

---

# Almost Optimal Algorithms for Two-player Zero-Sum Markov Games with Linear Function Approximation

---

Zixiang Chen<sup>1</sup> Dongruo Zhou<sup>1</sup> Quanquan Gu<sup>1</sup>

## Abstract

We study reinforcement learning for two-player zero-sum Markov games with simultaneous moves in the finite-horizon setting, where the transition kernel of the underlying Markov games can be parameterized by a linear function over the current state, both players' actions and the next state. In particular, we assume that we can control both players and aim to find the Nash Equilibrium by minimizing the duality gap. We propose an algorithm Nash-UCRL based on the principle "Optimism-in-Face-of-Uncertainty". Our algorithm only needs to find a Coarse Correlated Equilibrium (CCE), which is computationally very efficient. Specifically, we show that Nash-UCRL can provably achieve an  $\tilde{O}(\sqrt{d^2 H^2 T})$  regret, where  $d$  is the linear function dimension,  $H$  is the length of the game and  $T$  is the total number of steps in the game. To assess the optimality of our algorithm, we also prove an  $\Omega(\sqrt{d^2 H^2 T})$  lower bound on the regret. Our upper bound matches the lower bound up to logarithmic factors, which suggests the optimality of our algorithm.

## 1. Introduction

Multi-agent reinforcement learning (MARL) has achieved tremendous practical success across a wide range of machine learning tasks. For Markov games with large state and action spaces, it is natural to use linear function approximation. In particular, Xie et al. (2020) proposed an OMNI-VI algorithm for Markov games where the transition kernel and reward function possess a linear structure, and achieved an  $\tilde{O}(\sqrt{d^3 H^3 T})$  regret, with  $d$  being the dimension of the linear structure. However, as we will show in this paper, there is still a gap between the upper and lower bounds of existing algorithms for Markov games with linear structures.

---

<sup>\*</sup>Equal contribution <sup>1</sup>Department of Computer Science, University of California, Los Angeles, USA. Correspondence to: Quanquan Gu <qgu@cs.ucla.edu>.

In this paper, we consider Markov games with a linear mixture structure and design a minimax optimal algorithm for learning zero-sum Markov games based on the principle of "Optimism-in-Face-of-Uncertainty" without assuming the access to the generative model or well-explored behavior policy. We summarize the contributions of our work as follows:

- We propose a Nash-UCRL algorithm for general Markov games (which can be specialized to turn-based games) that can provably achieve an  $\tilde{O}(\sqrt{d^2 H^2 T})$  upper bound on the regret, where  $d$  is the dimension of linear mixture structure,  $H$  is the length of the game, and  $T$  the total number of steps in the Markov game.
- To access the optimality of our algorithm Nash-UCRL, we also prove a regret lower bound  $\Omega(\sqrt{d^2 H^2 T})$ . Our upper bound matches the lower bound up to logarithmic factors, which suggests the optimality of our algorithm.

**Notations** We use lower case letters to denote scalars, lower and upper case bold letters to denote vectors and matrices. We use  $\|\cdot\|$  to indicate Euclidean norm, and for a semi-positive definite matrix  $\Sigma$  and any vector  $\mathbf{x}$ ,  $\|\mathbf{x}\|_{\Sigma} := \|\Sigma^{1/2}\mathbf{x}\| = \sqrt{\mathbf{x}^{\top}\Sigma\mathbf{x}}$ . For a real value  $x$  and an interval  $[a, b]$ , we use  $[x]_{[a, b]}$  to indicate the projection of  $x$  onto  $[a, b]$ . We also use the standard  $O$  and  $\Omega$  notations. We say  $a_n = O(b_n)$  if and only if  $\exists C > 0, N > 0, \forall n > N, a_n \leq Cb_n$ ;  $a_n = \Omega(b_n)$  if  $a_n \geq Cb_n$ . The notation  $\tilde{O}$  is used to hide logarithmic factors.

## 2. Preliminaries

**Simultaneous-move MG.** Two-player zero-sum Markov game (MG) (Shapley, 1953; Littman, 1994) is a generalization of the standard Markov decision process (MDP). Formally, we denote a two-player zero-sum simultaneous-moves episodic Markov Game by a tuple  $M(\mathcal{S}, \mathcal{A}_{\max}, \mathcal{A}_{\min}, H, \{r_h\}_{h=1}^H, \{\mathbb{P}_h\}_{h=1}^H)$ .  $\mathcal{S}$  is a countable state space,  $\mathcal{A}_{\max}, \mathcal{A}_{\min}$  are the finite action spaces of the max-player and the min-player respectively.  $H$  is the length of the game/episode. For simplicity, we assume the reward function for the max-player  $\{r_h\}_{h=1}^H$  is deterministic and known function  $r_h : \mathcal{S} \times \mathcal{A}_{\max} \times \mathcal{A}_{\min} \rightarrow [-1, 1]$ .

$\mathbb{P}_h(s'|s, a, b)$  is the transition probability function which denotes the probability for state  $s$  to transit to state  $s'$  given players' action pair  $(a, b)$  at step  $h$ .

We now define the stochastic policies, which give distributions over the actions. A policy  $\pi = \{\pi_h : \mathcal{S} \rightarrow \Delta_{\mathcal{A}_{\max}}\}_{h=1}^H$  is a collection of functions which map a state  $s \in \mathcal{S}$  to a distribution of actions. Here  $\Delta_{\mathcal{A}_{\max}}$  is the probability simplex over action set  $\mathcal{A}_{\max}$ . Similarly, we can define a policy  $\nu$  for the min-player. We use the notation  $\pi_h(a|s)$  and  $\nu_h(b|s)$  to present the probability of taking action  $a$  or  $b$  for state  $s$  at step  $h$  under Markov policy  $\pi, \nu$  respectively. We use the notation  $Q_h^{\pi, \nu} : \mathcal{S} \times \mathcal{A}_{\max} \times \mathcal{A}_{\min} \rightarrow \mathbb{R}$  to present the action-value function (a.k.a.,  $Q$  function), and the notation  $V_h^{\pi, \nu} : \mathcal{S} \rightarrow \mathbb{R}$  to present the value function..

**Learning Objective.** The goal of the max-player is to maximize the total rewards. The goal of the min-player is to minimize the total rewards that the max-player will get because this is a zero-mean game. In other words, the max-player wants to maximize  $V_h^{\pi, \nu}(\cdot)$  by choosing a good policy  $\pi$ , while the min-player wants to minimize  $V_h^{\pi, \nu}(\cdot)$  by choose a good policy  $\nu$ . Accordingly, we can define the action-value function and the value function when the max-player gives the best response to a fixed policy  $\nu$  of the min-player:

$$Q_h^{*, \nu}(s, a, b) = \max_{\pi} Q_h^{\pi, \nu}(s, a, b),$$

$$V_h^{*, \nu}(s) = \max_{\pi} V_h^{\pi, \nu}(s).$$

By symmetry, we can also define  $Q_h^{\pi, *}(s, a, b), V_h^{\pi, *}(s)$ . A Nash Equilibrium (NE) of the game is a pair of policies  $\pi^*, \nu^*$  such that

$$V_1^{\pi^*, \nu^*}(s) = V_1^{\pi^*, *}(s) = V_1^{*, \nu^*}(s), \text{ for all } s \in \mathcal{S}. \quad (2.1)$$

For most applications, they are the ultimate solutions we want to pursue. We can measure the suboptimality of learned policies  $\{\pi^k, \nu^k\}_k$  by the gap between their performance and the performance of the optimal strategy (i.e., Nash equilibrium) when playing against the best responses respectively:

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

Accordingly, we aim to design a learning algorithm that outputs a sequence  $\{\pi^k, \nu^k\}_k$  based on past information, and minimize the regret. This regret has been widely used in previous work that studies the offline learning of two-player game (Jin et al., 2018; Xie et al., 2020; Liu et al., 2020).

**Episodic Linear Mixture Markov Games.** In this work, we consider a class of MGs called *linear mixture MGs*, inspired by the linear mixture/kernel MDPs studied in Modi et al. (2020); Jia et al. (2020); Ayoub et al. (2020) for the single-agent RL. Linear mixture MGs assume that at each

step  $h$ , the transition probability function  $\mathbb{P}_h(s'|s, a, b)$  is a linear combination of  $d$  feature mappings  $\phi_i(s'|s, a, b)$ , i.e.,

$$\mathbb{P}_h(s'|s, a, b) = \sum_{i=1}^d \theta_{i,h} \phi_i(s'|s, a, b),$$

where each feature mapping  $\phi_i(s'|s, a, b)$  is a function defined on the state-action-action-state pair  $(s, a, b, s') \in \mathcal{S} \times \mathcal{A}_{\max} \times \mathcal{A}_{\min} \times \mathcal{S}$ . For the sake of simplicity, we use a vector function  $\phi = [\phi_1, \dots, \phi_d] \in \mathbb{R}^d$  to denote the collection of  $\phi_i$ . After proper normalization, we assume  $\phi$  satisfy that for any bounded function  $V : \mathcal{S} \rightarrow [-1, 1]$  and any tuple  $(s, a, b) \in \mathcal{S} \times \mathcal{A}_{\max} \times \mathcal{A}_{\min}$ , we have

$$\|\phi_V(s, a, b)\|_2 \leq 1, \quad (2.2)$$

where  $\phi_V(s, a, b) = \sum_{s' \in \mathcal{S}} \phi(s'|s, a, b) V(s')$ . Formally, we define linear mixture MGs as follows:

**Definition 2.1.**  $M(\mathcal{S}, \mathcal{A}_{\max}, \mathcal{A}_{\min}, H, \{r_h\}_{h=1}^H, \{\mathbb{P}_h\}_{h=1}^H)$  is called a time inhomogeneous, episodic  $B$ -bounded linear mixture MG if there exist  $H$  vectors  $\theta_h \in \mathbb{R}^d$  satisfying for any  $h \in [H]$ ,  $\|\theta_h\|_2 \leq B$ , and feature mapping  $\phi$  satisfying (2.2), such that  $\mathbb{P}_h(s'|s, a, b) = \langle \phi(s'|s, a, b), \theta_h \rangle$  for any state-action-action-state triplet  $(s, a, b, s')$  and any step  $h$ . We denote the linear mixture MG by  $M_{\theta}$  for simplicity.

In this paper, we assume the underlying linear mixture MG is parameterized by  $\{\theta_h^*\}_{h=1}^H$ , denoted by  $M_{\theta^*}$ .

### 3. Algorithm

In this section, we propose our algorithm Nash-UCRL in Algorithm 1. Due to the space limit, we only show the detailed update rules for the max-player in Algorithm 1. For completeness, we will present the full algorithm as Algorithm 2 in Appendix C, and a turn-based version as Algorithm 3 in Appendix D. All the parameters corresponding to the max-player are marked by an overline, while the parameters for the min-player are marked by an underline. For any function  $V : \mathcal{S} \rightarrow \mathbb{R}$  we introduce the shorthands:

$$[\mathbb{P}_h V](s, a, b) = \mathbb{E}_{s' \sim \mathbb{P}_h(\cdot|s, a, b)} V(s'),$$

$$[\mathbb{V}_h V](s, a, b) = [\mathbb{P}_h^* V^2](s, a, b) - ([\mathbb{P}_h V](s, a, b))^2,$$

where  $V^2$  stands for the function whose value at  $s$  is  $V^2(s)$ .

To achieve the near-minimax optimality of solving a linear mixture MG, Nash-UCRL adopts the following three techniques, which we will introduce in sequence.

**Value-targeted regression** To find the NE of an MG, it suffices to find good estimates of the optimal value functions  $V_h^{*, \nu^k}$  and  $V_h^{\pi^k, *}$ . By the Bellman optimality equations and the definition of linear mixture MGs, it is sufficient to estimate the underlying unknown parameter  $\theta_h^*$  up to

**Algorithm 1** Nash-UCRL

---

```

1: Input: Regularization parameter  $\lambda$ , number of episode  $K$ , number of horizon  $H$ .
2: For any  $h$ ,  $\bar{\Sigma}_{1,h}^{(i)} \leftarrow \Sigma_{1,h}^{(i)} \leftarrow \lambda \mathbf{I}$ ;  $\bar{\mathbf{b}}_{1,h}^{(i)} \leftarrow \mathbf{b}_{1,h}^{(i)} \leftarrow \mathbf{0}$ ;  $\bar{\theta}_{1,h}^{(i)} \leftarrow \theta_{1,h}^{(i)} \leftarrow \mathbf{0}$ , for  $i \in \{0, 1\}$ .
3: for  $k = 1, \dots, K$  do
4:    $\bar{V}_{k,H+1}(\cdot) \leftarrow 0$ ,  $\underline{V}_{k,H+1}(\cdot) \leftarrow 0$ 
5:   for  $h = H, \dots, 1$  do
6:     Set  $\bar{Q}_{k,h}(\cdot, \cdot, \cdot) \leftarrow \left[ r_h(\cdot, \cdot, \cdot) + \langle \bar{\theta}_{k,h}^{(0)}, \phi_{\bar{V}_{k,h+1}}(\cdot, \cdot, \cdot) \rangle + \beta_k^{(0)} \|\bar{\Sigma}_{k,h}^{(0)}\|^{-1/2} \phi_{\bar{V}_{k,h+1}}(\cdot, \cdot, \cdot) \|_2 \right]_{[-H, H]}$ .
7:     Set  $\underline{Q}_{k,h}(\cdot, \cdot, \cdot) \leftarrow \left[ r_h(\cdot, \cdot, \cdot) + \langle \theta_{k,h}^{(0)}, \phi_{\underline{V}_{k,h+1}}(\cdot, \cdot, \cdot) \rangle - \beta_k^{(0)} \|\Sigma_{k,h}^{(0)}\|^{-1/2} \phi_{\underline{V}_{k,h+1}}(\cdot, \cdot, \cdot) \|_2 \right]_{[-H, H]}$ .
8:     for  $s \in \mathcal{S}$  do
9:       Let  $\mu_h^k(\cdot, \cdot | s) = \epsilon\text{-CCE}(\bar{Q}_{k,h}(s, \cdot, \cdot), \underline{Q}_{k,h}(s, \cdot, \cdot))$ .
10:       $\bar{V}_{k,h}(s) = \mathbb{E}_{(a,b) \sim \mu_h^k(\cdot, \cdot | s)} \bar{Q}_{k,h}(s, a, b)$ ,  $\underline{V}_{k,h}(s) = \mathbb{E}_{(a,b) \sim \mu_h^k(\cdot, \cdot | s)} \underline{Q}_{k,h}(s, a, b)$ 
11:       $\pi_h^k(\cdot | s) = \mathcal{P}_{\max} \mu_h^k(\cdot, \cdot | s)$ ,  $\nu_h^k(\cdot | s) = \mathcal{P}_{\min} \mu_h^k(\cdot, \cdot | s)$ 
12:    end for
13:  end for
14:  Receives  $s_1^k$ 
15:  for  $h = 1, \dots, H$  do
16:    Take action  $(a_h^k, b_h^k) \sim \mu_h^k(\cdot, \cdot | s_h^k)$  and central controller receives  $s_{h+1}^k \sim \mathbb{P}(\cdot | s_h^k, a_h^k, b_h^k)$ .
17:    Set  $\bar{E}_{k,h}$  as in (3.3),  $\bar{\sigma}_{k,h}$  as in (3.2), and  $\underline{E}_{k,h}, \underline{\sigma}_{k,h}$  in similar ways.
18:    Update  $\bar{\Sigma}_{k+1,h}^{(i)}, \bar{\mathbf{b}}_{k+1,h}^{(i)}, \Sigma_{k+1,h}^{(i)}, \mathbf{b}_{k+1,h}^{(i)}, i = 0, 1$ . (See Algorithm 2)
19:    Set  $\bar{\theta}_{k+1,h}^{(i)} \leftarrow [\bar{\Sigma}_{k+1,h}^{(i)}]^{-1} \bar{\mathbf{b}}_{k+1,h}^{(i)}$ ,  $\theta_{k+1,h}^{(i)} \leftarrow [\Sigma_{k+1,h}^{(i)}]^{-1} \mathbf{b}_{k+1,h}^{(i)}, i = 0, 1$ 
20:  end for
21: end for

```

---

good accuracy. Inspired by the UCRL with “value-targeted regression” (UCRL-VTR) proposed by Jia et al. (2020); Ayoub et al. (2020), Nash-UCRL uses a supervised learning framework to learn  $\theta_h^*$ . In the sequel, we introduce how the VTR framework works at episode  $k$  and step  $h$ . At the beginning of episode  $k$ , Nash-UCRL maintains two estimated value functions: optimistic value function  $\bar{V}_{k,h+1}$  for the max-player, which overestimates the optimal value function  $V_h^{*, \nu^k}$ , and optimistic value function  $\underline{V}_{k,h+1}$  for the min-player, which underestimates the value function  $V_h^{\pi^k, *}$ . We focus on the overestimate  $\bar{V}_{k,h+1}$  first. To encourage the agent to explore, Nash-UCRL constructs an optimistic action-value function  $\bar{Q}_{k,h}$  following the “optimism-in-the-face-of-uncertainty” principle as showed in Line 6 of Algorithm 1. Finally, Nash-UCRL constructs the optimistic value functions  $\bar{V}_{k,h}, \underline{V}_{k,h}$  and the policy  $\mu_h^k$  based on  $\bar{Q}_{k,h}, \underline{Q}_{k,h}$  for the current episode and step (which will be specified later).

**Coarse Correlated Equilibrium (CCE).** After we get  $\bar{Q}_{k,h}(s, \cdot, \cdot)$  for the max-player and  $\underline{Q}_{k,h}(s, \cdot, \cdot)$  for the min-player, we solve a general-sum matrix game to find the Coarse Correlated Equilibrium (CCE), following Xie et al. (2020). Here we give the formal definition of CCE:

**Definition 3.1** (Moulin & Vial 1978; Aumann 1987). Given two payoff matrices  $Q_{\max}, Q_{\min} \in \mathbb{R}^{|\mathcal{A}_{\max}| \times |\mathcal{A}_{\min}|}$ , we denote

the  $\epsilon$ -Coarse Correlated Equilibrium ( $\epsilon$ -CCE) as a joint distribution  $\sigma$  over  $\mathcal{A}_{\max}$  and  $\mathcal{A}_{\min}$  satisfying that

$$\begin{aligned} \mathbb{E}_{(a,b) \sim \sigma} Q_{\max}(a, b) &\geq \max_{a' \in \mathcal{A}_{\max}} \mathbb{E}_{b \sim \mathcal{P}_{\min} \sigma} Q_{\max}(a', b) - \epsilon, \\ \mathbb{E}_{(a,b) \sim \sigma} Q_{\min}(a, b) &\leq \min_{b' \in \mathcal{A}_{\min}} \mathbb{E}_{a \sim \mathcal{P}_{\max} \sigma} Q_{\min}(a, b') + \epsilon. \end{aligned}$$

Nash-UCRL computes the distribution  $\mu_h^k(\cdot, \cdot | s)$ , a  $\epsilon$ -CCE of  $\bar{Q}_{k,h}, \underline{Q}_{k,h}$  for each state  $s$  in Line 9. Then Nash-UCRL selects the value functions  $\bar{V}_{k,h}, \underline{V}_{k,h}$  as the expectation of  $\bar{Q}_{k,h}, \underline{Q}_{k,h}$  over the policies  $\mu_h^k$  as in Line 10 of Algorithm 1. After obtaining  $\mu^k$ , Nash-UCRL sets  $\pi_h^k(\cdot | s)$  as the marginal distribution for the max-player  $\mathcal{P}_{\max} \mu_h^k(\cdot, \cdot | s)$ , and sets  $\nu_h^k(\cdot | s)$  as the marginal distribution for the min-player  $\mathcal{P}_{\min} \mu_h^k(\cdot, \cdot | s)$ .

**Weighted linear regression for value function estimation.** Now we specify how to construct the estimators  $\bar{\theta}_{k,h}^{(0)}, \theta_{k,h}^{(0)}$ . We set the estimator  $\bar{\theta}_{k,h}$  as the minimizer to a *weighted ridge regression* problem with square loss over context-target pairs  $(\phi_{\bar{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k), \bar{V}_{k,h+1}(s_{h+1}^k))$  as follows:

$$\begin{aligned} \bar{\theta}_{k,h}^{(0)} &= \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \lambda \|\theta\|_2^2 + \sum_{j=1}^{k-1} [\langle \phi_{\bar{V}_{j,h+1}}(s_h^j, a_h^j, b_h^j), \theta \rangle \\ &\quad - \bar{V}_{j,h+1}(s_{h+1}^j) ]^2 / \bar{\sigma}_{j,h}^2, \end{aligned} \quad (3.1)$$

where  $\bar{\sigma}_{j,h}^2$  is an appropriate upper bound on the variance of the value function  $[\mathbb{V}_h \bar{V}_{j,h+1}](s_h^j, a_h^j, b_h^j)$ . In particular, we construct  $\bar{\sigma}_{k,h}^2$  as follows

$$\bar{\sigma}_{k,h} = \sqrt{\max\{H^2/d, \mathbb{V}_{k,h+1}^{\text{est}}(s_h^k, a_h^k, b_h^k) + \bar{E}_{k,h}\}}, \quad (3.2)$$

where  $[\mathbb{V}_{k,h}^{\text{est}} \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k)$  is a scalar-valued empirical estimate for the variance of the value function  $\bar{V}_{k,h+1}$  under the transition probability  $\mathbb{P}_h(\cdot | s_h^k, a_h^k, b_h^k)$ , and  $\bar{E}_{k,h}$  is an offset term that is used to guarantee that  $\bar{\sigma}_{k,h}^2$  upper bounds  $[\mathbb{V}_h \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k)$  with high probability.

Weighted ridge regression (3.1) has a closed-form solution  $\bar{\theta}_{k,h}^{(0)} = [\bar{\Sigma}_{k,h}^{(0)}]^{-1} \bar{\mathbf{b}}_{k,h}^{(0)}$  based on the covariance matrix  $\bar{\Sigma}_{k,h}^{(0)}$  and correlation vector  $\bar{\mathbf{b}}_{k,h}^{(0)}$ .  $\bar{\Sigma}_{k,h}^{(0)}$  and  $\bar{\mathbf{b}}_{k,h}^{(0)}$  can be updated in the online fashion, and their update rules are deferred to Appendix C. We deferred the construction of the empirical variance term  $[\mathbb{V}_{k,h}^{\text{est}} \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k)$  to Appendix C.

Finally, by the standard self-normalized concentration inequality for vector-valued martingales of Abbasi-Yadkori et al. (2011), we can show that, with high probability,  $\bar{\sigma}_{j,h}^2$  upper bounds  $[\mathbb{V}_h \bar{V}_{j,h+1}](s_h^j, a_h^j, b_h^j)$  if we select  $\bar{E}_{k,h}$  as follows

$$\begin{aligned} \bar{E}_{k,h} = & \min \{ H^2, \beta_k^{(1)} \| [\bar{\Sigma}_{k,h}^{(1)}]^{-1/2} \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \|_2 \} \\ & + \min \{ H^2, 2H\beta_k^{(2)} \| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k) \|_2 \}. \end{aligned} \quad (3.3)$$

Here  $\beta_k^{(1)}, \beta_k^{(2)}$  together with  $\beta_k^{(0)}$  in Line 6 are all constants setting as follows

$$\begin{aligned} \beta_k^{(0)} &= 16\sqrt{d \log(1+k/\lambda) \log(4k^2 H/\delta)} \\ &\quad + 8\sqrt{d \log(4k^2 H/\delta)} + \sqrt{\lambda} B \\ \beta_k^{(1)} &= 16\sqrt{dH^4 \log(1+KH^4/(d\lambda)) \log(4k^2 H/d\delta)} \\ &\quad + 8H^2 \log(4k^2 H/\delta) + \sqrt{\lambda} B \\ \beta_k^{(2)} &= 16d\sqrt{\log(1+k/\lambda) \log(4k^2 H/\delta)} \\ &\quad + 8\sqrt{d \log(4k^2 H/\delta)} + \sqrt{\lambda} B. \end{aligned}$$

## 4. Main Results

In this section, we present the main theoretical results. We first present the regret of Nash-UCRL.

**Theorem 4.1.** Setting  $\lambda = 1/B^2$ ,  $\epsilon = O(HT^{-1/2})$ , then with probability at least  $1 - 5\delta$ , the regret of Algorithm 1  $\text{Regret}(M_{\theta^*}, K)$  is bounded by

$$\tilde{O}(\sqrt{d^2 H^2 + dH^3 \sqrt{T}} + d^2 H^3 + d^3 H^2),$$

where  $T = KH$ .

Theorem 4.1 suggests that when  $d \geq H$  and  $T \geq d^4 H^2$ , the regret of Nash-UCRL is bounded by  $\tilde{O}(dH\sqrt{T})$ .

**Remark 4.2.** Our Nash-UCRL also enjoys a finite sample complexity. By the standard online-to-batch conversion (Xie et al., 2020), we can show that Nash-UCRL is guaranteed to find an  $\epsilon$ -approximate NE, i.e.,  $(\pi, \nu)$  satisfying  $V_1^{*,\nu} - V_1^{\pi,*} \leq \epsilon$ , within  $\tilde{O}((d^2 H^3 + dH^4)/\epsilon^2)$  episodes.

Here, we present a lower bound for linear mixture MGs. It has been shown in Zhou et al. (2020) that the regret lower bound for learning linear mixture MDPs is  $\Omega(dH\sqrt{T})$ , from which we can prove a lower bound for learning linear mixture MGs, since MDPs can be regarded as a special case of MGs with one dummy player, i.e.,  $\mathbb{P}_h(s'|s, a, b) = \mathbb{P}_h(s'|s, a)$  and  $r_h(s, a, b) = r_h(s, a)$ . Formally, we have the following lower bound:

**Theorem 4.3** (Regret lower bound, Theorem 5.6 in Zhou et al. 2020). Let  $B > 1$  and  $K \geq \max\{(d-1)^2 H/2, (d-1)/(32H(B-1))\}$ ,  $d \geq 4$ ,  $H \geq 3$ . Then for any algorithm there exists an episodic,  $B$ -bounded linear mixture MG  $M_{\theta^*}$  such that the expected regret of first  $T$  rounds is lower bounded as follows:

$$\mathbb{E}[\text{Regret}(M_{\theta^*}, K)] \geq \Omega(dH\sqrt{T}),$$

where  $T = KH$ .

**Remark 4.4.** When  $d \geq H$  and  $T \geq d^4 H^2$ , the regret of Nash-UCRL matches the lower bound up to logarithmic factors. Therefore, Nash-UCRL is nearly minimax optimal.

**Remark 4.5.** Based on a similar argument made in Zhou et al. (2020), we can show that the same lower bound holds for the Markov games with linear structures studied in Xie et al. (2020). Recall that the best-known algorithm for learning MGs with linear structures is OMNI-VI (Xie et al., 2020), which has an  $\tilde{O}(\sqrt{d^3 H^3 T})$  regret. This suggests that there is still a gap that needs to be closed for learning MGs with linear structure (Xie et al., 2020).

## 5. Conclusions

In this paper, we proposed the first provably optimal algorithm for learning two-player zero-sum Markov games with linear function approximation and without assuming access to the generative model. Specifically, we show that Nash-UCRL can provably achieve an  $\tilde{O}(\sqrt{d^2 H^2 T})$  regret, where  $d$  is the linear function dimension,  $H$  is the length of the game/episode, and  $T$  is the total number of steps in the Markov game. We also prove an  $\tilde{\Omega}(\sqrt{d^2 H^2 T})$  lower bound on the regret. Our upper bound matches the lower bound up to logarithmic factors, which suggests the optimality of our algorithm.

## References

- Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pp. 2312–2320, 2011.
- Aumann, R. J. Correlated equilibrium as an expression of bayesian rationality. *Econometrica: Journal of the Econometric Society*, pp. 1–18, 1987.
- Ayoub, A., Jia, Z., Szepesvari, C., Wang, M., and Yang, L. F. Model-based reinforcement learning with value-targeted regression. *arXiv preprint arXiv:2006.01107*, 2020.
- Azar, M. G., Osband, I., and Munos, R. Minimax regret bounds for reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 263–272. JMLR. org, 2017.
- Azuma, K. Weighted sums of certain dependent random variables. *Tohoku Mathematical Journal, Second Series*, 19(3):357–367, 1967.
- Bai, Y. and Jin, C. Provable self-play algorithms for competitive reinforcement learning. In *International Conference on Machine Learning*, pp. 551–560. PMLR, 2020.
- Cui, Q. and Yang, L. F. Minimax sample complexity for turn-based stochastic game. *arXiv preprint arXiv:2011.14267*, 2020.
- Jia, Z., Yang, L., Szepesvari, C., and Wang, M. Model-based reinforcement learning with value-targeted regression. In *LADC*, 2020.
- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. Is Q-learning provably efficient? In *Advances in Neural Information Processing Systems*, pp. 4868–4878, 2018.
- Littman, M. L. Markov games as a framework for multi-agent reinforcement learning. In *Machine learning proceedings 1994*, pp. 157–163. Elsevier, 1994.
- Liu, Q., Yu, T., Bai, Y., and Jin, C. A sharp analysis of model-based reinforcement learning with self-play. *arXiv preprint arXiv:2010.01604*, 2020.
- Modi, A., Jiang, N., Tewari, A., and Singh, S. Sample complexity of reinforcement learning using linearly combined model ensembles. In *International Conference on Artificial Intelligence and Statistics*, pp. 2010–2020. PMLR, 2020.
- Moulin, H. and Vial, J.-P. Strategically zero-sum games: the class of games whose completely mixed equilibria cannot be improved upon. *International Journal of Game Theory*, 7(3-4):201–221, 1978.
- Shapley, L. S. Stochastic games. *Proceedings of the national academy of sciences*, 39(10):1095–1100, 1953.
- Xie, Q., Chen, Y., Wang, Z., and Yang, Z. Learning zero-sum simultaneous-move markov games using function approximation and correlated equilibrium. *arXiv preprint arXiv:2002.07066*, 2020.
- Zhou, D., Gu, Q., and Szepesvari, C. Nearly minimax optimal reinforcement learning for linear mixture markov decision processes. *arXiv preprint arXiv:2012.08507*, 2020.



## A. Proof of Results in Section 4

We let  $\mathbb{P}$  be the distribution over  $(\mathcal{S} \times \mathcal{A}_{\max} \times \mathcal{A}_{\min})^{\mathbb{N}}$  induced by the episodic MG  $M$ , and further denote the sample space  $\Omega = (\mathcal{S} \times \mathcal{A}_{\max} \times \mathcal{A}_{\min})^{\mathbb{N}}$ . Thus, we work with the probability space given by the triplet  $(\Omega, \mathcal{F}, \mathbb{P})$ , where  $\mathcal{F}$  is the product  $\sigma$ -algebra generated by the discrete  $\sigma$ -algebras underlying  $\mathcal{S}$ ,  $\mathcal{A}_{\max}$  and  $\mathcal{A}_{\min}$ .

For  $1 \leq k \leq K$ ,  $1 \leq h \leq H$ , let  $\mathcal{F}_{k,h}$  be the  $\sigma$ -algebra generated by the random variables representing the state-action-action pairs up to and including those that appear stage  $h$  of episode  $k$ . That is,  $\mathcal{F}_{k,h}$  is generated by

$$\begin{aligned} & s_1^1, a_1^1, b_1^1, \dots, s_h^1, a_h^1, b_h^1, \dots, s_H^1, a_H^1, b_H^1, \\ & s_1^2, a_1^2, b_1^2, \dots, s_h^2, a_h^2, b_h^2, \dots, s_H^2, a_H^2, b_H^2, \\ & \vdots \\ & s_1^k, a_1^k, b_1^k, \dots, s_h^k, a_h^k, b_h^k. \end{aligned}$$

### A.1. Proof of Theorem 4.1

Nash-UCRL constructs  $\bar{\theta}_{k,h}^{(0)}$  as the estimator of  $\theta_h^*$  based on linear regression on  $(\phi_{\bar{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k), \bar{V}_{k,h+1}(s_{h+1}^k))$  (as mentioned in section 3). Due to the randomness of  $s_{h+1}^k$ ,  $\bar{\theta}_{k,h}^{(0)}$  can not estimate  $\theta_h^*$  exactly. Therefore Nash-UCRL also constructs an ellipsoid  $\bar{\mathcal{C}}_{k,h}^{(0)}$  centered at  $\bar{\theta}_{k,h}^{(0)}$  as the confidence set, which contains  $\theta_h^*$  with high probability:

$$\bar{\mathcal{C}}_{k,h}^{(0)} := \left\{ \theta : \left\| \left[ \bar{\Sigma}_{k,h}^{(0)} \right]^{1/2} (\theta - \bar{\theta}_{k,h}^{(0)}) \right\|_2 \leq \beta_k^{(0)} \right\}. \quad (\text{A.1})$$

Here  $\left[ \bar{\Sigma}_{k,h}^{(0)} \right]^{1/2}$  is the ‘‘covariance matrix’’ of the context  $\phi_{\bar{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k)$ , and  $\beta_k^{(0)}$  is the radius of the confidence set. We first show that under a specific parameter choice, our constructed confidence sets  $\bar{\mathcal{C}}_{k,h}^{(0)}$  and  $\underline{\mathcal{C}}_{k,h}^{(0)}$  include  $\theta_h^*$  with high probability, and the estimated variances  $\mathbb{V}^{\text{est}} \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k)$  and  $\mathbb{V}^{\text{est}} \underline{V}_{k,h+1}(s_h^k, a_h^k, b_h^k)$  deviate from the true variances by at most the offset terms  $\bar{E}_{k,h}$ ,  $\underline{E}_{k,h}$ .

**Lemma A.1.** Setting  $\beta_k^{(0)}$  in (A.1) and  $\beta_k^{(1)}, \beta_k^{(2)}$  in (3.3) to

$$\begin{aligned} \beta_k^{(0)} &= 16\sqrt{d \log(1+k/\lambda) \log(4k^2 H/\delta)} + 8\sqrt{d} \log(4k^2 H/\delta) + \sqrt{\lambda} B \\ \beta_k^{(1)} &= 16\sqrt{d H^4 \log(1+KH^4/(d\lambda)) \log(4k^2 H/d\delta)} + 8H^2 \log(4k^2 H/\delta) + \sqrt{\lambda} B \\ \beta_k^{(2)} &= 16d\sqrt{\log(1+k/\lambda) \log(4k^2 H/\delta)} + 8\sqrt{d} \log(4k^2 H/\delta) + \sqrt{\lambda} B, \end{aligned}$$

then with probability at least  $1 - 3\delta$ , we have  $\theta_h^* \in \bar{\mathcal{C}}_{k,h}^{(0)} \cap \underline{\mathcal{C}}_{k,h}^{(0)}$ . In addition, we have

$$\begin{aligned} |\mathbb{V}^{\text{est}} \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) - \mathbb{V} \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k)| &\leq \bar{E}_{k,h} \\ |\mathbb{V}^{\text{est}} \underline{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) - \mathbb{V} \underline{V}_{k,h+1}(s_h^k, a_h^k, b_h^k)| &\leq \underline{E}_{k,h} \end{aligned}$$

Let the event  $\mathcal{E}$  denote the event when the conclusion of Lemma A.1 holds. Then Lemma A.1 suggests that  $\mathbb{P}(\mathcal{E}) \geq 1 - 3\delta$ . We introduce another two events in the following lemma.

**Lemma A.2.** Denote events  $\mathcal{E}_1$  and  $\mathcal{E}_2$  as follows

$$\begin{aligned} \mathcal{E}_1 &= \left\{ \forall h' \in [H], \sum_{k=1}^K \sum_{h=h'}^H \left[ [\mathbb{P}_h \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) - [\mathbb{P}_h \underline{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) \right. \right. \\ &\quad \left. \left. - \bar{V}_{k,h+1}(s_{h+1}^k) + \underline{V}_{k,h+1}(s_{h+1}^k) \right] \leq 8H \sqrt{2T \log(H/\delta)} \right\} \\ \mathcal{E}_2 &= \left\{ \sum_{k=1}^K \sum_{h=1}^H \mathbb{V}_h V_{h+1}^{\mu^k}(s_h^k, a_h^k, b_h^k) \leq 3(HT + H^3 \log(1/\delta)) \right\}. \end{aligned}$$

Then we have  $\mathbb{P}(\mathcal{E}_1) \geq 1 - \delta$  and  $\mathbb{P}(\mathcal{E}_2) \geq 1 - \delta$ .

We now present three lemmas based on  $\mathcal{E}, \mathcal{E}_1, \mathcal{E}_2$ . The following lemma shows that  $\bar{Q}$  and  $\bar{V}$  provide the good UCB for the best response of the max-player and  $\underline{Q}$  and  $\underline{V}$  provide the good LCB for the best response of the min-player.

**Lemma A.3.** Suppose the event  $\mathcal{E}$  hold, then we have for any  $s, a, b, k, h$  following inequalities hold,

$$\underline{Q}_{k,h}(s, a, b) - (H - h + 1)\epsilon \leq Q_h^{\pi^k, *}(s, a, b) \leq Q_h^{*, \nu^k}(s, a, b) \leq \bar{Q}_{k,h}(s, a, b) + (H - h + 1)\epsilon,$$

and

$$\underline{V}_{k,h}(s) - (H - h + 2)\epsilon \leq V_h^{\pi^k, *}(s) \leq V_h^{*, \nu^k}(s) \leq \bar{V}_{k,h}(s) + (H - h + 2)\epsilon.$$

**Lemma A.4.** Suppose the events  $\mathcal{E} \cap \mathcal{E}_1$  hold, then we have

$$\begin{aligned} \sum_{k=1}^K [\bar{V}_{k,1}(s_{k,1}) - \underline{V}_{k,1}(s_{k,1})] &\leq 4\beta_K^{(0)} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2 + \sigma_{k,h}^2} \sqrt{2Hd \log(1 + K/\lambda)} \\ &\quad + 8H \sqrt{2T \log(H/\delta)}, \\ \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h[\bar{V}_{k,h+1} - \underline{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) &\leq 4\beta_K^{(0)} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2 + \sigma_{k,h}^2} \sqrt{2H^3d \log(1 + K/\lambda)} \\ &\quad + 8H^2 \sqrt{2T \log(H/\delta)}, \end{aligned}$$

**Lemma A.5.** Suppose the events  $\mathcal{E} \cap \mathcal{E}_2$  hold, then we have

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2 &\leq H^2T/d + 3(HT + H^3 \log(1/\delta)) + 4H \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h[\bar{V}_{k,h+1} - V_{h+1}^{\mu^k}] \\ &\quad + 2\beta_K^{(2)} \sqrt{T} \sqrt{2dH \log(1 + KH^4/(d\lambda))} + 7\beta_K^{(1)} H^2 \sqrt{T} \sqrt{2dH \log(1 + K/\lambda)} \\ \sum_{k=1}^K \sum_{h=1}^H \sigma_{k,h}^2 &\leq H^2T/d + 3(HT + H^3 \log(1/\delta)) + 4H \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h[V_{h+1}^{\mu^k} - \underline{V}_{k,h+1}] \\ &\quad + 2\beta_K^{(2)} \sqrt{T} \sqrt{2dH \log(1 + KH^4/(d\lambda))} + 7\beta_K^{(1)} H^2 \sqrt{T} \sqrt{2dH \log(1 + K/\lambda)} \end{aligned}$$

With all these lemmas, we can now give the proof of Theorem 4.1.

*Proof of Theorem 4.1.* By definition of Regret we have that

$$\begin{aligned} \text{Regret}(K) &= \sum_{k=1}^K V_1^{*, \nu^k}(s_1^k) - \sum_{k=1}^K V_1^{\pi^k, *}(s_1^k) \\ &\leq \sum_{k=1}^K \bar{V}_{k,1}(s_{k,1}) - \sum_{k=1}^K \underline{V}_{k,1}(s_{k,1}) + 4KH\epsilon \\ &\leq 4\beta_K^{(0)} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2 + \sigma_{k,h}^2} \sqrt{2Hd \log(1 + K/\lambda)} + 8H \sqrt{2T \log(H/\delta)} + 4KH\epsilon \\ &= \tilde{O}\left(d\sqrt{H} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2 + \sigma_{k,h}^2} + H\sqrt{T}\right), \end{aligned} \tag{A.2}$$

where the first inequality is by Lemma A.3, the second inequality is by the bound of accumulated difference between the UCB and LCB in Lemma A.4, the last inequality is due to  $\epsilon = O(H/\sqrt{T})$ ,  $\lambda = 1/B^2$  and the choice of  $\beta_K^{(0)} = \tilde{O}(\sqrt{d})$  in Lemma A.1.

Now we bound  $\sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2 + \underline{\sigma}_{k,h}^2$ ,

$$\sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2 + \underline{\sigma}_{k,h}^2 \tag{A.3}$$

$$\begin{aligned} &\leq 2H^2T/d + 6(HT + H^3 \log(1/\delta)) + 4H \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h[\bar{V}_{k,h+1} - \underline{V}_{k,h+1}] \\ &\quad + 4\beta_K^{(2)}\sqrt{T}\sqrt{2dH \log(1 + KH^4/(d\lambda))} + 14\beta_K^{(1)}H^2\sqrt{T}\sqrt{2dH \log(1 + K/\lambda)} \\ &\leq 2H^2T/d + 6(HT + H^3 \log(1/\delta)) \\ &\quad + 4H \left( 4\beta_K^{(0)} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2 + \underline{\sigma}_{k,h}^2} \sqrt{2H^3d \log(1 + K/\lambda)} + 8H^2\sqrt{2T \log(H/\delta)} \right) \\ &\quad + 4\beta_K^{(2)}\sqrt{T}\sqrt{2dH \log(1 + KH^4/(d\lambda))} + 14\beta_K^{(1)}H^2\sqrt{T}\sqrt{2dH \log(1 + K/\lambda)} \\ &= \tilde{O} \left( \sqrt{\sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2 + \underline{\sigma}_{k,h}^2} \sqrt{d^2H^5} + H^2T/d + TH + \sqrt{T}d^{1.5}H^{2.5} + H^3\sqrt{T} \right) \end{aligned} \tag{A.4}$$

where the first inequality is by Lemma A.5, the second inequality is by Lemma A.4 and the last inequality is due to the choice of  $\beta_K^{(0)} = \tilde{O}(\sqrt{d})$  in Lemma A.1,  $\lambda = 1/B^2$ ,

$$\begin{aligned} \beta_K^{(1)} &= 16\sqrt{dH^4 \log(1 + kH^4/d\lambda) \log(8k^2H/\delta)} + 8H^2 \log(8k^2H/\delta) + \sqrt{\lambda}B = \tilde{O}(dH^2) \\ \beta_K^{(2)} &= 16d\sqrt{\log(1 + k/\lambda) \log(8k^2H/\delta)} + 8\sqrt{d} \log(8k^2H/\delta) + \sqrt{\lambda}B = \tilde{O}(d). \end{aligned}$$

Therefore by the fact that  $x \leq a\sqrt{x} + b \Rightarrow x \leq 2a^2 + b$ , (A.3) suggests that

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2 + \underline{\sigma}_{k,h}^2 &= \tilde{O}(d^2H^5 + H^2T/d + TH + \sqrt{T}d^{1.5}H^{2.5} + H^3\sqrt{T}) \\ &= \tilde{O}(d^2H^5 + d^4H^3 + TH + H^2T/d), \end{aligned} \tag{A.5}$$

where the inequality holds by  $\sqrt{T}d^{1.5}H^{2.5} \leq (TH^2/4d + d^4H^3)/2$  and  $H^3\sqrt{T} \leq (d^2H^5 + H^2T/d)/2$ . Plugging (A.5) into (A.2) we have

$$\text{Regret}(M_{\theta^*}, K) = \tilde{O}(\sqrt{d^2H^2 + dH^3\sqrt{T}} + d^2H^3 + d^3H^2),$$

which finishing the proof.  $\square$

## A.2. Proof of Theorem 4.3

*Proof of Theorem 4.3.* For any algorithm, we need to construct a hard-to-learn episodic, B-bounded linear mixture Markov game. We make the min-player dummy: the action of the min-player won't affect the transition ability or reward function. So there exists  $\mathbb{P}_h(\cdot|\cdot, \cdot)$  and  $\tilde{r}_h(\cdot, \cdot)$  such that for any state-action-action-state pair  $s', a, b, s$  we have that  $\mathbb{P}_h(s'|s, a, b) = \tilde{\mathbb{P}}_h(s'|s, a)$  and  $r_h(s, a, b) = \tilde{r}_h(s, a)$ . Thus we can get a new MDP  $\tilde{M}(\mathcal{S}, \mathcal{A}_{\max}, H, \{\tilde{r}_h\}, \{\tilde{\mathbb{P}}_h\})$ . We further have  $V_h^{\pi, *}(s) = \tilde{V}_h^{\pi}(s)$  and  $V_h^{*, \nu}(s) = \tilde{V}_h^{*}(s)$ . The regret of two-player game can be reduced to the standard regret for single agent reinforcement learning setting. In particular,

$$\text{Regret}(M_{\theta^*}, K) = \sum_{k=1}^K V_1^{*, \nu^k}(s_1^k) - \sum_{k=1}^K V_1^{\pi^k, *}(s_1^k)$$



$$= \sum_{k=1}^K \tilde{V}_1^*(s_1^k) - \sum_{k=1}^K \tilde{V}_1^{\pi^k}(s_1^k).$$

Notice that  $\tilde{r}_h \in [-1, 1]$  rather than  $[0, 1]$ , we can shift the reward by  $(1 + \tilde{r}_h)/2$  to make it standard if necessary. Now recall the Theorem 5.6 in (Zhou et al., 2020), there exists an episodic,  $B$ -bounded linear mixture MDP  $\tilde{M}(\mathcal{S}, \mathcal{A}_{\max}, H, \{\tilde{r}_h\}, \{\tilde{\mathbb{P}}_h\})$  with feature  $\tilde{\phi}(\cdot, \cdot)$  parameterized by  $\Theta = (\theta_1, \dots, \theta_H)$  such that the expected regret is lower bounded as follows:

$$\mathbb{E}_{\Theta} \text{Regret}(\tilde{M}_{\Theta}, K) \geq \Omega(dH\sqrt{T}),$$

where  $T = KH$  and  $\mathbb{E}_{\Theta}$  denotes the expectation over the probability distribution generated by the interconnection of the algorithm and the MDP.

Now we only need to extend the MDP feature  $\tilde{\phi}(\cdot, \cdot)$  to the Markov game feature  $\phi(\cdot, \cdot, \cdot)$ . In particular, we set

$$\phi(s'|s, a, b) = \tilde{\phi}(s'|s, a), \forall s' \in \mathcal{S}, s \in \mathcal{S}, a \in \mathcal{A}_{\max}, b \in \mathcal{A}_{\min},$$

then we know that  $\phi(\cdot, \cdot, \cdot)$  satisfies (2.2) because by the definition of linear mixture MDP in (Zhou et al., 2020), we know that  $\tilde{\phi}(\cdot, \cdot)$  satisfies for any bounded function  $V : \mathcal{S} \rightarrow [0, 1]$ ,

$$\|\tilde{\phi}_V(s, a)\|_2 \leq 1,$$

where  $\tilde{\phi}_V(s, a) = \sum_{s' \in \mathcal{S}} \tilde{\phi}(s'|s, a)V(s')$ .

□

## B. Proof of Lemmas in Appendix A

### B.1. Proof of Lemma A.1

For simplicity we denote the following confident sets:

$$\begin{aligned} \bar{\mathcal{C}}_{k,h}^{(0)} &= \left\{ \theta : \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{1/2} (\theta - \bar{\theta}_{k,h}^{(0)}) \right\|_2 \leq \beta_k^{(0)} \right\}, \underline{\mathcal{C}}_{k,h}^{(0)} = \left\{ \theta : \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{1/2} (\theta - \underline{\theta}_{k,h}^{(0)}) \right\|_2 \leq \beta_k^{(0)} \right\}, \\ \bar{\mathcal{C}}_{k,h}^{(1)} &= \left\{ \theta : \left\| [\bar{\Sigma}_{k,h}^{(1)}]^{1/2} (\theta - \bar{\theta}_{k,h}^{(1)}) \right\|_2 \leq \beta_k^{(1)} \right\}, \underline{\mathcal{C}}_{k,h}^{(1)} = \left\{ \theta : \left\| [\underline{\Sigma}_{k,h}^{(1)}]^{1/2} (\theta - \underline{\theta}_{k,h}^{(1)}) \right\|_2 \leq \beta_k^{(1)} \right\}, \\ \bar{\mathcal{C}}_{k,h}^{(2)} &= \left\{ \theta : \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{1/2} (\theta - \bar{\theta}_{k,h}^{(0)}) \right\|_2 \leq \beta_k^{(2)} \right\}, \underline{\mathcal{C}}_{k,h}^{(2)} = \left\{ \theta : \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{1/2} (\theta - \underline{\theta}_{k,h}^{(0)}) \right\|_2 \leq \beta_k^{(2)} \right\}. \end{aligned}$$

By the selection  $\beta_k^{(0)} < \beta_k^{(2)}$  in Lemma A.1, we have that  $\bar{\mathcal{C}}_{k,h}^{(0)} \subset \bar{\mathcal{C}}_{k,h}^{(2)}$  and  $\underline{\mathcal{C}}_{k,h}^{(0)} \subset \underline{\mathcal{C}}_{k,h}^{(2)}$ . We first use standard self-normalized tail inequality to show that  $\theta_h^*$  is included in  $\bar{\mathcal{C}}_{k,h}^{(1)} \cap \bar{\mathcal{C}}_{k,h}^{(2)}$  with high probability. Based on that we can further decrease  $\beta_k^{(2)}$  to  $\beta_k^{(1)}$  without significantly increasing the probability of the bad event when  $\theta_h^* \notin \bar{\mathcal{C}}_{k,h}^{(0)}$  or  $\theta_h^* \notin \underline{\mathcal{C}}_{k,h}^{(0)}$ .

We start with the following Bernstein-type self-normalized concentration inequality.

**Lemma B.1** (Bernstein inequality in (Zhou et al., 2020)). Let  $\{\mathcal{G}_t\}_{t=1}^{\infty}$  be a filtration,  $\{\mathbf{x}_t, \eta_t\}_{t \geq 1}$  a stochastic process so that  $\mathbf{x}_t \in \mathbb{R}^d$  is  $\mathcal{G}_t$ -measurable and  $\eta_t \in \mathbb{R}$  is  $\mathcal{G}_{t+1}$ -measurable. Fix  $R, L, \sigma, \lambda > 0$ ,  $\mu^* \in \mathbb{R}^d$ . For  $t \geq 1$  let  $y_t = \langle \mu^*, \mathbf{x}_t \rangle + \eta_t$  and suppose that  $\eta_t, \mathbf{x}_t$  also satisfy

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

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

$$\forall t > 0, \left\| \sum_{i=1}^t \mathbf{x}_i \eta_i \right\|_{\mathbf{Z}_t^{-1}} \leq \beta_t, \|\mu_t - \mu^*\|_{\mathbf{Z}_t} \leq \beta_t + \sqrt{\lambda} \|\mu^*\|_2, \quad (\text{B.1})$$

where for  $t \geq 1$ ,  $\mu_t = \mathbf{Z}_t^{-1} \mathbf{b}_t$ ,  $\mathbf{Z}_t = \lambda \mathbf{I} + \sum_{i=1}^t \mathbf{x}_i \mathbf{x}_i^\top$ ,  $\mathbf{b}_t = \sum_{i=1}^t y_i \mathbf{x}_i$  and

$$\beta_t = 8\sigma \sqrt{d \log(1 + tL^2/(d\lambda)) \log(4t^2/\delta)} + 4R \log(4t^2/\delta).$$

**Lemma B.2.**

$$\begin{aligned}
 & |\mathbb{V}^{\text{est}} \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) - \mathbb{V} \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k)| \\
 & \leq \min \left\{ H^2, \left\| [\bar{\Sigma}_{k,h}^{(1)}]^{-1/2} \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \right\|_2 \left\| [\bar{\Sigma}_{k,h}^{(1)}]^{1/2} (\bar{\theta}_{k,h}^{(1)} - \theta_h^*) \right\|_2 \right\} \\
 & \quad + \min \left\{ H^2, 2H \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \right\|_2 \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{1/2} (\bar{\theta}_{k,h}^{(0)} - \theta_h^*) \right\|_2 \right\}.
 \end{aligned}$$

$$\begin{aligned}
 & |\mathbb{V}^{\text{est}} \underline{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) - \mathbb{V} \underline{V}_{k,h+1}(s_h^k, a_h^k, b_h^k)| \\
 & \leq \min \left\{ H^2, \left\| [\underline{\Sigma}_{k,h}^{(1)}]^{-1/2} \phi_{\underline{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \right\|_2 \left\| [\underline{\Sigma}_{k,h}^{(1)}]^{1/2} (\underline{\theta}_{k,h}^{(1)} - \theta_h^*) \right\|_2 \right\} \\
 & \quad + \min \left\{ H^2, 2H \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \right\|_2 \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{1/2} (\underline{\theta}_{k,h}^{(0)} - \theta_h^*) \right\|_2 \right\}.
 \end{aligned}$$

*Proof.* For simplicity, we only prove the results for the max-player.

By the triangle inequality we have that

$$\begin{aligned}
 & |\mathbb{V}^{\text{est}} \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) - \mathbb{V} \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k)| \\
 & \leq \underbrace{\left| \langle \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \theta_h^* \rangle - [\langle \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \bar{\theta}_{k,h}^{(1)} \rangle]_{[0, H^2]} \right|}_{I_1} \\
 & \quad + \underbrace{\left| (\langle \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \theta_h^* \rangle)^2 - [\langle \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \bar{\theta}_{k,h}^{(0)} \rangle]_{[-H, H]}^2 \right|}_{I_2}. \tag{B.2}
 \end{aligned}$$

We first bound  $I_1$ . Because  $\langle \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \theta_h^* \rangle \in [0, H^2]$ , we have that

$$\begin{aligned}
 I_1 & \leq \left| \langle \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \theta_h^* \rangle - \langle \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \bar{\theta}_{k,h}^{(1)} \rangle \right| \\
 & \leq \left\| [\bar{\Sigma}_{k,h}^{(1)}]^{-1/2} \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \right\|_2 \left\| [\bar{\Sigma}_{k,h}^{(1)}]^{1/2} (\bar{\theta}_{k,h}^{(1)} - \theta_h^*) \right\|_2,
 \end{aligned}$$

where the first inequality is by the property of projection, the second inequality holds due to Cauchy-Schwarz. We also have that  $I_1 \leq H^2$  since both terms in  $I_1$  belongs to the interval  $[0, H^2]$ , so we have that

$$I_1 \leq \min \left\{ H^2, \left\| [\bar{\Sigma}_{k,h}^{(1)}]^{-1/2} \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \right\|_2 \left\| [\bar{\Sigma}_{k,h}^{(1)}]^{1/2} (\bar{\theta}_{k,h}^{(1)} - \theta_h^*) \right\|_2 \right\}, \tag{B.3}$$

For the term  $I_2$ ,

$$\begin{aligned}
 I_2 & = \left| \langle \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \theta_h^* \rangle - [\langle \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \bar{\theta}_{k,h}^{(0)} \rangle]_{[-H, H]} \right| \\
 & \quad \cdot \left| \langle \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \theta_h^* \rangle + [\langle \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \bar{\theta}_{k,h}^{(0)} \rangle]_{[-H, H]} \right| \\
 & \leq 2H \left| \langle \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \theta_h^* \rangle - \langle \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \bar{\theta}_{k,h}^{(0)} \rangle \right| \\
 & \leq 2H \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \right\|_2 \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{1/2} (\bar{\theta}_{k,h}^{(0)} - \theta_h^*) \right\|_2,
 \end{aligned}$$

where the first inequality holds since both terms in this line lies in  $[-H, H]$ , the second inequality holds since the Cauchy-Schwarz inequality. We also have that  $I_2 \leq H^2$ , so we have that

$$I_2 \leq \min \left\{ H^2, 2H \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \right\|_2 \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{1/2} (\bar{\theta}_{k,h}^{(0)} - \theta_h^*) \right\|_2 \right\}. \tag{B.4}$$

Plugging (B.4) and (B.3) into (B.2) gets

$$|\mathbb{V}^{\text{est}} \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) - \mathbb{V} \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k)|$$

$$\leq \min \left\{ H^2, \left\| [\bar{\Sigma}_{k,h}^{(1)}]^{-1/2} \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \right\|_2 \left\| [\bar{\Sigma}_{k,h}^{(1)}]^{1/2} (\bar{\theta}_{k,h}^{(1)} - \theta_h^*) \right\|_2 \right\} \\ + \min \left\{ H^2, 2H \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \right\|_2 \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{1/2} (\bar{\theta}_{k,h}^{(0)} - \theta_h^*) \right\|_2 \right\}.$$

□

Now we present the Proof of Lemma A.1.

*Proof of Lemma A.1.* For simplicity, we only prove the results for the max-player. Fix  $h \in [H]$ .

We first show that with probability at least  $1 - \delta/(2H)$ ,  $\left\| [\bar{\Sigma}_{k,h}^{(0)}]^{1/2} (\bar{\theta}_{k,h}^{(0)} - \theta_h^*) \right\|_2 \leq \beta_k^{(2)}$ . To show this, we apply Lemma B.1. Let  $\mathbf{x}_i = \bar{\sigma}_{i,h}^{-1} \phi_{\bar{V}_{i,h+1}}(s_h^i, a_h^i, b_h^i)$  and  $\eta_i = \bar{\sigma}_{i,h}^{-1} \bar{V}_{i,h+1}(s_{h+1}^i) - \bar{\sigma}_{i,h}^{-1} \langle \phi_{\bar{V}_{i,h+1}}(s_h^i, a_h^i, b_h^i), \theta_h^* \rangle$ ,  $\mathcal{G}_i = \mathcal{F}_{i,h}$ ,  $\mu^* = \theta_h^*$ ,  $y_i = \langle \mu^*, \mathbf{x}_i \rangle + \eta_i$ ,  $\mathbf{Z}_i = \lambda \mathbf{I} + \sum_{i'=1}^i \mathbf{x}_{i'} \mathbf{x}_{i'}^\top$ ,  $\mathbf{b}_i = \sum_{i'=1}^i \mathbf{x}_{i'} y_{i'}$  and  $\mu_i = \mathbf{Z}_i^{-1} \mathbf{b}_i$ . Then it can be verified that  $y_i = \bar{\sigma}_{i,h}^{-1} \bar{V}_{i,h+1}(s_{h+1}^i)$  and  $\mu_i = \bar{\theta}_{i+1,h}^{(0)}$ . Moreover, we have that

$$\|\mathbf{x}_i\|_2 \leq \bar{\sigma}_{i,h}^{-1} H \leq \sqrt{d}, \quad |\eta_i| \leq \bar{\sigma}_{i,h}^{-1} 2H \leq 2\sqrt{d}, \quad \mathbb{E}[\eta_i | \mathcal{G}_i] = 0, \quad \mathbb{E}[\eta_i^2 | \mathcal{G}_i] \leq 4d,$$

where we apply  $\|\phi_{\bar{V}_{i,h+1}}(\cdot, \cdot, \cdot)\|_2 \leq H$ ,  $\bar{V}_{i,h+1} \in [-H, H]$  and  $\bar{\sigma}_{i,h} \geq H/\sqrt{d}$ . Since we also have that  $\mathbf{x}_i$  is  $\mathcal{G}_i$  measurable and  $\eta_i$  is  $\mathcal{G}_{i+1}$  measurable, by Lemma B.1, we obtain that with probability at least  $1 - \delta/(2H)$ , for all  $k \leq K$ ,  $\left\| [\bar{\Sigma}_{k,h}^{(0)}]^{1/2} (\bar{\theta}_{k,h}^{(0)} - \theta_h^*) \right\|_2$  is bounded by

$$16d\sqrt{\log(1+k/\lambda)\log(8k^2H/\delta)} + 8\sqrt{d}\log(8k^2H/\delta) + \sqrt{\lambda}B = \beta_k^{(2)}, \quad (\text{B.5})$$

implying that with probability at least  $1 - \delta/(2H)$ , for any  $k \leq K$ ,  $\theta_h^* \in \bar{\mathcal{C}}_{k,h}^{(2)}$ .

An argument, which is analogous to the one just used (except that now the range of the “noise” matches the range of “squared values” and is thus bounded by  $H^2$ , rather than being bounded by  $\sqrt{d}$ ) gives that with probability at least  $1 - \delta/(2H)$ , for any  $k \leq K$  we have  $\left\| [\bar{\Sigma}_{k,h}^{(1)}]^{1/2} (\bar{\theta}_{k,h}^{(1)} - \theta_h^*) \right\|_2$  bounded by

$$16\sqrt{dH^4\log(1+kH^4/(d\lambda))\log(8k^2H/\delta)} + 8H^2\log(8k^2H/\delta) + \sqrt{\lambda}B = \beta_k^{(1)}, \quad (\text{B.6})$$

implying that with probability at least  $1 - \delta/(2H)$ , for any  $k \leq K$ ,  $\theta_h^* \in \bar{\mathcal{C}}_{k,h}^{(1)}$ .

We now show that  $\theta_h^* \in \bar{\mathcal{C}}_{k,h}^{(0)}$  with high probability. We again apply Lemma B.1. Let  $\mathbf{x}_i = \bar{\sigma}_{i,h}^{-1} \phi_{\bar{V}_{i,h+1}}(s_h^i, a_h^i, b_h^i)$  and

$$\eta_i = \bar{\sigma}_{i,h}^{-1} \mathbb{1}\{\theta_h^* \in \bar{\mathcal{C}}_{i,h}^{(1)} \cap \bar{\mathcal{C}}_{i,h}^{(2)}\} [\bar{V}_{i,h+1}(s_{h+1}^i) - \langle \phi_{\bar{V}_{i,h+1}}(s_h^i, a_h^i, b_h^i), \theta_h^* \rangle],$$

$\mathcal{G}_i = \mathcal{F}_{i,h}$ ,  $\mu^* = \theta_h^*$ ,  $y_i = \langle \mu^*, \mathbf{x}_i \rangle + \eta_i$ ,  $\mathbf{Z}_i = \lambda \mathbf{I} + \sum_{i'=1}^i \mathbf{x}_{i'} \mathbf{x}_{i'}^\top$ ,  $\mathbf{b}_i = \sum_{i'=1}^i \mathbf{x}_{i'} y_{i'}$  and  $\mu_i = \mathbf{Z}_i^{-1} \mathbf{b}_i$ . Still we have that  $\|\mathbf{x}_i\|_2 \leq \bar{\sigma}_{i,h}^{-1} H \leq \sqrt{d}$ . Because  $\mathbb{1}\{\theta_h^* \in \bar{\mathcal{C}}_{i,h}^{(1)} \cap \bar{\mathcal{C}}_{i,h}^{(2)}\}$  is  $\mathcal{G}_i$ -measurable, we have  $\mathbb{E}[\eta_i | \mathcal{G}_i] = 0$ . We also have  $|\eta_i| \leq \bar{\sigma}_{i,h}^{-1} 2H \leq 2\sqrt{d}$  since  $|\bar{V}_{i,h+1}(\cdot)| \leq H$  and  $\bar{\sigma}_{i,h} \geq H/\sqrt{d}$ . To get better bound  $\beta_k^{(0)}$  rather than  $\beta_k^{(2)}$  in (B.5), we need more careful computation of  $\mathbb{E}[\eta_i^2 | \mathcal{G}_i]$  as follows,

$$\mathbb{E}[\eta_i^2 | \mathcal{G}_i] = \bar{\sigma}_{i,h}^{-2} \mathbb{1}\{\theta_h^* \in \bar{\mathcal{C}}_{i,h}^{(1)} \cap \bar{\mathcal{C}}_{i,h}^{(2)}\} [\mathbb{V}_h \bar{V}_{i,h+1}](s_h^i, a_h^i, b_h^i) \\ \leq \bar{\sigma}_{i,h}^{-2} \mathbb{1}\{\theta_h^* \in \bar{\mathcal{C}}_{i,h}^{(1)} \cap \bar{\mathcal{C}}_{i,h}^{(2)}\} \left[ [\mathbb{V}_{i,h}^{\text{est}} \bar{V}_{i,h+1}](s_h^i, a_h^i, b_h^i) \right. \\ \left. + \min \left\{ H^2, \left\| [\bar{\Sigma}_{i,h}^{(1)}]^{-1/2} \phi_{\bar{V}_{i,h+1}^2}(s_h^i, a_h^i, b_h^i) \right\|_2 \left\| [\bar{\Sigma}_{i,h}^{(1)}]^{1/2} (\bar{\theta}_{i,h}^{(1)} - \theta_h^*) \right\|_2 \right\} \right. \\ \left. + \min \left\{ H^2, 2H \left\| [\bar{\Sigma}_{i,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{i,h+1}}(s_h^i, a_h^i, b_h^i) \right\|_2 \left\| [\bar{\Sigma}_{i,h}^{(0)}]^{1/2} (\bar{\theta}_{i,h}^{(0)} - \theta_h^*) \right\|_2 \right\} \right] \\ \leq \bar{\sigma}_{i,h}^{-2} \left[ [\mathbb{V}_{i,h}^{\text{est}} \bar{V}_{i,h+1}](s_h^i, a_h^i, b_h^i) + \min \left\{ H^2, \beta_i^{(1)} \left\| [\bar{\Sigma}_{i,h}^{(1)}]^{-1/2} \phi_{\bar{V}_{i,h+1}^2}(s_h^i, a_h^i, b_h^i) \right\|_2 \right\} \right]$$

$$\begin{aligned}
 & + \min \left\{ H^2, 2H\beta_i^{(2)} \left\| [\Sigma_{i,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{i,h+1}}(s_h^i, a_h^i, b_h^i) \right\|_2 \right\} \\
 & = 1,
 \end{aligned}$$

where the first inequality holds due to Lemma B.2, the second inequality holds due to the indicator function, the last equality holds due to the definition of  $\bar{\sigma}_{i,h}$ . Then, by Lemma B.1, with probability at least  $1 - \delta/(2H)$ ,  $\forall k \leq K$ ,

$$\|\mu_k - \mu^*\|_{\mathbf{z}_i} \leq 16\sqrt{d \log(1 + k/\lambda) \log(8k^2 H/\delta)} + 8\sqrt{d} \log(8k^2 H/\delta) + \sqrt{\lambda} B = \beta_k^{(0)}, \quad (\text{B.7})$$

where the equality uses the definition of  $\beta_k^{(0)}$ . Let  $\mathcal{E}'$  be the event when  $\theta_h^* \in \cap_{k \leq K} \bar{\mathcal{C}}_{k,h}^{(1)} \cap \bar{\mathcal{C}}_{k,h}^{(2)}$  and (B.7) hold. By the union bound,  $\mathbb{P}(\mathcal{E}') \geq 1 - 3\delta/(2H)$ .

We now show that  $\theta_h^* \in \bar{\mathcal{C}}_{k,h}^{(0)}$  holds on  $\mathcal{E}'$ . For this note that on  $\mathcal{E}'$ , for any  $k \leq K$ ,  $\mu_k = \bar{\theta}_{k+1,h}^{(0)}$  and for any  $i \leq K$ ,

$$\begin{aligned}
 y_i &= \bar{\sigma}_{i,h}^{-1} (\langle \theta_h^*, \phi_{\bar{V}_{i,h+1}}(s_h^i, a_h^i, b_h^i) \rangle + \mathbb{1}_{\{\theta_h^* \in \bar{\mathcal{C}}_{i,h}^{(1)} \cap \bar{\mathcal{C}}_{i,h}^{(2)}\}} [\bar{V}_{i,h+1}(s_{h+1}^i) \\
 & \quad - \langle \phi_{\bar{V}_{i,h+1}}(s_h^i, a_h^i, b_h^i), \theta^* \rangle]) \\
 &= \bar{\sigma}_{i,h}^{-1} \bar{V}_{i,h+1}(s_{h+1}^i),
 \end{aligned}$$

which implies the claim. Therefore, by the definition of  $\mathcal{C}_{k,h}^{(0)}$ , we get that on  $\mathcal{E}'$ ,  $\theta_h^* \in \cap_{k \leq K} \bar{\mathcal{C}}_{k,h}^{(0)} \cap \bar{\mathcal{C}}_{k,h}^{(1)}$ . Moreover,  $\mathbb{P}(\mathcal{E}') \geq 1 - 3\delta/(2H)$ . Finally, taking union bound over  $h$  shows that with probability at least  $1 - 3\delta/2$ , for all  $h \in [H]$ ,

$$\theta_h^* \in \cap_{k \leq K} \bar{\mathcal{C}}_{k,h}^{(1)} \cap \bar{\mathcal{C}}_{k,h}^{(2)} \quad (\text{B.8})$$

To finish our proof, it is thus sufficient to show that on the event when (B.8) holds, it also holds that

$$|[\mathbb{V}_{k,h}^{\text{est}} \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) - [\mathbb{V}_h \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k)| \leq \bar{E}_{k,h}.$$

However, by the definition of  $\bar{E}_{k,h}$ , this is immediate from substituting (B.5), (B.6) into Lemma B.2.  $\square$

## B.2. Proof of Lemma A.2

We first present the Azuma-Hoeffding inequality:

**Lemma B.3** (Azuma-Hoeffding inequality, (Azuma, 1967)). Let  $M > 0$  be a constant. Let  $\{x_i\}_{i=1}^n$  be a martingale difference sequence with respect to a filtration  $\{\mathcal{G}_i\}_i$  ( $\mathbb{E}[x_i|\mathcal{G}_i] = 0$  a.s. and  $x_i$  is  $\mathcal{G}_{i+1}$ -measurable) such that for all  $i \in [n]$ ,  $|x_i| \leq M$  holds almost surely. Then, for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$ , we have

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

*Proof of Lemma A.2.* To prove  $\mathbb{P}(\mathcal{E}_1) \geq 1 - \delta$ , we apply the Azuma-Hoeffding inequality (Lemma B.3). Fix  $h' \in H$ , set  $x_{k,h} = [\mathbb{P}_h \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) - [\mathbb{P}_h \underline{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) - [\bar{V}_{k,h+1}(s_{h+1}^k) - \underline{V}_{k,h+1}(s_{h+1}^k)]$ .  $x_{1,h'}, \dots, x_{1,H}, x_{2,h'}, \dots, x_{2,H}, \dots, x_{K,h'}, \dots, x_{K,H}$  forms a martingale difference sequence of which the absolute value is bounded by  $8H$  and length no greater than  $T = KH$ . Thus with probability at least  $1 - \delta/H$ , we have

$$\begin{aligned}
 & \sum_{k=1}^K \sum_{h=h'}^H \left[ [\mathbb{P}_h \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) - [\mathbb{P}_h \underline{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) - \bar{V}_{k,h+1}(s_{h+1}^k) + \underline{V}_{k,h+1}(s_{h+1}^k) \right] \\
 & \leq 8H \sqrt{2T \log(H/\delta)}.
 \end{aligned}$$

Take union bound for  $h' \in [H]$ , we get  $\mathbb{P}(\mathcal{E}_1) \geq 1 - \delta$ .

$\mathbb{P}(\mathcal{E}_2) \geq 1 - \delta$  holds due to the Lemma C.5 in (Jin et al., 2018) or Lemma 8 in (Azar et al., 2017).  $\square$

### B.3. Proof of Lemma A.3

Following Lemma directly from the definition of  $\epsilon$ -CCE,

**Lemma B.4.** For each  $(k, h, s)$ ,  $\mu_h^k(\cdot, \cdot | s)$ ,  $\pi_h^k(\cdot | s)$ ,  $\nu_h^k(\cdot | s)$  satisfy that

$$\begin{aligned} \mathbb{E}_{(a,b) \sim \mu_h^k(\cdot, \cdot | s)} [\bar{Q}_{k,h}(s, a, b)] &\geq \mathbb{E}_{b \sim \nu_h^k(s)} [\bar{Q}_{k,h}(s, a', b)] - \epsilon, \forall a' \in \mathcal{A}_{\max} \\ \mathbb{E}_{(a,b) \sim \mu_h^k(\cdot, \cdot | s)} [Q_{k,h}(s, a, b)] &\leq \mathbb{E}_{a \sim \pi_h^k(s)} [Q_{k,h}(s, a, b')] - \epsilon, \forall b' \in \mathcal{A}_{\min} \end{aligned}$$

*Proof of Lemma A.3.* For simplicity, we only prove the following UCB by induction,

$$Q_h^{*,\nu^k}(s, a, b) \leq \bar{Q}_{k,h}(s, a, b) + (H - h + 1)\epsilon, V_h^{*,\nu^k}(s) \leq \bar{V}_{k,h}(s) + (H - h + 2)\epsilon. \quad (\text{B.9})$$

The base case  $h = H + 1$  holds trivially since the terminal cost is zero. Now we assume that the bounds (B.9) holds for step  $h + 1$ . That is,

$$Q_{h+1}^{*,\nu^k}(s, a, b) \leq \bar{Q}_{k,h+1}(s, a, b) + (H - h)\epsilon, V_{h+1}^{*,\nu^k}(s) \leq \bar{V}_{k,h+1}(s) + (H - h + 1)\epsilon. \quad (\text{B.10})$$

If  $\bar{Q}_{k,h}(s, a, b) \geq H$ , then it is obvious to have  $Q_h^{*,\nu^k}(s, a, b) \leq \bar{Q}_{k,h}(s, a, b) + (H - h)\epsilon$ , otherwise we have that

$$\begin{aligned} &\bar{Q}_{k,h}(s, a, b) - Q_h^{*,\nu^k}(s, a, b) \\ &= \langle \bar{\theta}_{k,h}^{(0)}, \phi_{\bar{V}_{k,h+1}} \rangle + \beta_k^{(0)} \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} \right\|_2 - \langle \theta_h^*, \phi_{\bar{V}_{k,h+1}} \rangle \\ &\quad + \mathbb{P}_h \bar{V}_{k,h+1}(s, a, b) - \mathbb{P}_h V_{h+1}^{*,\nu^k}(s) \\ &\geq \beta_k^{(0)} \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} \right\|_2 - \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{1/2} (\bar{\theta}_{k,h}^{(0)} - \theta_h^*) \right\|_2 \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} \right\|_2 \\ &\quad + \mathbb{P}_h \bar{V}_{k,h+1}(s) - \mathbb{P}_h V_{h+1}^{*,\nu^k}(s) \\ &\geq \mathbb{P}_h \bar{V}_{k,h+1}(s) - \mathbb{P}_h V_{h+1}^{*,\nu^k}(s) \\ &\geq -(H - h + 1)\epsilon, \end{aligned} \quad (\text{B.11})$$

where the first inequality holds due to Cauchy-Schwarz inequality, the second inequality holds since the assumption that  $\theta_h^* \in \bar{\mathcal{C}}_{k,h}^{(0)}$  on event  $\mathcal{E}$ , the third inequality holds by the induction assumption. Finally, let  $\text{br}(\nu_h^k(\cdot | s))$  denote the best response to  $\nu_h^k(\cdot | s)$  with respect to  $Q_h^{*,\nu^k}(s, \cdot, \cdot)$  such that

$$\text{br}(\nu_h^k(\cdot | s)) = \underset{\sigma \in \Delta_{\mathcal{A}_{\max}}}{\operatorname{argmax}} \mathbb{E}_{a \sim \sigma, b \sim \nu_h^k(\cdot | s)} Q_h^{*,\nu^k}(s, a, b).$$

Then we have that

$$\begin{aligned} \bar{V}_{k,h}(s) &= \mathbb{E}_{(a,b) \sim \mu_h^k(\cdot, \cdot | s)} [\bar{Q}_{k,h}(s, a, b)] \\ &\geq \mathbb{E}_{a' \sim \text{br}(\nu_h^k(\cdot | s)), b \sim \nu_h^k(\cdot | s)} [\bar{Q}_{k,h}(s, a', b)] - \epsilon \\ &\geq \mathbb{E}_{a' \sim \text{br}(\nu_h^k(\cdot | s)), b \sim \nu_h^k(\cdot | s)} [Q_h^{*,\nu^k}(s, a', b)] - (H - h + 2)\epsilon \\ &= V_h^{*,\nu^k}(s) - (H - h + 2)\epsilon, \end{aligned}$$

where the the first equality is by the property of  $\epsilon$ -CCE in Lemma B.4, the second inequality is by (B.11), the last inequality is due to the Bellman equation. Therefore, our proof ends.  $\square$

### B.4. Proof of Lemma A.4

*Proof of Lemma A.4.*

$$\bar{V}_{k,h}(s_h^k) - \underline{V}_{k,h}(s_h^k)$$

$$\begin{aligned}
 &= \langle \bar{\theta}_{k,h}^{(0)}, \phi_{\bar{V}_{k,h+1}} \rangle + \beta_k^{(0)} \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} \right\|_2 - \langle \underline{\theta}_{k,h}^{(0)}, \phi_{\underline{V}_{k,h+1}} \rangle + \beta_k^{(0)} \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}} \right\|_2 \\
 &= \langle \theta_h^*, \phi_{\bar{V}_{k,h+1}} \rangle + \langle \bar{\theta}_{k,h}^{(0)} - \theta_h^*, \phi_{\bar{V}_{k,h+1}} \rangle + \beta_k^{(0)} \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} \right\|_2 \\
 &\quad - \langle \theta_h^*, \phi_{\underline{V}_{k,h+1}} \rangle - \langle \underline{\theta}_{k,h}^{(0)} - \theta_h^*, \phi_{\underline{V}_{k,h+1}} \rangle + \beta_k^{(0)} \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}} \right\|_2 \\
 &\leq \langle \theta_h^*, \phi_{\bar{V}_{k,h+1}} \rangle + \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{1/2} (\bar{\theta}_{k,h}^{(0)} - \theta_h^*) \right\|_2 \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} \right\|_2 + \beta_k^{(0)} \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} \right\|_2 \\
 &\quad - \langle \theta_h^*, \phi_{\underline{V}_{k,h+1}} \rangle + \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{1/2} (\underline{\theta}_{k,h}^{(0)} - \theta_h^*) \right\|_2 \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}} \right\|_2 + \beta_k^{(0)} \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}} \right\|_2 \\
 &\leq [\mathbb{P}_h \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) + 2\beta_k^{(0)} \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} \right\|_2 \\
 &\quad - [\mathbb{P}_h \underline{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) + 2\beta_k^{(0)} \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}} \right\|_2,
 \end{aligned}$$

where the first equation is by the definition of  $\bar{V}_{k,h}(s_h^k)$ ,  $\underline{V}_{k,h}(s_h^k)$  and the second inequality is due to Cauchy-Schwarz inequality, the last inequality is by  $\theta_h^* \in \bar{\mathcal{C}}_{k,h}^{(0)} \cap \underline{\mathcal{C}}_{k,h}^{(0)}$  on the event  $\mathcal{E}$ .

Meanwhile, since  $\bar{V}_{k,h}(s_h^k) - \underline{V}_{k,h}(s_h^k) \leq 2H$ , we have that

$$\begin{aligned}
 \bar{V}_{k,h}(s_h^k) - \underline{V}_{k,h}(s_h^k) &\leq \min \left\{ 2H, [\mathbb{P}_h \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) + 2\beta_k^{(0)} \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} \right\|_2 \right. \\
 &\quad \left. - [\mathbb{P}_h \underline{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) + 2\beta_k^{(0)} \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}} \right\|_2 \right\} \\
 &\leq \min \left\{ 4H, 2\beta_k^{(0)} \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} \right\|_2 + 2\beta_k^{(0)} \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}} \right\|_2 \right\} \\
 &\quad + [\mathbb{P}_h \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) - [\mathbb{P}_h \underline{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) \\
 &\leq \min \left\{ 4H, 2\beta_k^{(0)} \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} \right\|_2 \right\} \\
 &\quad + \min \left\{ 4H, 2\beta_k^{(0)} \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}} \right\|_2 \right\} \\
 &\quad + [\mathbb{P}_h \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) - [\mathbb{P}_h \underline{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) \\
 &\leq 2\beta_k^{(0)} \bar{\sigma}_{k,h} \min \left\{ 1, \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} / \bar{\sigma}_{k,h} \right\|_2 \right\} \\
 &\quad + 2\beta_k^{(0)} \underline{\sigma}_{k,h} \min \left\{ 1, \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}} / \underline{\sigma}_{k,h} \right\|_2 \right\} \\
 &\quad + [\mathbb{P}_h \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) - [\mathbb{P}_h \underline{V}_{k,h+1}](s_h^k, a_h^k, b_h^k),
 \end{aligned}$$

where the second inequality holds because  $[\mathbb{P}_h \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) - [\mathbb{P}_h \underline{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) \geq -2H$ , the last inequality holds since  $\beta_k^{(0)} \bar{\sigma}_{k,h} \geq 2H$ ,  $\beta_k^{(0)} \underline{\sigma}_{k,h} \geq 2H$ . Subtracting  $\bar{V}_{k,h+1}(s_{h+1}^k) - \underline{V}_{k,h+1}(s_{h+1}^k)$  from the both side, we can further get,

$$\begin{aligned}
 &\bar{V}_{k,h}(s_h^k) - \underline{V}_{k,h}(s_h^k) - [\bar{V}_{k,h+1}(s_h^k) - \underline{V}_{k,h+1}(s_h^k)] \\
 &\leq 2\beta_k^{(0)} \bar{\sigma}_{k,h} \min \left\{ 1, \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} / \bar{\sigma}_{k,h} \right\|_2 \right\} \\
 &\quad + 2\beta_k^{(0)} \underline{\sigma}_{k,h} \min \left\{ 1, \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}} / \underline{\sigma}_{k,h} \right\|_2 \right\} \\
 &\quad + [\mathbb{P}_h \bar{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) - [\mathbb{P}_h \underline{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) \\
 &\quad - [\bar{V}_{k,h+1}(s_{h+1}^k) - \underline{V}_{k,h+1}(s_{h+1}^k)], \tag{B.12}
 \end{aligned}$$

Taking summation of (B.12) from  $k = 1 \dots K$  and  $h = h' \dots H$ , we have following inequality holds

$$\begin{aligned}
 &\sum_{k=1}^K [\bar{V}_{k,h'}(s_{h'}^k) - \underline{V}_{k,h'}(s_{h'}^k)] \\
 &\leq 2\beta_k^{(0)} \sum_{k=1}^K \sum_{h=h'}^H \bar{\sigma}_{k,h} \min \left\{ 1, \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} / \bar{\sigma}_{k,h} \right\|_2 \right\}
 \end{aligned}$$



$$\begin{aligned}
 & + 2\beta_k^{(0)} \sum_{k=1}^K \sum_{h=h'}^H \underline{\sigma}_{k,h} \min \left\{ 1, \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}} / \underline{\sigma}_{k,h} \right\|_2 \right\} \\
 & + \sum_{k=1}^K \sum_{h=h'}^H \left[ \mathbb{P}_h \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) - \mathbb{P}_h \underline{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) \right. \\
 & \quad \left. - [\bar{V}_{k,h+1}(s_{h+1}^k) - \underline{V}_{k,h+1}(s_{h+1}^k)] \right] \\
 & \leq 2\beta_k^{(0)} \sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h} \min \left\{ 1, \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} / \bar{\sigma}_{k,h} \right\|_2 \right\} \\
 & \quad + 2\beta_k^{(0)} \sum_{k=1}^K \sum_{h=1}^H \underline{\sigma}_{k,h} \min \left\{ 1, \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}} / \underline{\sigma}_{k,h} \right\|_2 \right\} + 8H\sqrt{2T \log(H/\delta)} \\
 & \leq 2\beta_k^{(0)} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \min \left\{ 1, \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}} / \bar{\sigma}_{k,h} \right\|_2^2 \right\}} \\
 & \quad + 2\beta_k^{(0)} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \underline{\sigma}_{k,h}^2} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \min \left\{ 1, \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}} / \underline{\sigma}_{k,h} \right\|_2^2 \right\}} \\
 & \quad + 8H\sqrt{2T \log(H/\delta)} \\
 & \leq 2\beta_K^{(0)} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2} \sqrt{2Hd \log(1 + K/\lambda)} \\
 & \quad + 2\beta_K^{(0)} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \underline{\sigma}_{k,h}^2} \sqrt{2Hd \log(1 + K/\lambda)} + 8H\sqrt{2T \log(H/\delta)} \\
 & \leq 4\beta_K^{(0)} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2 + \underline{\sigma}_{k,h}^2} \sqrt{2Hd \log(1 + K/\lambda)} + 8H\sqrt{2T \log(H/\delta)}, \tag{B.13}
 \end{aligned}$$

where the first inequality holds since  $\bar{V}_{k,H+1} = \underline{V}_{k,H+1} = 0$ , the second inequality holds on event  $\mathcal{E}_1$ , the third inequality holds due to Cauchy-Schwarz inequality, the fourth inequality holds due to Azuma Hoeffding inequality with the fact that  $\left\| \phi_{\bar{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k) / \bar{\sigma}_{k,h} \right\|_2 \leq \left\| \phi_{\bar{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k) \right\|_2 \cdot \sqrt{d}/H \leq \sqrt{d}$ ,  $\left\| \phi_{\underline{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k) / \underline{\sigma}_{k,h} \right\|_2 \leq \left\| \phi_{\underline{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k) \right\|_2 \cdot \sqrt{d}/H \leq \sqrt{d}$ , the last inequality is by the fact that  $\sqrt{a} + \sqrt{b} \leq 2\sqrt{a+b}$ . (B.13) holds for any  $h'$ , then we have following inequality holds

$$\begin{aligned}
 & \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h [\bar{V}_{k,h+1} - \underline{V}_{k,h+1}](s_h^k, a_h^k, b_h^k) \\
 & = \sum_{k=1}^K \sum_{h=1}^H [\bar{V}_{k,h} - \underline{V}_{k,h}](s_h^k) + \sum_{k=1}^K \sum_{h=1}^H \left[ \mathbb{P}_h \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) - \mathbb{P}_h \underline{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) \right. \\
 & \quad \left. - [\bar{V}_{k,h+1}(s_{h+1}^k) - \underline{V}_{k,h+1}(s_{h+1}^k)] \right] \\
 & \leq 4\beta_K^{(0)} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2 + \underline{\sigma}_{k,h}^2} \sqrt{2H^3d \log(1 + K/\lambda)} + 8H^2\sqrt{2T \log(H/\delta)},
 \end{aligned}$$

where the inequality holds due to (B.13) and on event  $\mathcal{E}_1$ .  $\square$

### B.5. Proof of Lemma A.5

To estimate the variance in weighted ridge regression we need the following lemma, which is similar to the Lemma A.3 but without tolerant error  $\epsilon$ .

**Lemma B.5.** Suppose the event  $\mathcal{E}$  hold. Then we have for any  $s, a, b, k, h$  following inequalities hold,

$$\underline{Q}_{k,h}(s, a, b) \leq Q_h^{\mu^k}(s, a, b) \leq \overline{Q}_{k,h}(s, a, b),$$

and

$$\underline{V}_{k,h}(s) \leq V_h^{\mu^k}(s) \leq \overline{V}_{k,h}(s).$$

*Proof.* For simplicity, we only prove the following UCB by induction,

$$Q_h^{\mu^k}(s, a, b) \leq \overline{Q}_{k,h}(s, a, b), V_h^{\mu^k}(s) \leq \overline{V}_{k,h}(s).$$

The base case  $h = H + 1$  holds trivially since the terminal cost is zero. Now we assume that the bounds (B.9) holds for step  $h + 1$ . That is,

$$Q_{h+1}^{\mu^k}(s, a, b) \leq \overline{Q}_{k,h+1}(s, a, b), V_{h+1}^{\mu^k}(s) \leq \overline{V}_{k,h+1}(s).$$

If  $\overline{Q}_{k,h}(s, a, b) \geq H$ , then it is obvious to have  $Q_h^{\mu^k}(s, a, b) \leq \overline{Q}_{k,h}(s, a, b)$ , otherwise we have that

$$\begin{aligned} & \overline{Q}_{k,h}(s, a, b) - Q_h^{\mu^k}(s, a, b) \\ &= \langle \bar{\theta}_{k,h}^{(0)}, \phi_{\overline{V}_{k,h+1}} \rangle + \beta_k^{(0)} \left\| [\overline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\overline{V}_{k,h+1}} \right\|_2 - \langle \theta_h^*, \phi_{\overline{V}_{k,h+1}} \rangle \\ & \quad + \mathbb{P}_h \overline{V}_{k,h+1}(s, a, b) - \mathbb{P}_h V_{h+1}^{*,\nu^k}(s) \\ & \geq \beta_k^{(0)} \left\| [\overline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\overline{V}_{k,h+1}} \right\|_2 - \left\| [\overline{\Sigma}_{k,h}^{(0)}]^{1/2} (\bar{\theta}_{k,h}^{(0)} - \theta_h^*) \right\|_2 \left\| [\overline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\overline{V}_{k,h+1}} \right\|_2 \\ & \quad + \mathbb{P}_h \overline{V}_{k,h+1}(s) - \mathbb{P}_h V_{h+1}^{\mu^k}(s) \\ & \geq \mathbb{P}_h \overline{V}_{k,h+1}(s) - \mathbb{P}_h V_{h+1}^{\mu^k}(s) \\ & \geq 0, \end{aligned} \tag{B.14}$$

where the first inequality holds due to Cauchy-Schwarz inequality, the second inequality holds since the assumption that  $\theta_h^* \in \overline{\mathcal{C}}_{k,h}^{(0)}$  in event  $\mathcal{E}$ , the third inequality holds by the induction assumption.

Then we have that

$$\begin{aligned} \overline{V}_{k,h}(s) &= \mathbb{E}_{(a,b) \sim \mu_h^k(\cdot, \cdot | s)} [\overline{Q}_{k,h}(s, a, b)] \\ &\geq \mathbb{E}_{(a,b) \sim \mu_h^k(\cdot, \cdot | s)} [Q_h^{\mu^k}(s, a', b)] \\ &= V_h^{\mu^k}(s), \end{aligned}$$

where the inequality is by (B.14), the last inequality is due to the Bellman equation. Therefore, our proof ends.  $\square$

**Lemma B.6.** (Lemma 11 in (Abbasi-Yadkori et al., 2011)). For any  $\{\mathbf{x}_t\}_{t=1}^T \subset \mathbb{R}^d$  satisfying that  $\|\mathbf{x}_t\|_2 \leq L$ , let  $\mathbf{A}_0 = \lambda \mathbf{I}$  and  $\mathbf{A}_t = \mathbf{A}_0 + \sum_{i=1}^t \mathbf{x}_i \mathbf{x}_i^\top$ , then we have

$$\sum_{t=1}^T \min\{1, \|\mathbf{x}_t\|_{\mathbf{A}_{t-1}^{-1}}^2\} \leq 2d \log \frac{d\lambda + TL^2}{d\lambda}.$$

*Proof of Lemma A.5.* Suppose the event in Lemma 9.1 holds, we have the following results:

$$\sum_{k=1}^K \sum_{h=1}^H \overline{\sigma}_{k,h}^2 = \sum_{k=1}^K \sum_{h=1}^H [H^2/d + \mathbb{V}_{k,h}^{\text{est}} \overline{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) + \overline{E}_{k,h}]$$

$$\begin{aligned}
 &= H^2 T / d + \underbrace{\sum_{k=1}^K \sum_{h=1}^H [\mathbb{V}_h \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) - \mathbb{V}_h V_{h+1}^{\mu^k}(s_h^k, a_h^k, b_h^k)]}_{I_1} \\
 &\quad + \underbrace{2 \sum_{k=1}^K \sum_{h=1}^H \bar{E}_{k,h}}_{I_2} + \underbrace{\sum_{k=1}^K \sum_{h=1}^H \mathbb{V}_h V_{h+1}^{\mu^k}(s_h^k, a_h^k, b_h^k)}_{I_3} \\
 &\quad + \underbrace{\sum_{k=1}^K \sum_{h=1}^H [\mathbb{V}_h^{\text{est}} \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) - \mathbb{V}_h \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) - \bar{E}_{k,h}]}_{I_4}, \tag{B.15}
 \end{aligned}$$

where the first equation is by the definition of  $\bar{\sigma}_{k,h}$ . To bound  $I_1$ , we have

$$\begin{aligned}
 I_1 &= \sum_{k=1}^K \sum_{h=1}^H [\mathbb{P}_h \bar{V}_{k,h+1}^2(s_h^k, a_h^k, b_h^k) - \mathbb{P}_h [V_{h+1}^{\mu^k}]^2(s_h^k, a_h^k, b_h^k)] \\
 &\quad - \sum_{k=1}^K \sum_{h=1}^H [(\mathbb{P}_h \bar{V}_{k,h+1})^2(s_h^k, a_h^k, b_h^k) - (\mathbb{P}_h V_{h+1}^{\mu^k})^2(s_h^k, a_h^k, b_h^k)] \\
 &\leq \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h [(\bar{V}_{k,h+1} - V_{h+1}^{\mu^k})(\bar{V}_{k,h+1} + V_{h+1}^{\mu^k})](s_h^k, a_h^k, b_h^k) \\
 &\quad - \sum_{k=1}^K \sum_{h=1}^H [(\mathbb{P}_h \bar{V}_{k,h+1} - \mathbb{P}_h V_{h+1}^{\mu^k})(\mathbb{P}_h \bar{V}_{k,h+1} + \mathbb{P}_h V_{h+1}^{\mu^k})](s_h^k, a_h^k, b_h^k) \\
 &\leq 4H \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h [\bar{V}_{k,h+1} - V_{h+1}^{\mu^k}](s_h^k, a_h^k, b_h^k) \\
 &= 4H \sum_{k=1}^K \sum_{h=1}^H \mathbb{P}_h [\bar{V}_{k,h+1} - V_{h+1}^{\mu^k}](s_h^k, a_h^k, b_h^k),
 \end{aligned}$$

where the first inequality is by  $|\bar{V}_{k,h+1}|, |V_{h+1}^{\mu^k}| \leq H$ , and the second inequality is by  $\bar{V}_{k,h+1} - V_{h+1}^{\mu^k} \geq 0$  due to Lemma B.5. To bound  $I_2$ , we have

$$\begin{aligned}
 I_2 &\leq 2 \sum_{k=1}^K \sum_{h=1}^H \beta_k^{(1)} \min \{1, \left\| [\bar{\Sigma}_{k,h}^{(1)}]^{-1/2} \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \right\|_2 \} \\
 &\quad + 4H \sum_{k=1}^K \sum_{h=1}^H \beta_k^{(2)} \bar{\sigma}_{k,h} \min \{1, \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k) / \bar{\sigma}_{k,h} \right\|_2 \} \\
 &\leq 2\beta_K^{(1)} \sqrt{T} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \min \{1, \left\| [\bar{\Sigma}_{k,h}^{(1)}]^{-1/2} \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \right\|_2^2 \}} \\
 &\quad + 7\beta_K^{(1)} H^2 \sqrt{T} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \min \{1, \left\| [\bar{\Sigma}_{k,h}^{(1)}]^{-1/2} \phi_{\bar{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k) \right\|_2^2 / \bar{\sigma} \}} \\
 &\leq 2\beta_K^{(2)} \sqrt{T} \sqrt{2dH \log(1 + KH^4/(d\lambda))} + 7\beta_K^{(1)} H^2 \sqrt{T} \sqrt{2dH \log(1 + K/\lambda)},
 \end{aligned}$$

where the first inequality holds due to  $\beta_k^{(1)} \geq H^2$  and  $\beta_k^{(2)} \bar{\sigma}_{k,h} \geq \sqrt{d} \cdot H / \sqrt{d} = H$ , the second inequality holds due to Cauchy-Schwartz inequality,  $\beta_k^{(1)} \leq \beta_K^{(1)}$ ,  $\beta_k^{(2)} \leq \beta_K^{(2)}$ ,

$$\bar{\sigma}_{k,h}^2 = \max \{ H^2 / d, \mathbb{V}_h^{\text{est}} \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) + \bar{E}_{k,h} \} \leq \max \{ H^2 / d, H^2 + 2H^2 \} = 3H^2,$$

the third inequality holds due to Lemma B.6. Next we bound  $I_3$ , since event  $\mathcal{E}_2$  holds, we have

$$I_3 \leq 3(HT + H^3 \log(1/\delta)).$$

Finally, due to on event  $\mathcal{E}$ , we have  $I_4 \leq 0$ . We finish the proof by substituting  $I_1, I_2, I_3, I_4$  into (B.15).  $\square$

## C. Full Version of Algorithm 1

In this section, we present the full version of Algorithm 1 in Algorithm 2.

---

### Algorithm 2 Nash-UCRL

---

```

1: Input: Regularization parameter  $\lambda$ , Number of episode  $K$ , number of horizon  $H$ .
2: For any  $h$ ,  $\bar{\Sigma}_{1,h}^{(i)} \leftarrow \underline{\Sigma}_{1,h}^{(i)} \leftarrow \lambda \mathbf{I}$ ;  $\bar{\mathbf{b}}_{1,h}^{(i)} \leftarrow \underline{\mathbf{b}}_{1,h}^{(i)} \leftarrow \mathbf{0}$ ;  $\bar{\boldsymbol{\theta}}_{1,h}^{(i)} \leftarrow \underline{\boldsymbol{\theta}}_{1,h}^{(i)} \leftarrow \mathbf{0}$ , for  $i \in \{0, 1\}$ .
3: for  $k = 1, \dots, K$  do
4:    $\bar{V}_{k,H+1}(\cdot) \leftarrow 0$ ,  $\underline{V}_{k,H+1}(\cdot) \leftarrow 0$ 
5:   for  $h = H, \dots, 1$  do
6:     Set  $\bar{Q}_{k,h}(\cdot, \cdot, \cdot)$  and  $\underline{Q}_{k,h}(\cdot, \cdot, \cdot)$  as in (C.1).
7:     for  $s \in \mathcal{S}$  do
8:       Let  $\mu_h^k(\cdot, \cdot | s) = \epsilon\text{-CCE}(\bar{Q}_{k,h}(s, \cdot, \cdot), \underline{Q}_{k,h}(s, \cdot, \cdot))$ .
9:        $\bar{V}_{k,h}(s) = \mathbb{E}_{(a,b) \sim \mu_h^k(\cdot, \cdot | s)} \bar{Q}_{k,h}(s, a, b)$ ,  $\underline{V}_{k,h}(s) = \mathbb{E}_{(a,b) \sim \mu_h^k(\cdot, \cdot | s)} \underline{Q}_{k,h}(s, a, b)$ 
10:       $\pi_h^k(\cdot | s) = \mathcal{P}_{\max} \mu_h^k(\cdot, \cdot | s)$ ,  $\nu_h^k(\cdot | s) = \mathcal{P}_{\min} \mu_h^k(\cdot, \cdot | s)$ 
11:    end for
12:  end for
13:  receives  $s_1^k$ 
14:  for  $h = 1, \dots, H$  do
15:    Take action  $a_h^k \sim \pi_h^k(s_h^k)$  and  $b_h^k \sim \nu_h^k(s_h^k)$  and receives  $s_{h+1}^k \sim \mathbb{P}(\cdot | s_h^k, a_h^k, b_h^k)$ .
16:    Set  $\bar{V}_{k,h+1}^{\text{est}}(s_h^k, a_h^k, b_h^k)$  and  $\underline{V}_{k,h+1}^{\text{est}}(s_h^k, a_h^k, b_h^k)$  as in (C.2).
17:    Set  $\bar{E}_{k,h}, \underline{E}_{k,h}, \bar{\boldsymbol{\sigma}}_{k,h}, \underline{\boldsymbol{\sigma}}_{k,h}, \bar{\Sigma}_{k+1,h}^{(0)}, \underline{\Sigma}_{k+1,h}^{(0)}, \bar{\mathbf{b}}_{k+1,h}^{(0)}, \underline{\mathbf{b}}_{k+1,h}^{(0)}, \bar{\Sigma}_{k+1,h}^{(1)}, \underline{\Sigma}_{k+1,h}^{(1)}, \bar{\mathbf{b}}_{k+1,h}^{(1)}, \underline{\mathbf{b}}_{k+1,h}^{(1)}$  as defined in (C.3).
18:    Set  $\bar{\boldsymbol{\theta}}_{k+1,h}^{(i)} \leftarrow [\bar{\Sigma}_{k+1,h}^{(i)}]^{-1} \bar{\mathbf{b}}_{k+1,h}^{(i)}$ ,  $\underline{\boldsymbol{\theta}}_{k+1,h}^{(i)} \leftarrow [\underline{\Sigma}_{k+1,h}^{(i)}]^{-1} \underline{\mathbf{b}}_{k+1,h}^{(i)}$ ,  $i = 0, 1$ 
19:  end for
20: end for

```

---

#### Update of optimistic action-value function:

$$\begin{aligned}
 \bar{Q}_{k,h}(\cdot, \cdot, \cdot) &\leftarrow \left[ r_h(\cdot, \cdot, \cdot) + \langle \bar{\boldsymbol{\theta}}_{k,h}^{(0)}, \phi_{\bar{V}_{k,h+1}}(\cdot, \cdot, \cdot) \rangle + \beta_k^{(0)} \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\bar{V}_{k,h+1}}(\cdot, \cdot, \cdot) \right\|_2 \right]_{[-H, H]} \\
 \underline{Q}_{k,h}(\cdot, \cdot, \cdot) &\leftarrow \left[ r_h(\cdot, \cdot, \cdot) + \langle \underline{\boldsymbol{\theta}}_{k,h}^{(0)}, \phi_{\underline{V}_{k,h+1}}(\cdot, \cdot, \cdot) \rangle - \beta_k^{(0)} \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \phi_{\underline{V}_{k,h+1}}(\cdot, \cdot, \cdot) \right\|_2 \right]_{[-H, H]}.
 \end{aligned} \tag{C.1}$$

#### Update of variance estimation:

$$\begin{aligned}
 \bar{V}_{k,h+1}^{\text{est}}(s_h^k, a_h^k, b_h^k) &\leftarrow [\langle \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \bar{\boldsymbol{\theta}}_{k,h}^{(1)} \rangle]_{[0, H^2]} - [\langle \phi_{\bar{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k), \bar{\boldsymbol{\theta}}_{k,h}^{(0)} \rangle]_{[-H, H]}^2, \\
 \underline{V}_{k,h+1}^{\text{est}}(s_h^k, a_h^k, b_h^k) &\leftarrow [\langle \phi_{\underline{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k), \underline{\boldsymbol{\theta}}_{k,h}^{(1)} \rangle]_{[0, H^2]} - [\langle \phi_{\underline{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k), \underline{\boldsymbol{\theta}}_{k,h}^{(0)} \rangle]_{[-H, H]}^2.
 \end{aligned} \tag{C.2}$$

#### Update of other parameters:

$$\begin{aligned}
 \bar{E}_{k,h} &= \min \{ H^2, \beta_k^{(1)} \left\| [\bar{\Sigma}_{k,h}^{(1)}]^{-1/2} \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \right\|_2 \} \\
 &\quad + \min \{ H^2, 2H \beta_k^{(2)} \left\| \bar{\Sigma}_{k,h}^{(0)-1/2} \phi_{\bar{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k) \right\|_2 \} \\
 \underline{E}_{k,h} &= \min \{ H^2, \beta_k^{(1)} \left\| [\underline{\Sigma}_{k,h}^{(1)}]^{-1/2} \phi_{\underline{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \right\|_2 \}
 \end{aligned}$$

$$\begin{aligned}
 & + \min \left\{ H^2, 2H\beta_k^{(2)} \left\| \underline{\Sigma}_{k,h}^{(0)-1/2} \phi_{V_{k,h+1}}(s_h^k, a_h^k, b_h^k) \right\|_2 \right\}, \\
 \bar{\sigma}_{k,h} &= \sqrt{\max\{H^2/4d, \mathbb{V}^{\text{est}} \bar{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) + \bar{E}_{k,h}\}}, \\
 \underline{\sigma}_{k,h} &= \sqrt{\max\{H^2/4d, \mathbb{V}^{\text{est}} \underline{V}_{k,h+1}(s_h^k, a_h^k, b_h^k) + \underline{E}_{k,h}\}}, \\
 \bar{\Sigma}_{k+1,h}^{(0)} &\leftarrow \bar{\Sigma}_{k,h}^{(0)} + \bar{\sigma}_{k,h}^{-2} \phi_{\bar{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k) \phi_{\bar{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k)^\top, \\
 \underline{\Sigma}_{k+1,h}^{(0)} &\leftarrow \underline{\Sigma}_{k,h}^{(0)} + \underline{\sigma}_{k,h}^{-2} \phi_{\underline{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k) \phi_{\underline{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k)^\top, \\
 \bar{\mathbf{b}}_{k+1,h}^{(0)} &= \bar{\mathbf{b}}_{k,h}^{(0)} + \bar{\sigma}_{k,h}^{-2} \phi_{\bar{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k) \bar{V}_{k,h+1}(s_{h+1}^k), \\
 \underline{\mathbf{b}}_{k+1,h}^{(0)} &= \underline{\mathbf{b}}_{k,h}^{(0)} + \underline{\sigma}_{k,h}^{-2} \phi_{\underline{V}_{k,h+1}}(s_h^k, a_h^k, b_h^k) \underline{V}_{k,h+1}(s_{h+1}^k), \\
 \bar{\Sigma}_{k+1,h}^{(1)} &\leftarrow \bar{\Sigma}_{k,h}^{(1)} + \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k)^\top, \\
 \underline{\Sigma}_{k+1,h}^{(1)} &\leftarrow \underline{\Sigma}_{k,h}^{(1)} + \phi_{\underline{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \phi_{\underline{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k)^\top, \\
 \bar{\mathbf{b}}_{k+1,h}^{(1)} &= \bar{\mathbf{b}}_{k,h}^{(1)} + \phi_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \bar{V}_{k,h+1}^2(s_{h+1}^k), \\
 \underline{\mathbf{b}}_{k+1,h}^{(1)} &= \underline{\mathbf{b}}_{k,h}^{(1)} + \phi_{\underline{V}_{k,h+1}^2}(s_h^k, a_h^k, b_h^k) \underline{V}_{k,h+1}^2(s_{h+1}^k). \tag{C.3}
 \end{aligned}$$

## D. Extensions to Turn-based Games

In this section, we extend our algorithm and results to turn-based Markov games.

**Turn-based MGs** A two-player zero-sum turn-based episodic MG is denoted by a tuple

$M(\mathcal{S}, \mathcal{A}, H, \{r_h\}_{h=1}^H, \{\mathbb{P}_h\}_{h=1}^H)$ , where  $\mathcal{S} = \mathcal{S}_{\max} \cup \mathcal{S}_{\min}$ ,  $\mathcal{S}_{\max}$  ( $\mathcal{S}_{\min}$ ) are the states where the max (min)-player plays,  $\mathcal{S}_{\max} \cap \mathcal{S}_{\min} = \emptyset$ . Note that the partition of state space suggests that at each step, only one player can play.  $\mathcal{A}$  is the action space,  $H$  is the length of game/episode,  $r_h : \mathcal{S} \times \mathcal{A} \rightarrow [-1, 1]$  is the reward function,  $\mathbb{P}_h(s'|s, a)$  denotes the transition probability for the max (min)-player ( $s \in \mathcal{S}_{\max}$  or  $\mathcal{S}_{\min}$ ) to take action  $a$  and transit to next state  $s'$ . Similar to the linear mixture MGs, we can define linear mixture turn-based MGs as follows.

**Definition D.1.**  $M(\mathcal{S}, \mathcal{A}, H, \{r_h\}_{h=1}^H, \{\mathbb{P}_h\}_{h=1}^H)$  is called a time inhomogeneous, episodic  $B$ -bounded linear mixture turn-based Markov game if there exist  $\{\theta_h\}_{h=1}^H \subset \mathbb{R}^d$  and  $\tilde{\phi}(s'|s, a) \in \mathbb{R}^d$  satisfying

$$\|\theta_h\|_2 \leq B, \quad \forall V : \mathcal{S} \rightarrow [-1, 1], \quad \left\| \sum_{s' \in \mathcal{S}} \tilde{\phi}(s'|s, a) V(s') \right\|_2 \leq 1,$$

such that  $\mathbb{P}_h(s'|s, a) = \langle \phi(s'|s, a), \theta_h \rangle$  for any state-action-state triplet  $(s, a, s')$  and any step  $h$ .

Based on above definition, we show that any turn-based linear mixture MG can be regarded as a special case of linear mixture simultaneous-move MG. In fact, for any turn-based linear mixture MG with feature mapping  $\tilde{\phi}(\cdot|\cdot, \cdot)$  and reward  $\tilde{r}_h(\cdot, \cdot)$ , we can define the corresponding linear mixture simultaneous-move MG with feature mapping  $\phi(\cdot|\cdot, \cdot)$  and reward  $r_h(\cdot, \cdot, \cdot)$  as follows: for each  $s \in \mathcal{S}_{\max}$ ,

$$\phi(s'|s, a, b) = \tilde{\phi}_h(s'|s, a), \quad r_h(s'|s, a, b) = \tilde{r}_h(s'|s, a),$$

and for each  $s \in \mathcal{S}_{\min}$ ,

$$\phi(s'|s, a, b) = \tilde{\phi}_h(s'|s, b), \quad r_h(s'|s, a, b) = \tilde{r}_h(s'|s, b).$$

Therefore, we can still use Algorithm 1 to find the Nash equilibrium. Notice that for the turn-based game, at each step only one player can take action. Thus, the  $\epsilon$ -CCE routine in Line 9 of Algorithm 1 needs be replaced by two separate subroutines: taking  $\pi_h^k$  and  $\nu_h^k$  as greedy policies w.r.t.  $\bar{Q}_{k,h}$  and  $\underline{Q}_{k,h}$ . For completeness, we present the turn-based version of Algorithm 1 as Algorithm 3.

By Theorem 4.1, we immediately have that the regret of our turn-based algorithm is also bounded by

$$\tilde{O}(\sqrt{d^2 H^2 + d H^3} \sqrt{T} + d^2 H^3 + d^3 H^2),$$

where  $T = KH$ . Similarly, we can show that if  $d \geq H$  and  $T \geq d^4 H^2$ , our turn-based algorithm is nearly minimax optimal.

Instead of dividing by state (Cui & Yang, 2020; Xie et al., 2020), some papers (Bai & Jin, 2020) consider dividing the steps  $[H]$  into  $\mathcal{H}_{\max}$  and  $\mathcal{H}_{\min}$ . This setting looks different at first glance but can actually be written as the turn-based MGs in this section by augmenting the state  $s$  by a latent variable  $h$ , i.e.,  $\tilde{s} = (s, h)$ , to indicate the max-player's step and the min-player's step.

---

**Algorithm 3** Turn-based Nash-UCRL
 

---

```

1: For any  $h$ ,  $\bar{\Sigma}_{1,h}^{(i)} \leftarrow \underline{\Sigma}_{1,h}^{(i)} \leftarrow \lambda \mathbf{I}$ ;  $\bar{\mathbf{b}}_{1,h}^{(i)} \leftarrow \underline{\mathbf{b}}_{1,h}^{(i)} \leftarrow \mathbf{0}$ ;  $\bar{\boldsymbol{\theta}}_{1,h}^{(i)} \leftarrow \underline{\boldsymbol{\theta}}_{1,h}^{(i)} \leftarrow \mathbf{0}$ , for  $i \in \{0, 1\}$ .
2: for  $k = 1, \dots, K$  do
3:    $\bar{V}_{k,H+1}(\cdot) \leftarrow 0$ ,  $\underline{V}_{k,H+1}(\cdot) \leftarrow 0$ 
4:   for  $h = H, \dots, 1$  do
5:     Set  $\bar{Q}_{k,h}(\cdot, \cdot)$  and  $\underline{Q}_{k,h}(\cdot, \cdot)$  as in (D.1).
6:     for  $s \in \mathcal{S}_{\max}$  do
7:        $\pi_h^k(\cdot|s) = \max_{a \in \mathcal{A}} \bar{Q}_{k,h}(s, a)$ ,  $\bar{V}_{k,h}(s) = \mathbb{E}_{a \sim \pi_h^k(\cdot|s)} \bar{Q}_{k,h}(s, a)$ .
8:     end for
9:     for  $s \in \mathcal{S}_{\min}$  do
10:       $\nu_h^k(\cdot|s) = \min_{b \in \mathcal{A}} \underline{Q}_{k,h}(s, b)$ ,  $\underline{V}_{k,h}(s) = \mathbb{E}_{b \sim \nu_h^k(\cdot|s)} \underline{Q}_{k,h}(s, b)$ .
11:    end for
12:  end for
13:  receives  $s_{k,1}$ 
14:  for  $h = 1, \dots, H$  do
15:    if  $s_h^k \in \mathcal{S}_{\max}$  then
16:      Take action  $a_h^k \sim \pi_h^k(\cdot|s_h^k)$  and receives  $s_{h+1}^k \sim \mathbb{P}(\cdot|s_h^k, a_h^k)$ .
17:    else
18:      Take action  $a_h^k \sim \nu_h^k(\cdot|s_h^k)$  and receives  $s_{h+1}^k \sim \mathbb{P}(\cdot|s_h^k, a_h^k)$ .
19:    end if
20:    Set  $\nabla^{\text{est}} \bar{V}_{k,h+1}(s_h^k, a_h^k)$  and  $\nabla^{\text{est}} \underline{V}_{k,h+1}(s_h^k, a_h^k)$  as in (D.2).
21:    Set  $\bar{E}_{k,h}, \underline{E}_{k,h}, \bar{\sigma}_{k,h}, \underline{\sigma}_{k,h}, \bar{\Sigma}_{k+1,h}^{(0)}, \underline{\Sigma}_{k+1,h}^{(0)}, \bar{\mathbf{b}}_{k+1,h}^{(0)}, \underline{\mathbf{b}}_{k+1,h}^{(0)}, \bar{\Sigma}_{k+1,h}^{(1)}, \underline{\Sigma}_{k+1,h}^{(1)}, \bar{\mathbf{b}}_{k+1,h}^{(1)}, \underline{\mathbf{b}}_{k+1,h}^{(1)}$  as defined in (D.3).
22:    Set  $\bar{\boldsymbol{\theta}}_{k+1,h}^{(i)} \leftarrow [\bar{\Sigma}_{k+1,h}^{(i)}]^{-1} \bar{\mathbf{b}}_{k+1,h}^{(i)}$ ,  $\underline{\boldsymbol{\theta}}_{k+1,h}^{(i)} \leftarrow [\underline{\Sigma}_{k+1,h}^{(i)}]^{-1} \underline{\mathbf{b}}_{k+1,h}^{(i)}$ ,  $i = 0, 1$ 
23:  end for
24: end for
    
```

---

**Update of optimistic action-value function:**

$$\begin{aligned}
 \bar{Q}_{k,h}(\cdot, \cdot) &\leftarrow \min\{H, \tilde{r}_h(\cdot, \cdot) + \langle \bar{\boldsymbol{\theta}}_{k,h}^{(0)}, \tilde{\boldsymbol{\phi}}_{\bar{V}_{k,h+1}}(\cdot, \cdot) \rangle + \beta_k^{(0)} \left\| [\bar{\Sigma}_{k,h}^{(0)}]^{-1/2} \tilde{\boldsymbol{\phi}}_{\bar{V}_{k,h+1}}(\cdot, \cdot) \right\|_2\} \\
 \underline{Q}_{k,h}(\cdot, \cdot) &\leftarrow \max\{-H, \tilde{r}_h(\cdot, \cdot) + \langle \underline{\boldsymbol{\theta}}_{k,h}^{(0)}, \tilde{\boldsymbol{\phi}}_{\underline{V}_{k,h+1}}(\cdot, \cdot) \rangle - \beta_k^{(0)} \left\| [\underline{\Sigma}_{k,h}^{(0)}]^{-1/2} \tilde{\boldsymbol{\phi}}_{\underline{V}_{k,h+1}}(\cdot, \cdot) \right\|_2\}.
 \end{aligned} \tag{D.1}$$

**Update of variance estimation:**

$$\begin{aligned}
 \nabla^{\text{est}} \bar{V}_{k,h+1}(s_h^k, a_h^k) &\leftarrow [\langle \tilde{\boldsymbol{\phi}}_{\bar{V}_{k,h+1}}^2(s_h^k, a_h^k), \bar{\boldsymbol{\theta}}_{k,h}^{(1)} \rangle]_{[0, H^2]} - [\langle \tilde{\boldsymbol{\phi}}_{\bar{V}_{k,h+1}}(s_h^k, a_h^k), \bar{\boldsymbol{\theta}}_{k,h}^{(0)} \rangle]_{[-H, H]}^2, \\
 \nabla^{\text{est}} \underline{V}_{k,h+1}(s_h^k, a_h^k) &\leftarrow [\langle \tilde{\boldsymbol{\phi}}_{\underline{V}_{k,h+1}}^2(s_h^k, a_h^k), \underline{\boldsymbol{\theta}}_{k,h}^{(1)} \rangle]_{[0, H^2]} - [\langle \tilde{\boldsymbol{\phi}}_{\underline{V}_{k,h+1}}(s_h^k, a_h^k), \underline{\boldsymbol{\theta}}_{k,h}^{(0)} \rangle]_{[-H, H]}^2.
 \end{aligned} \tag{D.2}$$

**Update of other parameters:**

$$\begin{aligned}
 \bar{E}_{k,h} &= \min\{H^2, \beta_k^{(1)} \left\| [\bar{\Sigma}_{k,h}^{(1)}]^{-1/2} \tilde{\boldsymbol{\phi}}_{\bar{V}_{k,h+1}}^2(s_h^k, a_h^k) \right\|_2\} \\
 &\quad + \min\{H^2, 2H\beta_k^{(2)} \left\| \bar{\Sigma}_{k,h}^{(0)-1/2} \tilde{\boldsymbol{\phi}}_{\bar{V}_{k,h+1}}(s_h^k, a_h^k) \right\|_2\}, \\
 \underline{E}_{k,h} &= \min\{H^2, \beta_k^{(1)} \left\| [\underline{\Sigma}_{k,h}^{(1)}]^{-1/2} \tilde{\boldsymbol{\phi}}_{\underline{V}_{k,h+1}}^2(s_h^k, a_h^k) \right\|_2\}
 \end{aligned}$$



$$\begin{aligned}
 & + \min \{ H^2, 2H\beta_k^{(2)} \left\| \underline{\Sigma}_{k,h}^{(0)-1/2} \tilde{\Phi}_{\underline{V}_{k,h+1}}(s_h^k, a_h^k) \right\|_2 \}, \\
 \bar{\sigma}_{k,h} &= \sqrt{\max\{H^2/d, \mathbb{V}^{\text{est}} \bar{V}_{k,h+1}(s_h^k, a_h^k) + \bar{E}_{k,h}\}}, \\
 \underline{\sigma}_{k,h} &= \sqrt{\max\{H^2/d, \mathbb{V}^{\text{est}} \underline{V}_{k,h+1}(s_h^k, a_h^k) + \underline{E}_{k,h}\}}, \\
 \bar{\Sigma}_{k+1,h}^{(0)} &\leftarrow \bar{\Sigma}_{k,h}^{(0)} + \bar{\sigma}_{k,h}^{-2} \tilde{\Phi}_{\bar{V}_{k,h+1}}(s_h^k, a_h^k) \tilde{\Phi}_{\bar{V}_{k,h+1}}(s_h^k, a_h^k)^\top \\
 \underline{\Sigma}_{k+1,h}^{(0)} &\leftarrow \underline{\Sigma}_{k,h}^{(0)} + \underline{\sigma}_{k,h}^{-2} \tilde{\Phi}_{\underline{V}_{k,h+1}}(s_h^k, a_h^k) \tilde{\Phi}_{\underline{V}_{k,h+1}}(s_h^k, a_h^k)^\top \\
 \bar{\mathbf{b}}_{k+1,h}^{(0)} &= \bar{\mathbf{b}}_{k,h}^{(0)} + \bar{\sigma}_{k,h}^{-2} \tilde{\Phi}_{\bar{V}_{k,h+1}}(s_h^k, a_h^k) \bar{V}_{k,h+1}(s_{k,h+1}^k) \\
 \underline{\mathbf{b}}_{k+1,h}^{(0)} &= \underline{\mathbf{b}}_{k,h}^{(0)} + \underline{\sigma}_{k,h}^{-2} \tilde{\Phi}_{\underline{V}_{k,h+1}}(s_h^k, a_h^k) \underline{V}_{k,h+1}(s_{k,h+1}^k) \\
 \bar{\Sigma}_{k+1,h}^{(1)} &\leftarrow \bar{\Sigma}_{k,h}^{(1)} + \tilde{\Phi}_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k) \tilde{\Phi}_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k)^\top \\
 \underline{\Sigma}_{k+1,h}^{(1)} &\leftarrow \underline{\Sigma}_{k,h}^{(1)} + \tilde{\Phi}_{\underline{V}_{k,h+1}^2}(s_h^k, a_h^k) \tilde{\Phi}_{\underline{V}_{k,h+1}^2}(s_h^k, a_h^k)^\top \\
 \bar{\mathbf{b}}_{k+1,h}^{(1)} &= \bar{\mathbf{b}}_{k,h}^{(1)} + \tilde{\Phi}_{\bar{V}_{k,h+1}^2}(s_h^k, a_h^k) \bar{V}_{k,h+1}^2(s_{h+1}^k), \\
 \underline{\mathbf{b}}_{k+1,h}^{(1)} &= \underline{\mathbf{b}}_{k,h}^{(1)} + \tilde{\Phi}_{\underline{V}_{k,h+1}^2}(s_h^k, a_h^k) \underline{V}_{k,h+1}^2(s_{h+1}^k).
 \end{aligned} \tag{D.3}$$