH, B. and LEVINE, S. (2021). Maximum entropy rl (provably) solves some robust rl problems. *arXiv preprint arXiv:2103.06257* . 1
- GEIST, M., SCHERRER, B. and PIETQUIN, O. (2019). A theory of regularized markov decision processes. In *International Conference on Machine Learning*. PMLR. 1
- GHOSH, D., ATIA, G. K. and WANG, Y. (2025). Provably near-optimal distributionally robust reinforcement learning in online settings. *arXiv preprint arXiv:2508.03768* . 3
- HE, Y., LIU, Z., WANG, W. and XU, P. (2025). Sample complexity of distributionally robust off-dynamics reinforcement learning with online interaction. In *Forty-second International Conference on Machine Learning*. 1, 2, 3, 4, 5, 8, 9, 22, 25, 27

- HUSAIN, H., CIOSEK, K. and TOMIOKA, R. (2021). Regularized policies are reward robust. In *International Conference on Artificial Intelligence and Statistics*. PMLR. [1](#)
- IYENGAR, G. N. (2005). Robust dynamic programming. *Mathematics of Operations Research* **30** 257–280. [1](#), [2](#)
- 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. [9](#), [54](#)
- LIU, Z., WANG, W. and XU, P. (2024). Upper and lower bounds for distributionally robust off-dynamics reinforcement learning. *arXiv preprint arXiv:2409.20521* . [2](#)
- LIU, Z. and XU, P. (2024a). Distributionally robust off-dynamics reinforcement learning: Provable efficiency with linear function approximation. In *International Conference on Artificial Intelligence and Statistics*. PMLR. [2](#), [3](#), [5](#), [8](#), [9](#), [13](#), [14](#), [35](#), [54](#)
- LIU, Z. and XU, P. (2024b). Minimax optimal and computationally efficient algorithms for distributionally robust offline reinforcement learning. In *Advances in Neural Information Processing Systems*, vol. 37. [3](#)
- LIU, Z. and XU, P. (2025). Linear mixture distributionally robust markov decision processes. *arXiv preprint arXiv:2505.18044* . [3](#)
- LU, M., ZHONG, H., ZHANG, T. and BLANCHET, J. (2024). Distributionally robust reinforcement learning with interactive data collection: Fundamental hardness and near-optimal algorithms. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. [2](#), [3](#), [5](#), [8](#)
- MANNOR, S., MEBEL, O. and XU, H. (2016). Robust mdps with k-rectangular uncertainty. *Mathematics of Operations Research* **41** 1484–1509. [2](#)
- MAURER, A. and PONTIL, M. (2009). Empirical bernstein bounds and sample variance penalization. *arXiv preprint arXiv:0907.3740* . [53](#)
- MNIH, V., BADIA, A. P., MIRZA, M., GRAVES, A., LILICRAP, T., HARLEY, T., SILVER, D. and KAVUKCUOGLU, K. (2016). Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*. PmLR. [1](#)
- NILIM, A. and EL GHAOU, L. (2005). Robust control of markov decision processes with uncertain transition matrices. *Operations Research* **53** 780–798. [1](#), [2](#)
- OUYANG, L., WU, J., JIANG, X., ALMEIDA, D., WAINWRIGHT, C., MISHKIN, P., ZHANG, C., AGARWAL, S., SLAMA, K., RAY, A. ET AL. (2022). Training language models to follow instructions with human feedback. *Advances in neural information processing systems* **35** 27730–27744. [1](#)
- PANAGANTI, K. and KALATHIL, D. (2022). Sample complexity of robust reinforcement learning with a generative model. In *International Conference on Artificial Intelligence and Statistics*. PMLR. [3](#), [53](#)
- PANAGANTI, K., WIERMAN, A. and MAZUMDAR, E. (2024). Model-free robust  $\phi$ -divergence reinforcement learning using both offline and online data. In *Proceedings of the 41st International Conference on Machine Learning*. [2](#), [3](#), [4](#)

- PANAGANTI, K., XU, Z., KALATHIL, D. and GHAVAMZADEH, M. (2022). Robust reinforcement learning using offline data. *Advances in Neural Information Processing Systems* **35** 32211–32224. [3](#)
- ROLFE, J. T. (2016). Discrete variational autoencoders. *arXiv preprint arXiv:1609.02200* . [14](#)
- SCHULMAN, J., LEVINE, S., ABBEEL, P., JORDAN, M. and MORITZ, P. (2015). Trust region policy optimization. In *International conference on machine learning*. PMLR. [1](#)
- SCHULMAN, J., WOLSKI, F., DHARIWAL, P., RADFORD, A. and KLIMOV, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* . [1](#)
- SHI, L. and CHI, Y. (2024). Distributionally robust model-based offline reinforcement learning with near-optimal sample complexity. *Journal of Machine Learning Research* **25** 1–91. [3](#)
- SHI, L., LI, G., WEI, Y., CHEN, Y., GEIST, M. and CHI, Y. (2024). The curious price of distributional robustness in reinforcement learning with a generative model. *Advances in Neural Information Processing Systems* **36**. [3](#)
- TANG, C., LIU, Z. and XU, P. (2025). Robust offline reinforcement learning with linearly structured  $\mathbb{F}$ -divergence regularization. In *Forty-second International Conference on Machine Learning*. [3](#), [4](#), [9](#), [36](#), [41](#), [46](#)
- VAN HANDEL, R. (2014). Probability in high dimension. Tech. rep. [53](#)
- VERSHYNIN, R. (2018). *High-dimensional probability: An introduction with applications in data science*, vol. 47. Cambridge university press. [53](#)
- VIEILLARD, N., KOZUNO, T., SCHERRER, B., PIETQUIN, O., MUNOS, R. and GEIST, M. (2020). Leverage the average: an analysis of kl regularization in reinforcement learning. *Advances in Neural Information Processing Systems* **33** 12163–12174. [1](#)
- WANG, H., SHI, L. and CHI, Y. (2024a). Sample complexity of offline distributionally robust linear markov decision processes. In *Reinforcement Learning Conference*. [3](#)
- WANG, S., SI, N., BLANCHET, J. and ZHOU, Z. (2023). A finite sample complexity bound for distributionally robust q-learning. In *International Conference on Artificial Intelligence and Statistics*. PMLR. [3](#)
- WANG, S., SI, N., BLANCHET, J. and ZHOU, Z. (2024b). Sample complexity of variance-reduced distributionally robust q-learning. *Journal of Machine Learning Research* **25** 1–77. [3](#)
- WEISSMAN, T., ORDENTLICH, E., SEROUSSI, G., VERDU, S. and WEINBERGER, M. J. (2003). Inequalities for the l1 deviation of the empirical distribution. *Hewlett-Packard Labs, Tech. Rep* 125. [53](#)
- WIESEMANN, W., KUHN, D. and RUSTEM, B. (2013). Robust markov decision processes. *Mathematics of Operations Research* **38** 153–183. [2](#)
- XU, H. and MANNOR, S. (2006). The robustness-performance tradeoff in markov decision processes. *Advances in Neural Information Processing Systems* **19**. [2](#)

- YANG, W., LI, X. and ZHANG, Z. (2019). A regularized approach to sparse optimal policy in reinforcement learning. *Advances in Neural Information Processing Systems* **32**. 1
- YANG, W., WANG, H., KOZUNO, T., JORDAN, S. M. and ZHANG, Z. (2023). Robust markov decision processes without model estimation. *arXiv preprint arXiv:2302.01248* . 1, 2, 3, 4
- YANG, W., ZHANG, L. and ZHANG, Z. (2022). Toward theoretical understandings of robust markov decision processes: Sample complexity and asymptotics. *The Annals of Statistics* **50** 3223–3248. 3
- ZHANG, R., HU, Y. and LI, N. (2024). Soft robust mdps and risk-sensitive mdps: Equivalence, policy gradient, and sample complexity. In *The Twelfth International Conference on Learning Representations*. 1, 2, 3, 4
- ZHAO, H., YE, C., XIONG, W., GU, Q. and ZHANG, T. (2025). Logarithmic regret for online kl-regularized reinforcement learning. *arXiv preprint arXiv:2502.07460* . 1
- ZHOU, Z., ZHOU, Z., BAI, Q., QIU, L., BLANCHET, J. and GLYNN, P. (2021). Finite-sample regret bound for distributionally robust offline tabular reinforcement learning. In *International Conference on Artificial Intelligence and Statistics*. PMLR. 3