

REPRESENTATION LEARNING FOR GENERAL-SUM LOW-RANK MARKOV GAMES

Chengzhuo Ni

Princeton University
cn10@princeton.edu

Yuda Song

Carnegie Mellon University
yudas@andrew.cmu.edu

Xuezhou Zhang

Princeton University
xz7392@princeton.edu

Zihan Ding

Princeton University
zihand@princeton.edu

Chi Jin

Princeton University
chij@princeton.edu

Mengdi Wang

Princeton University
mengdiw@princeton.edu

ABSTRACT

We study multi-agent general-sum Markov games with nonlinear function approximation. We focus on low-rank Markov games whose transition matrix admits a hidden low-rank structure on top of an unknown non-linear representation. The goal is to design an algorithm that (1) finds an ε -equilibrium policy sample efficiently without prior knowledge of the environment or the representation, and (2) permits a deep-learning friendly implementation. We leverage representation learning and present a model-based and a model-free approach to construct an effective representation from collected data. For both approaches, the algorithm achieves a sample complexity of $\text{poly}(H, d, A, 1/\varepsilon)$, where H is the game horizon, d is the dimension of the feature vector, A is the size of the joint action space and ε is the optimality gap. When the number of players is large, the above sample complexity can scale exponentially with the number of players in the worst case. To address this challenge, we consider Markov Games with a factorized transition structure and present an algorithm that escapes such exponential scaling. To our best knowledge, this is the first sample-efficient algorithm for multi-agent general-sum Markov games that incorporates (non-linear) function approximation. We accompany our theoretical result with a neural network-based implementation of our algorithm and evaluate it against the widely used deep RL baseline, DQN with fictitious play.

1 INTRODUCTION

Multi-agent reinforcement learning (MARL) studies the problem where multiple agents learn to make sequential decisions in an unknown environment to maximize their (own) cumulative rewards. Recently, MARL has achieved remarkable empirical success, such as in traditional games like GO (Silver et al., 2016, 2017) and Poker (Moravčík et al., 2017), real-time video games such as Starcraft and Dota 2 (Vinyals et al., 2019; Berner et al., 2019), decentralized controls or multi-agent robotics systems (Brambilla et al., 2013) and autonomous driving (Shalev-Shwartz et al., 2016).

On the theoretical front, however, provably sample-efficient algorithms for Markov games have been largely restricted to either two-player zero-sum games (Bai et al., 2020; Xie et al., 2020; Chen et al., 2021; Jin et al., 2021c) or general-sum games with small and finite state and action spaces (Bai and Jin, 2020; Liu et al., 2021; Jin et al., 2021b). These algorithms typically do not permit a scalable implementation applicable to real-world games, due to either (1) they only work for tabular or linear Markov games which are too restrictive to model real-world games, or (2) the ones that do handle rich non-linear function approximation (Jin et al., 2021c) are not computationally efficient. This motivates us to ask the following question:

Can we design an efficient algorithm that (1) provably learns multi-player general-sum Markov games with rich nonlinear function approximation and (2) permits scalable implementations?

This paper presents the first positive answer to the above question. In particular, we make the following contributions:

1. We design a new centralized self-play meta algorithm for multi-agent low-rank Markov games: **General Representation Learning for Multi-player General-sum Markov Game (GERL_MG2)**. We present a model-based and a model-free instantiation of GERL_MG2 which differ by the way function approximation is used, and a clean and unified analysis for both approaches.
2. We show that the model-based variant requires access to an MLE oracle and a NE/CE/CCE oracle for matrix games, and enjoys a $\tilde{O}(H^6 d^4 A^2 \log(|\Phi||\Psi|)/\varepsilon^2)$ sample complexity to learn an ε -NE/CE/CCE equilibrium policy, where d is the dimension of the feature vector, A is the size of the joint action space, H is the game horizon, Φ and Ψ are the function classes for the representation and emission process. The model-free variant replaces model-learning with solving a minimax optimization problem, and enjoys a sample complexity of $\tilde{O}(H^6 d^4 A^3 M \log(|\Phi|)/\varepsilon^2)$ for a slightly restricted class of Markov game with latent block structure.
3. Both of the above algorithms have sample complexities scaling with the joint action space size, which is exponential in the number of players. This unfavorable scaling is referred to as the *curse of multi-agent*, and is unavoidable in the worst case under general function approximation. We consider a spatial factorization structure where the transition of each player’s local state is directly affected only by at most $L = O(1)$ players in its adjacency. Given this additional structure, we provide an algorithm that achieves $\tilde{O}(M^4 H^6 d^{2(L+1)} \tilde{A}^{2(L+1)}/\varepsilon^2)$ sample complexity, where \tilde{A} is the size of a single player’s action space, thus escaping the exponential scaling to the number of agents.
4. Finally, we provide an efficient implementation of our reward-free algorithm, and show that it achieves superior performance against traditional deep RL baselines without principled representation learning.

1.1 RELATED WORKS

Markov games Markov games (Littman, 1994; Shapley, 1953) is an extensively used framework introduced for game playing with sequential decision making. Previous works (Littman, 1994; Hu and Wellman, 2003; Hansen et al., 2013) studied how to find the Nash equilibrium of a Markov game when the transition matrix and reward function are known. When the dynamic of the Markov game is unknown, recent works provide a line of finite-sample guarantees for learning Nash equilibrium in two-player zero-sum Markov games (Bai and Jin, 2020; Xie et al., 2020; Bai et al., 2020; Zhang et al., 2020; Liu et al., 2021; Jin et al., 2021c; Huang et al., 2021) and learning various equilibriums (including NE, CE, CCE, which are standard solution notions in games (Roughgarden, 2010)) in general-sum Markov games (Liu et al., 2021; Bai et al., 2021; Jin et al., 2021b). Some of the analyses in these works are based on the techniques for learning single-agent Markov Decision Processes (MDPs) (Azar et al., 2017; Jin et al., 2018, 2020).

RL with Function Approximation Function approximation in reinforcement learning has been extensively studied in recent years. For the single-agent Markov decision process, function approximation is adopted to achieve a better sample complexity that depends on the complexity of function approximators rather than the size of the state-action space. For example, (Yang and Wang, 2019; Jin et al., 2020; Zanette et al., 2020) considered the linear MDP model, where the transition probability function and reward function are linear in some feature mapping over state-action pairs. Another line of works (see, e.g., Jiang et al., 2017; Jin et al., 2021a; Du et al., 2021; Foster et al., 2021) studied the MDPs with general nonlinear function approximations.

When it comes to Markov game, (Chen et al., 2021; Xie et al., 2020; Jia et al., 2019) studied the Markov games with linear function approximations. Recently, (Huang et al., 2021) and (Jin et al., 2021c) proposed the first algorithms for two-player zero-sum Markov games with general function approximation, and provided a sample complexity governed by the minimax Eluder dimension. However, technical difficulties prevent extending these results to multi-player general-sum Markov games with nonlinear function approximation. The results for linear function approximation assume a known state-action feature, and are unable to solve the Markov games with a more general non-linear approximation where both the feature and function parameters are unknown. For the general function class works, their approaches rely heavily on the two-player nature, and it’s not clear how to apply their methods to the general multi-player setting.

Representation Learning in RL Our work is closely related to representation learning in single-agent RL, where the study mainly focuses on the low-rank MDPs. A low-rank MDP is strictly more

general than a linear MDP which assumes the representation is known a priori. Several related works studied low-rank MDPs with provable sample complexities. (Agarwal et al., 2020b; Ren et al., 2021) and (Uehara et al., 2021) consider the model-based setting, where the algorithm learns the representation with the model class of the transition probability given. (Modi et al., 2021) provided a representation learning algorithm under the model-free setting and proved its sample efficiency when the MDP satisfies the minimal reachability assumption. (Zhang et al., 2022) proposed a model-free method for the more restricted MDP class called Block MDP, but does not rely on the reachability assumption, which is also studied in papers including (Du et al., 2019) and (Misra et al., 2020). A concurrent work (Qiu et al., 2022) studies representation learning in RL with contrastive learning and extends their algorithm to the Markov game setting. However, their method requires strong data assumption and does not provide any practical implementation in the Markov game setting.

2 PROBLEM SETTINGS

A general-sum Markov game with M players is defined by a tuple $(\mathcal{S}, \{\mathcal{A}_i\}_{i=1}^M, P^*, \{r_i\}_{i=1}^M, H, d_1)$. Here \mathcal{S} is the state space, \mathcal{A}_i is the action space for player i , H is the time horizon of each episode and d_1 is the initial state distribution. We let $\mathcal{A} = \mathcal{A}_1 \times \dots \times \mathcal{A}_M$ and use $\mathbf{a} = (a_1, a_2, \dots, a_M)$ to denote the joint actions by all M players. Denote $\bar{A} = \max_i |\mathcal{A}_i|$ and $A = |\mathcal{A}|$. $P^* = \{P_h^*\}_{h=1}^H$ is a collection of transition probabilities, so that $P_h^*(\cdot|s, \mathbf{a})$ gives the distribution of the next state if actions \mathbf{a} are taken at state s and step h . And $r_i = \{r_{h,i}\}_{h=1}^H$ is a collection of reward functions, so that $r_{h,i}(s, \mathbf{a})$ gives the reward received by player i when actions \mathbf{a} are taken at state s and step h .

2.1 SOLUTION CONCEPTS

The policy of player i is denoted as $\pi_i := \{\pi_{h,i} : \mathcal{S} \rightarrow \Delta_{\mathcal{A}_i}\}_{h \in [H]}$. We denote the product policy of all the players as $\pi := \pi_1 \times \dots \times \pi_M$, here ‘‘product’’ means that conditioned on the same state, the action of each player is sampled independently according to their own policy. We denote the policy of all the players except the i th player as π_{-i} . We define $V_{h,i}^\pi(s)$ as the expected cumulative reward that will be received by the i th player if starting at state s at step h and all players follow policy π . For any strategy π_{-i} , there exists a best response of the i th player, which is a policy $\mu^\dagger(\pi_{-i})$ satisfying $V_{h,i}^{\mu^\dagger(\pi_{-i}), \pi_{-i}}(s) = \max_{\pi_i} V_{h,i}^{\pi_i, \pi_{-i}}(s)$ for any $(s, h) \in \mathcal{S} \times [H]$. We denote $V_{h,i}^{\dagger, \pi_{-i}} := V_{h,i}^{\mu^\dagger(\pi_{-i}), \pi_{-i}}$. Let $v_i^{\dagger, \pi_{-i}} := \mathbb{E}_{s \sim d_1} [V_{1,i}^{\dagger, \pi_{-i}}(s)]$, $v_i^\pi := \mathbb{E}_{s \sim d_1} [V_{1,i}^\pi(s)]$.

Definition 2.1 (NE). A product policy π is a Nash equilibrium (NE) if $v_i^\pi = v_i^{\dagger, \pi_{-i}}, \forall i \in [M]$. And we call π an ε -approximate NE if $\max_{i \in [M]} \{v_i^{\dagger, \pi_{-i}} - v_i^\pi\} < \varepsilon$.

The coarse correlated equilibrium (CCE) is a relaxed version of Nash equilibrium in which we consider general correlated policies instead of product policies.

Definition 2.2 (CCE). A correlated policy π is a CCE if $V_{h,i}^{\dagger, \pi_{-i}}(s) \leq V_{h,i}^\pi(s)$ for all $s \in \mathcal{S}, h \in [H], i \in [M]$. And we call π an ε -approximate CCE if $\max_{i \in [M]} \{v_i^{\dagger, \pi_{-i}} - v_i^\pi\} < \varepsilon$.

The correlated equilibrium (CE) is another relaxation of the Nash equilibrium. To define CE, we first introduce the concept of strategy modification: A strategy modification $\omega_i := \{\omega_{h,i}\}_{h \in [H]}$ for player i is a set of H functions from $\mathcal{S} \times \mathcal{A}_i$ to \mathcal{A}_i . Let $\Omega_i := \{\Omega_{h,i}\}_{h \in [H]}$ denote the set of all possible strategy modifications for player i . One can compose a strategy modification ω_i with any Markov policy π and obtain a new policy $\omega_i \circ \pi$ such that when policy π chooses to play $\mathbf{a} := (a_1, \dots, a_M)$ at state s and step h , policy $\omega_i \circ \pi$ will play $(a_1, \dots, a_{i-1}, \omega_{h,i}(s, a_i), a_{i+1}, \dots, a_M)$ instead.

Definition 2.3 (CE). A correlated policy π is a CE if $\max_{i \in [M]} \max_{\omega_i \in \Omega_i} V_{h,i}^{\omega_i \circ \pi}(s) \leq V_{h,i}^\pi(s)$ for all $(s, h) \in \mathcal{S} \times [H]$. And we call π an ε -approximate CE if $\max_{i \in [M]} \{\max_{\omega_i \in \Omega_i} v_i^{\omega_i \circ \pi} - v_i^\pi\} < \varepsilon$.

Remark 2.1. For general-sum Markov Games, we have $\{\text{NE}\} \subseteq \{\text{CE}\} \subseteq \{\text{CCE}\}$, so that they form a nested set of notions of equilibria (Roughgarden, 2010). While there exist algorithms to approximately compute the Nash equilibrium (Berg and Sandholm, 2017), the computation of NE for general-sum games in the worst case is still PPD-hard (Daskalakis, 2013). On the other hand, CCE and CE can be solved in polynomial time using linear programming (Examples include Papadimitriou and Roughgarden (2008); Blum et al. (2008)). Therefore, in this paper we study both NE and these weaker equilibrium concepts that permit more computationally efficient solutions.

2.2 LOW-RANK MARKOV GAMES

In this paper, we consider the class of low-rank Markov games. A Markov game is called a low-rank Markov game if the transition probability at any time step h has a latent low-rank structure.

Definition 2.4 (Low-Rank Markov Game). *We call a Markov game a low-rank Markov game if for any $s, s' \in \mathcal{S}, \mathbf{a} \in \mathcal{A}, h \in [H], i \in [M]$, we have $P_h^*(s'|s, \mathbf{a}) = \phi_h^*(s, \mathbf{a})^\top w_h^*(s')$, where $\|\phi_h^*(s, \mathbf{a})\|_2 \leq 1$ and $\|w_h^*(s')\|_2 \leq \sqrt{d}$ for all (s, \mathbf{a}, s') .*

A special case of low-rank Markov Game is the Block Markov game:

Definition 2.5 (Block Markov Game). *Consider any $h \in [H]$. A Block Markov game has an emission distribution $o_h(\cdot|z) \in \Delta_{\mathcal{S}}$ and a latent state space transition $T_h(z'|z, \mathbf{a})$, such that for any $s \in \mathcal{S}, o_h(s|z) > 0$ for a unique latent state $z \in \mathcal{Z}$, denoted as $\psi_h^*(s)$. Denote $Z = |\mathcal{Z}|$. Together with the ground truth decoder ψ_h^* , it defines the transitions $P_h^*(s'|s, \mathbf{a}) = \sum_{z' \in \mathcal{Z}} o_h(s'|z')T_h(z'|\psi_h^*(s), \mathbf{a})$.*

With the definition of the Block Markov game, one can naturally derive a feature vector that in addition takes the one-hot form: we just need to let the ground truth $\phi_h^*(s, \mathbf{a})$ at step h be a $Z \cdot A$ -dimensional vector $e_{(\psi^*(s), \mathbf{a})}$ where e_i is the i -th basis vector. Correspondingly, for any $s \in \mathcal{S}, w_h^*(s)$ is a $Z \cdot A$ dimensional vector such that the (z, \mathbf{a}) -th entry is $\sum_{z' \in \mathcal{Z}} o_h(s|z')T_h(z'|z, \mathbf{a})$. Then $P_h^*(s'|s, \mathbf{a}) = \phi_h^*(s, \mathbf{a})^\top w_h^*(s')$, so that the Block Markov game is a low-rank Markov game with rank $d = Z \cdot A$.

Learning Objective The goal of multi-agent reinforcement learning is to design algorithms for Markov games that find an ε -approximate equilibrium (NE, CCE, CE) from a small number of interactions with the environment. We focus on the low-rank Markov games whose feature vector ϕ^* and transition probability P^* are both *unknown*, and the goal is to identify a ε -approximate equilibrium policy with a number of interactions scaling polynomially with $d, A, H, \frac{1}{\varepsilon}$ and the log-cardinality of the function class, without depending on the number of raw states which could be infinite.

3 ALGORITHM DESCRIPTION

In this section, we present our algorithm GERL_MG2 (see Alg. 1). The algorithm mainly consists of two modules: the representation learning module and the planning module. We develop the representation learning module base on the past works on the single agent MDP (e.g. Agarwal et al. (2020b); Uehara et al. (2021); Modi et al. (2021)), and but modify them to work with UCB-style planning module. Here we denote $d_{P, h}^\pi$ as the state distribution under transition probability P and policy π at step h , and $U(\mathcal{A})$ as the uniform distribution over the joint action space.

3.1 REPRESENTATION LEARNING

In the representation learning module, the main goal is to learn a representation function $\hat{\phi}$ to approximate ϕ^* , using the data collected so far. In each episode, the algorithm first collects some new data using the policy derived from the previous episode. Note that in our data collection scheme, for each time step h , we maintain two buffers $\mathcal{D}_h^{(n)}$ and $\tilde{\mathcal{D}}_h^{(n)}$ of transition tuples (s, a, s') (line 7 of Alg. 1) which draw the state s from slightly different distributions. Based on the data collected in history, the representation learning module estimates the feature $\hat{\phi}^{(n)}$ and transition probability $\hat{P}^{(n)}$. We propose two versions of the representation learning algorithm (model-based, Alg. 2; model-free, Alg. 3) based on whether we are given the full model class \mathcal{M}_h of the transition probability, or only the function class of the state-action features Φ_h .

Model-based Representation Learning In the model-based setting, we assume the access to a *realizable* model class $\mathcal{M}_h = \{(w_h, \phi_h) : w_h \in \Psi_h, \phi_h \in \Phi_h\}, h \in [H]$ such that the true model is included in this class, i.e., $w_h^* \in \Psi_h, \phi_h^* \in \Phi_h, \forall h \in [H]$. Following the norm bounds on ϕ_h^*, w_h^* , we assume that the same norm bounds hold for our function approximator, i.e., for any $\phi_h \in \Phi_h, w_h \in \Psi_h$, we have $\|\phi_h(s, \mathbf{a})\|_2 \leq 1$ and $\|w_h(s')\|_2 \leq \sqrt{d}$ for all (s, \mathbf{a}, s') , and $\int \phi_h(s, \mathbf{a})^\top w_h(s') ds' = 1$. Given the dataset $\mathcal{D} := \mathcal{D}_h^{(n)} \cup \tilde{\mathcal{D}}_h^{(n)}$, MBREPLEARN learns the features and transition probability using maximum likelihood estimation (MLE):

$$\left(\hat{w}_h^{(n)}, \hat{\phi}_h^{(n)} \right) = \arg \max_{(w, \phi) \in \mathcal{M}_h} \mathbb{E}_{\mathcal{D}} [\log (\phi(s, \mathbf{a})^\top w(s'))], \quad \hat{P}_h^{(n)}(s'|s, \mathbf{a}) = \hat{\phi}_h^{(n)}(s, \mathbf{a})^\top \hat{w}_h^{(n)}(s').$$

Algorithm 1 General Representation Learning for Multi-player General-sum Low-Rank Markov Game with UCB-driven Exploration (GERL_MG2)

-
- 1: **Input:** Regularizer λ , iteration N , parameter $\{\alpha^{(n)}\}_{n=1}^N, \{\zeta^{(n)}\}_{n=1}^N$.
 - 2: Initialize $\pi^{(0)}$ to be uniform; set $\mathcal{D}_h^{(0)} = \emptyset, \tilde{\mathcal{D}}_h^{(0)} = \emptyset, \forall h \in [H]$.
 - 3: **for** episode $n = 1, 2, \dots, N$ **do**
 - 4: Set $\bar{V}_{H+1,i}^{(n)} \leftarrow 0, \underline{V}_{H+1,i}^{(n)} \leftarrow 0$
 - 5: **for** step $h = H, H-1 \dots, 1$ **do**
 - 6: Collect two triples $(s, \mathbf{a}, s'), (\tilde{s}', \tilde{\mathbf{a}}', \tilde{s}'')$ with
 $s \sim d_{P_{h^*}^*}^{\pi^{(n-1)}}, \mathbf{a} \sim U(\mathcal{A}), s' \sim P_h^*(s, \mathbf{a}),$
 $\tilde{s} \sim d_{P_{h^*}^*}^{\pi^{(n-1)}}, \tilde{\mathbf{a}} \sim U(\mathcal{A}), \tilde{s}' \sim P_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}), \tilde{\mathbf{a}}' \sim U(\mathcal{A}), \tilde{s}'' \sim P_h^*(\tilde{s}', \tilde{\mathbf{a}}').$
 - 7: Update datasets: $\mathcal{D}_h^{(n)} = \mathcal{D}_h^{(n-1)} \cup \{(s, \mathbf{a}, s')\}, \tilde{\mathcal{D}}_h^{(n)} = \tilde{\mathcal{D}}_h^{(n-1)} \cup \{(\tilde{s}', \tilde{\mathbf{a}}', \tilde{s}'')\}.$
 - 8: Learn representation via model-based or model-free methods:
 $\phi_h^{(n)}, \hat{P}_h^{(n)} = \text{MBREPLEARN}(\mathcal{D}_h^{(n)} \cup \tilde{\mathcal{D}}_h^{(n)}, h)$ or $\text{MFREPLEARN}(\mathcal{D}_h^{(n)} \cup \tilde{\mathcal{D}}_h^{(n)}, h, \lambda)$
 - 9: Compute $\hat{\beta}_h^{(n)}$ from equation 5, for each $(s, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}, i \in [M]$, set

$$\bar{Q}_{h,i}^{(n)}(s, \mathbf{a}) \leftarrow r_{h,i}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \bar{V}_{h+1,i}^{(n)} \right)(s, \mathbf{a}) + \hat{\beta}_h^{(n)}(s, \mathbf{a})$$

$$\underline{Q}_{h,i}^{(n)}(s, \mathbf{a}) \leftarrow r_{h,i}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \underline{V}_{h+1,i}^{(n)} \right)(s, \mathbf{a}) - \hat{\beta}_h^{(n)}(s, \mathbf{a}).$$
 - 10: Compute $\pi_h^{(n)}$ from equation 2 or equation 3 or equation 4. For each $s \in \mathcal{S}, i \in [M]$, set

$$\bar{V}_{h,i}^{(n)}(s) \leftarrow \left(\mathbb{D}_{\pi_h^{(n)}} \bar{Q}_{h,i}^{(n)} \right)(s), \quad \underline{V}_{h,i}^{(n)}(s) \leftarrow \left(\mathbb{D}_{\pi_h^{(n)}} \underline{Q}_{h,i}^{(n)} \right)(s), \quad \forall s \in \mathcal{S}.$$
 - 11: **end for**
 - 12: Let $\Delta^{(n)} = \max_{i \in [M]} \left\{ \bar{v}_i^{(n)} - \underline{v}_i^{(n)} \right\} + 2H\sqrt{A\zeta^{(n)}}$, where $\bar{v}_i^{(n)} = \int_{\mathcal{S}} \bar{V}_{1,i}^{(n)}(s) d_1(s) ds$, and
 $\underline{v}_i^{(n)} = \int_{\mathcal{S}} \underline{V}_{1,i}^{(n)}(s) d_1(s) ds.$
 - 13: **end for**
 - 14: **Return** $\hat{\pi} = \pi^{(n^*)}$ where $n^* = \arg \min_{n \in [N]} \Delta^{(n)}.$
-

Model-free Representation Learning In the model-free setting, we are only given the function class of the feature vectors, Φ_h , which we assume also includes the true feature ϕ_h^* . Given the dataset $\mathcal{D} := \mathcal{D}_h^{(n)} \cup \tilde{\mathcal{D}}_h^{(n)}$, MFREPLEARN aims to learn a feature vector that is able to linearly fit the Bellman backup of any function $f(s)$ in an appropriately chosen discriminator function class \mathcal{F}_h . To be precise, we aim to optimize the following objective:

$$\min_{\phi \in \Phi_h} \max_{f \in \mathcal{F}_h} \left[\min_{\theta} \mathbb{E}_{\mathcal{D}} \left[(\phi(s, \mathbf{a})^\top \theta - f(s'))^2 \right] - \min_{\tilde{\theta}, \tilde{\phi} \in \Phi_h} \mathbb{E}_{\mathcal{D}}^{(n)} \left[(\tilde{\phi}(s, \mathbf{a})^\top \tilde{\theta} - f(s'))^2 \right] \right],$$

where the first term is the empirical squared loss and the second term is the conditional expectation of $f(s')$ given (s, \mathbf{a}) , subtracted for the purpose of bias reduction. Once we obtain an estimation $\hat{\phi}_h^{(n)}$, we can construct a *non-parametric* transition model defined as:

$$\hat{P}_h^{(n)}(s'|s, \mathbf{a}) = \hat{\phi}_h^{(n)}(s, \mathbf{a})^\top \left(\sum_{(\tilde{s}, \tilde{\mathbf{a}}) \in \mathcal{D}} \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}})^\top + \lambda I_d \right)^{-1} \sum_{(\tilde{s}, \tilde{\mathbf{a}}, \tilde{s}') \in \mathcal{D}} \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \mathbf{1}_{\tilde{s}'=s'}. \quad (1)$$

We show that doing Least-square Value Iteration (LSVI) is equivalent to doing model-based planning inside $\hat{P}_h^{(n)}$ (line 10 of Alg. 1), and thus the model-free algorithm can be analyzed in the same way as the model-based algorithm. In practice, for applications where the raw observation states are high-dimensional, e.g. images, estimating the transition is often much harder than estimating the one-directional feature function. In such cases, we expect the Ψ class to be much larger than the Φ class and the model-free approach to be more efficient.

3.2 PLANNING

Based on the feature vector and transition probability computed from the representation learning phase, a new policy $\pi^{(n+1)}$ is computed using the planning module. The planning phase is conducted

Algorithm 2 Model-based Representation Learning, MBRepLearn

-
- 1: **Input:** Dataset \mathcal{D} , step h .
 - 2: Compute $(\hat{w}, \hat{\phi}) := \arg \max_{(w, \phi) \in \mathcal{M}_h} \mathbb{E}_{\mathcal{D}} [\log w(s')^\top \phi(s, \mathbf{a})]$.
 - 3: **Return** $\hat{\phi}, \hat{P} : \hat{P}(s'|s, \mathbf{a}) = \hat{w}(s')^\top \hat{\phi}(s, \mathbf{a})$.
-

Algorithm 3 Model-free Representation Learning, MFRepLearn

-
- 1: **Input:** Dataset \mathcal{D} , step h , regularization λ .
 - 2: Denote least squares loss: $\mathcal{L}_{\lambda, \mathcal{D}}(\phi, \theta, f) := \mathbb{E}_{\mathcal{D}} \left[(\phi(s, \mathbf{a})^\top \theta - f(s'))^2 \right] + \lambda \|\theta\|_2^2$.
 - 3: Compute $\hat{\phi} = \arg \min_{\phi \in \Phi_h} \max_{f \in \mathcal{F}_h} [\min_{\theta} \mathcal{L}_{\lambda, \mathcal{D}}(\phi, \theta, f) - \min_{\tilde{\phi} \in \Phi_h, \tilde{\theta}} \mathcal{L}_{\lambda, \mathcal{D}}(\tilde{\phi}, \tilde{\theta}, f)]$
 - 4: **Return** $\hat{\phi}, \hat{P}$ where \hat{P} is calculated from equation 1.
-

with a Upper-Confidence-Bound (UCB) style approach, and we maintain both an optimistic and a pessimistic estimation of the value functions and the Q-value functions $\bar{V}_{h,i}^{(n)}, \underline{V}_{h,i}^{(n)}, \bar{Q}_{h,i}^{(n)}, \underline{Q}_{h,i}^{(n)}$, which are computed recursively through the Bellman's equation with the bonus function $\hat{\beta}_h^{(n)}$ (Line 9 and 10 of Alg. 1). Here the operator \mathbb{D} is defined by $(\mathbb{D}_{\pi} f)(s) := \mathbb{E}_{\mathbf{a} \sim \pi(s)} [f(s, \mathbf{a})], \forall f : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, and $\pi_h^{(n)}$ is the policy computed from M induced Q-value functions $\tilde{Q}_{h,i}^{(n)}$. For the model-based setting, we simply let $\tilde{Q}_{h,i}^{(n)}$ be the optimistic estimator $\bar{Q}_{h,i}^{(n)}$. For the model-free setting, for technical reasons, we instead let $\tilde{Q}_{h,i}^{(n)}$ be the nearest neighbor of $\bar{Q}_{h,i}^{(n)}$ in \mathcal{N}_h with respect to the $\|\cdot\|_\infty$ metric, where $\mathcal{N}_h \subseteq \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ is a properly designed set of functions, whose construction is deferred to the appendix.

Depending on the problem settings, the policy $\pi_h^{(n)}$ takes either one of the following formulations:

- For the NE, we compute $\pi_h^{(n)} = (\pi_{h,1}^{(n)}, \pi_{h,2}^{(n)}, \dots, \pi_{h,M}^{(n)})$ such that $\forall s \in \mathcal{S}, i \in [M]$,

$$\pi_{h,i}^{(n)}(\cdot|s) = \arg \max_{\pi_{h,i}} \left(\mathbb{D}_{\pi_{h,i}, \pi_{h,-i}^{(n)}} \tilde{Q}_{h,i}^{(n)} \right) (s). \quad (2)$$

- For the CCE, we compute $\pi_h^{(n)}$ such that $\forall s \in \mathcal{S}, i \in [M]$,

$$\max_{\pi_{h,i}} \left(\mathbb{D}_{\pi_{h,i}, \pi_{h,-i}^{(n)}} \tilde{Q}_{h,i}^{(n)} \right) (s) \leq \left(\mathbb{D}_{\pi^{(n)}} \tilde{Q}_{h,i}^{(n)} \right) (s). \quad (3)$$

- For CE, we compute $\pi_h^{(n)}$ such that $\forall s \in \mathcal{S}, i \in [M]$,

$$\max_{\omega_{h,i} \in \Omega_{h,i}} \left(\mathbb{D}_{\omega_{h,i} \circ \pi_h^{(n)}} \tilde{Q}_{h,i}^{(n)} \right) (s) \leq \left(\mathbb{D}_{\pi^{(n)}} \tilde{Q}_{h,i}^{(n)} \right) (s). \quad (4)$$

Without loss of generality we assume the solution to the above formulations is unique, if there are multiple solutions, one can always adopt a deterministic selection rule such that it always outputs the same policy given the same inputs.

Note that although the policy is computed using only the optimistic estimations, we still maintain a pessimistic estimator, which is used to estimate the optimality gap $\Delta^{(n)}$ of the current policy. The algorithm's output policy $\hat{\pi}$ is chosen to be the one with the minimum estimated optimality gap.

The bonus term $\hat{\beta}_h^{(n)}$ is a linear bandit style bonus computed using the learned feature $\hat{\phi}$:

$$\hat{\beta}_h^{(n)}(s, \mathbf{a}) := \min \{ \alpha^{(n)} \|\hat{\phi}_h^{(n)}(s, \mathbf{a})\|_{(\hat{\Sigma}_h^{(n)})^{-1}}, H \}. \quad (5)$$

where $\hat{\Sigma}_h^{(n)} := \sum_{(s, \mathbf{a}) \in \mathcal{D}_h^{(n)}} \hat{\phi}_h^{(n)}(s, \mathbf{a}) \hat{\phi}_h^{(n)}(s, \mathbf{a})^\top + \lambda I_d$ is the empirical covariance matrix.

4 THEORETICAL RESULTS

In this section, we provide the theoretical guarantees of the proposed algorithm for both the model-based and model-free approaches. We denote $|\mathcal{M}| := \max_{h \in [H]} |\mathcal{M}_h|$ and $|\Phi| := \max_{h \in [H]} |\Phi_h|$. The first theorem provides a guarantee of the sample complexity for the model-based method.

Algorithm 4 Model-based Representation Learning for Factored MG (MBREPLEARN_FACTOR)

-
- 1: **Input:** Dataset \mathcal{D} , step h .
 - 2: Compute $(\hat{w}_i, \hat{\phi}_i) := \arg \max_{(w, \phi) \in \mathcal{M}_{h,i}} \mathbb{E}_{\mathcal{D}} [\log w(s'_i)^\top \phi(s[Z_i], \mathbf{a}_i)]$, for each $i \in [M]$.
 - 3: **Return** $\{\hat{\phi}_i\}_{i=1}^M, \hat{P} : \hat{P}(s'|s, \mathbf{a}) = \prod_{i=1}^M (\hat{w}_i(s'_i)^\top \hat{\phi}_i(s[Z_i], \mathbf{a}_i))$.
-

Theorem 4.1 (PAC guarantee of Algorithm 1 (model-based)). *When Alg. 1 is applied with model-based representation learning algorithm Alg. 2, with parameters $\lambda = \Theta(d \log(N^H |\Phi|/\delta))$, $\alpha^{(n)} = \Theta\left(Hd\sqrt{A \log(|\mathcal{M}|HN/\delta)}\right)$, $\zeta^{(n)} = \Theta(n^{-1} \log(|\mathcal{M}|HN/\delta))$, by setting the number of episodes N to be at most*

$$O\left(H^6 d^4 A^2 \varepsilon^{-2} \log^2(HdA|\mathcal{M}|/\delta\varepsilon)\right),$$

with probability $1 - \delta$, the output policy $\hat{\pi}$ is an ε -approximate $\{\text{NE}, \text{CCE}, \text{CE}\}$.

Theorem 4.1 shows that GERL_MG2 can find an ε -approximate $\{\text{NE}, \text{CCE}, \text{CE}\}$ by running the algorithm for at most $\tilde{O}(H^6 d^4 A^2 \varepsilon^{-2})$ episodes, which depends polynomially on the parameters $H, d, A, \varepsilon^{-1}$ and only has a logarithmic dependency on the cardinality of the model class $|\mathcal{M}|$. In particular, when reducing the Markov game to the single-agent MDP setting, the sample complexity of the model-based approach matches the result provided in (Uehara et al., 2021), which is known to have the best sample complexity among all oracle efficient algorithms for low-rank MDPs.

For model-free representation learning, we have the following guarantee:

Theorem 4.2 (PAC guarantee of Algorithm 1 (model-free)). *When Alg. 1 is applied with model-free representation learning algorithm Alg. 3, and $\lambda = \Theta(d \log(N^H |\Phi|/\delta))$, $\alpha^{(n)} = \Theta\left(HAd\sqrt{M \log(dNHAM|\Phi|/\delta)}\right)$, $\zeta^{(n)} = \Theta(d^2 An^{-1} \log(dNHAM|\Phi|/\delta))$, and the Markov game is a Block Markov game. When we set the number of episodes N to be at most*

$$O\left(H^6 d^4 A^3 M \varepsilon^{-2} \log^2(HdAM|\Phi|/\delta\varepsilon)\right),$$

for an appropriately designed function class $\{\mathcal{N}_h\}_{h=1}^H$ and discriminator class $\{\mathcal{F}_h\}_{h=1}^H$, with probability $1 - \delta$, the output policy $\hat{\pi}$ is an ε -approximate $\{\text{NE}, \text{CCE}, \text{CE}\}$.

For the model-free block Markov game setting, the number of episodes required to find an ε -approximate $\{\text{NE}, \text{CCE}, \text{CE}\}$ becomes $\tilde{O}(H^6 d^4 A^3 M \varepsilon^{-2})$. While it has a worse dependency compared with the model-based approach, the advantage of the model-free approach is it doesn't require the full model class of the transition probability but only the model class of the feature vector, which applies to a wider range of RL problems.

The proofs of Theorem 4.1 and Theorem 4.2 are deferred to Appendix B and C. Theorem 4.1 and Theorem 4.2 show that GERL_MG2 learns low-rank Markov games in a statistically efficient and oracle-efficient manner. We also remark that our modular analysis can be of independent theoretical interest. Unlike prior works that make heavy distinctions between model-based and model-free approaches, e.g. (Liu et al., 2021), we show that both approaches can be analyzed in a unified manner.

5 FACTORED MARKOV GAMES

The result in Theorem 4.1 is tractable in games with a moderate number of players. However, in applications with a large number of players, such as the scenario of autonomous traffic control, the total number of players in the game can be so large that the joint action space size $A = \tilde{A}^M$ dominates all other factors in the sample complexity bound. This exponential scaling with the number of players is sometimes referred to as the *curse of multi-player*. The only known class of algorithms that overcomes this challenge in Markov games is V-learning (see, e.g., Bai et al., 2020; Jin et al., 2021b), a value-based method that fits the V-function rather than the Q-function, thus removing the dependency on the action space size. However, V-learning only works for tabular Markov games with finite state and action spaces. Extending V-learning to the function approximation setting is extremely

non-trivial, because even in the single agent setting, no known algorithm can achieve sample efficient learning in MDPs while only performing function approximation on the V-function.

In this section we take a different approach that relies on the following observation. In a setting where the number of agents is large, there is often a spatial correlation among the agents, such that each agent’s local state is only immediately affected by the agent’s own action and the states of agents in its adjacency. For example, in smart traffic control, a vehicle’s local environment is only immediately affected by the states of the vehicles around it. On the other hand, it takes time for the course of actions of a vehicle from afar to propagate its influence on the vehicle of reference. Such spatial structure motivates the definition of a *factored* Markov Game.

In a factored Markov Game, each agent i has its local state s_i , whose transition is affected by agent i ’s action \mathbf{a}_i and the state of the agents in its neighborhood Z_i . We remark that the factored Markov Game structure still allows an agent to be affected by all other agents in the long run, as long as the directed graph defined by the neighborhood sets Z_i is connected. In particular, we have

Definition 5.1 (Low-Rank Factored Markov Game). *We call a Markov game a low-rank factored Markov game if for any $s, s' \in \mathcal{S}, \mathbf{a} \in \mathcal{A}, h \in [H], i \in [M]$, we have*

$$P_h^*(s'|s, \mathbf{a}) = \prod_{i=1}^M [\phi_{h,i}^*(s[Z_i], \mathbf{a}_i)^\top w_{h,i}^*(s'_i)].$$

where $Z_i \subseteq [M]$, $\phi_{h,i}^*(s[Z_i], \mathbf{a}_i), w_{h,i}^*(s'_i) \in \mathbb{R}^d$, $\|\phi_{h,i}^*(s[Z_i], \mathbf{a}_i)\|_2 \leq 1$ and $\|w_{h,i}^*(s'_i)\|_2 \leq \sqrt{d}$ for all $(s[Z_i], \mathbf{a}_i, s'_i)$. We assume $|Z_i| \leq L, \forall i \in [M]$. And we are given a group of model classes $\mathcal{M}_{h,i}, h \in [H], i \in [M]$ such that $(\phi_{h,i}^*, w_{h,i}^*) \in \mathcal{M}_{h,i}$.

We are now ready to present our algorithm and result in the low-rank factored Markov Game setting. Surprisingly, the same algorithm GERL_MG2 works in this setting, with the representation learning module Alg. 2 replaced by Alg. 4, and a few changes of variables. For simplicity, we focus on the model-based version. Define $\bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) = \otimes_{j \in Z_i} \hat{\phi}_{h,j}^{(n)}(s[Z_j], \mathbf{a}_j) \in \mathbb{R}^{d|Z_i|}$ where \otimes means the Kronecker product. Let

$$\hat{\beta}_h^{(n)}(s, \mathbf{a}) := \sum_{i=1}^M \min\{\alpha^{(n)} \|\bar{\phi}_{h,i}^{(n)}(s, \mathbf{a})\|_{(\bar{\Sigma}_{h,i}^{(n)})^{-1}, H}\}, \Delta^{(n)} := \max_{i \in [M]} \{\bar{v}_i^{(n)} - \underline{v}_i^{(n)}\} + 2HM\sqrt{\tilde{A}\zeta^{(n)}}$$

where $\bar{\Sigma}_{h,i}^{(n)} := \sum_{(s, \mathbf{a}) \in \mathcal{D}_h^{(n)}} \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a})^\top + \lambda I_{d|Z_i|}$. Then, GERL_MG2 with $\bar{\phi}$ and the newly defined $\hat{\beta}^{(n)}, \Delta^{(n)}$ achieves the following guarantee:

Theorem 5.1 (PAC guarantee of GERL_MG2 in Low-Rank Factored Markov Game). *When Alg. 1 is applied with model-based representation learning algorithm Alg. 4, with $L = O(1)$ and parameters $\lambda = \Theta(Ld^L \log(NHM|\Phi|/\delta)), \alpha^{(n)} = \Theta(H\tilde{A}d^L \sqrt{L \log(|\mathcal{M}|HN M/\delta)}), \zeta^{(n)} = \Theta(n^{-1} \log(|\mathcal{M}|HN M/\delta))$, by setting the number of episodes N to be at most*

$$O\left(M^4 H^6 d^{2(L+1)^2} \tilde{A}^{2(L+1)} \varepsilon^{-2} \log^2(HdALM|\mathcal{M}|/\delta\varepsilon)\right),$$

with probability $1 - \delta$, the output policy $\hat{\pi}$ is an ε -approximate $\{\text{NE}, \text{CCE}, \text{CE}\}$.

Remark 5.1. *This sample complexity only scales with $\exp(L)$ where L is the degree of the connection graph, which is assumed to be $O(1)$ in Definition 5.1 and in general much smaller than the total number of agents in practice. We remark that the factored structure is also previously studied in single-agent tabular MDPs (examples include Chen et al. (2020); Kearns and Koller (1999); Guestrin et al. (2002, 2003); Strehl et al. (2007)). Chen et al. (2020) provided a lower-bound showing that the exponential dependency on L is unimprovable in the worst case. Therefore, our bound here is also nearly tight, upto polynomial factors.*

6 EXPERIMENT

In this section we investigate our algorithm with proof-of-concept empirical studies. We design our testing bed using rich observation Markov game with arbitrary latent transitions and rewards. To

Table 1: **Top:** Short Horizon (H=3) exploitability of the final policy of DQN and GERL_MG2. **Bottom:** Long Horizon (H=10) exploitability of the final policy of DQN and GERL_MG2. Note that lower exploitability implies that the policy is closer to the NE policy.

	H=3 Environment 1	H=3 Environment 2	H=3 Environment 3
DQN	0.0851 (0.1152)	0.0877 (0.1961)	0.0090 (0.0200)
GERL_MG2	0.0013 (0.0018)	0.0032 (0.0032)	0.0004 (0.0009)
	H=10 Environment 1	H=10 Environment 2	H=10 Environment 3
DQN	0.2730 (0.3270)	0.0340 (0.0760)	0.0320 (0.0170)
GERL_MG2	0.0780 (0.1560)	0.0070 (0.0160)	0.0060 (0.0130)

solve the rich observation Markov game, an algorithm must correctly decode the latent structure (thus learning the dynamics) as well as solve the latent Markov game to find the NE/CE/CCE strategies concurrently. Below, we first introduce the setup of the experiments and then make comparisons with prior baselines in the two-player zero-sum setting. We then follow by showing the efficiency of GERL_MG2 in the general-sum setting. All further experiment details can be found in Appendix F. Here we focus on the model-free version of GERL_MG2. Specifically, we implement Algorithm 3 with deep learning libraries (Paszke et al., 2017). We defer more details to Appendix F.2.

Block Markov game Block Markov game is a multi-agent extension of single agent Block MDP, as defined in Def. 2.5. We design our Block Markov game by first randomly generating a tabular Markov game with horizon H , 3 states, 2 players each with 3 actions, and random reward matrix $R_h \in (0, 1)^{3 \times 3^2 \times H}$ and random transition matrix $T_h(s_h, a_h) \in \Delta_{S_{h+1}}$. We provide more details (e.g., generation of rich observation) in Appendix F.1.

Zero-sum Markov game In this section we first show the empirical evaluations under the two-player zero-sum Markov game setting. For an environment with horizon H , the randomly generated matrix R denotes the reward for player 1 and $-R^\top$ denotes the reward for player 2, respectively. For the zero-sum game setting, we designed two variants of Block Markov games: one with short horizon ($H = 3$) and one with long horizon ($H = 10$). We show in the following that GERL_MG2 works in both settings where the other baseline could only work in the short horizon setting.

Baseline We adopt one open-sourced implementation of DQN (Silver et al., 2016) with fictitious self-play (Heinrich et al., 2015).

We keep track of the *exploitability of the returned strategy* to evaluate the practical performances of the baselines. In the zero-sum setting, we only need to fix one agent (e.g., agent 2), train the other single agent (the exploiter) to maximize its corresponding return until convergence, and report the difference between the returns of the exploiter and the final return of the final policies. We include the exploitability in Table 1. We provide training curves in Appendix F.3 for completeness. We note that compared with the Deep RL baseline, GERL_MG2 shows a faster and more stable convergence in both environments, where the baseline is unstable during training and has a much larger exploitability.

General-sum Markov game. In this section we move on to the general-sum setting. To our best knowledge, our algorithm is the only principled algorithm that can be implemented on scale under the general-sum setting. For the general sum setting, we can not just compare our returned value to the oracle NE values, because multiple NE/CCE values may exist. Instead, we keep track of the exploitability of the policy and plot the training curve on the exploitability in Fig. 2 (deferred to Appendix F). Note that in this case we need to test both policies since their reward matrices are independently sampled.

7 DISCUSSION AND FUTURE WORKS

In this paper, we present the first algorithm that solves general-sum Markov games under function approximation. We provide both a model-based and a model-free variant of the algorithm and present a unified analysis. Empirically, we show that our algorithm outperforms existing deep RL baselines in a general benchmark with rich observation. Future work includes evaluating more challenging benchmarks and extending beyond the low-rank Markov game structure.

REPRODUCIBILITY STATEMENT

For theory, we provide proof and additional results in the Appendix. For empirical results, we provide implementation and environment details and hyperparameters in the Appendix. We also submit anonymous code in the supplemental materials.

ACKNOWLEDGMENTS

Mengdi Wang acknowledges the support by NSF grants DMS-1953686, IIS-2107304, CMMI1653435, ONR grant 1006977, and <http://C3.AI>.
Chi Jin gratefully acknowledges the support by Office of Naval Research Grant N00014-22-1-2253.

REFERENCES

- Alekh Agarwal, Mikael Henaff, Sham Kakade, and Wen Sun. Pc-pg: Policy cover directed exploration for provable policy gradient learning. *Advances in Neural Information Processing Systems*, 33: 13399–13412, 2020a.
- Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. Flambe: Structural complexity and representation learning of low rank mdps. *Advances in neural information processing systems*, 33:20095–20107, 2020b.
- Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *International Conference on Machine Learning*, pages 263–272. PMLR, 2017.
- Yu Bai and Chi Jin. Provable self-play algorithms for competitive reinforcement learning. In *International conference on machine learning*, pages 551–560. PMLR, 2020.
- Yu Bai, Chi Jin, and Tiancheng Yu. Near-optimal reinforcement learning with self-play. *Advances in neural information processing systems*, 33:2159–2170, 2020.
- Yu Bai, Chi Jin, Huan Wang, and Caiming Xiong. Sample-efficient learning of stackelberg equilibria in general-sum games. *Advances in Neural Information Processing Systems*, 34, 2021.
- Kimmo Berg and Tuomas Sandholm. Exclusion method for finding nash equilibrium in multiplayer games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębniak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- Avrim Blum, MohammadTaghi Hajiaghayi, Katrina Ligett, and Aaron Roth. Regret minimization and the price of total anarchy. In *Proceedings of the fortieth annual ACM symposium on Theory of computing*, pages 373–382, 2008.
- Manuele Brambilla, Eliseo Ferrante, Mauro Birattari, and Marco Dorigo. Swarm robotics: a review from the swarm engineering perspective. *Swarm Intelligence*, 7(1):1–41, 2013.
- Xiaoyu Chen, Jiachen Hu, Lihong Li, and Liwei Wang. Efficient reinforcement learning in factored mdps with application to constrained rl. *arXiv preprint arXiv:2008.13319*, 2020.
- Zixiang Chen, Dongruo Zhou, and Quanquan Gu. Almost optimal algorithms for two-player markov games with linear function approximation. *arXiv preprint arXiv:2102.07404*, 2021.
- Constantinos Daskalakis. On the complexity of approximating a nash equilibrium. *ACM Transactions on Algorithms (TALG)*, 9(3):1–35, 2013.
- Simon Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudik, and John Langford. Provably efficient rl with rich observations via latent state decoding. In *International Conference on Machine Learning*, pages 1665–1674. PMLR, 2019.

- Simon Du, Sham Kakade, Jason Lee, Shachar Lovett, Gaurav Mahajan, Wen Sun, and Ruosong Wang. Bilinear classes: A structural framework for provable generalization in rl. In *International Conference on Machine Learning*, pages 2826–2836. PMLR, 2021.
- Dylan Foster and Alexander Rakhlin. Beyond ucb: Optimal and efficient contextual bandits with regression oracles. In *International Conference on Machine Learning*, pages 3199–3210. PMLR, 2020.
- Dylan J Foster, Sham M Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity of interactive decision making. *arXiv preprint arXiv:2112.13487*, 2021.
- Carlos Guestrin, Relu Patrascu, and Dale Schuurmans. Algorithm-directed exploration for model-based reinforcement learning in factored mdps. In *ICML*, pages 235–242. Citeseer, 2002.
- Carlos Guestrin, Daphne Koller, Ronald Parr, and Shobha Venkataraman. Efficient solution algorithms for factored mdps. *Journal of Artificial Intelligence Research*, 19:399–468, 2003.
- Thomas Dueholm Hansen, Peter Bro Miltersen, and Uri Zwick. Strategy iteration is strongly polynomial for 2-player turn-based stochastic games with a constant discount factor. *Journal of the ACM (JACM)*, 60(1):1–16, 2013.
- Johannes Heinrich, Marc Lanctot, and David Silver. Fictitious self-play in extensive-form games. In *International conference on machine learning*, pages 805–813. PMLR, 2015.
- Junling Hu and Michael P Wellman. Nash q-learning for general-sum stochastic games. *Journal of machine learning research*, 4(Nov):1039–1069, 2003.
- Baihe Huang, Jason D Lee, Zhaoran Wang, and Zhuoran Yang. Towards general function approximation in zero-sum markov games. *arXiv preprint arXiv:2107.14702*, 2021.
- Zeyu Jia, Lin F Yang, and Mengdi Wang. Feature-based q-learning for two-player stochastic games. *arXiv preprint arXiv:1906.00423*, 2019.
- Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Contextual decision processes with low bellman rank are pac-learnable. In *International Conference on Machine Learning*, pages 1704–1713. PMLR, 2017.
- Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is q-learning provably efficient? *Advances in neural information processing systems*, 31, 2018.
- Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR, 2020.
- Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman eluder dimension: New rich classes of rl problems, and sample-efficient algorithms. *Advances in Neural Information Processing Systems*, 34, 2021a.
- Chi Jin, Qinghua Liu, Yuanhao Wang, and Tiancheng Yu. V-learning—a simple, efficient, decentralized algorithm for multiagent rl. *arXiv preprint arXiv:2110.14555*, 2021b.
- Chi Jin, Qinghua Liu, and Tiancheng Yu. The power of exploiter: Provable multi-agent rl in large state spaces. *arXiv preprint arXiv:2106.03352*, 2021c.
- Michael Kearns and Daphne Koller. Efficient reinforcement learning in factored mdps. In *IJCAI*, volume 16, pages 740–747, 1999.
- Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In *Machine learning proceedings 1994*, pages 157–163. Elsevier, 1994.
- Qinghua Liu, Tiancheng Yu, Yu Bai, and Chi Jin. A sharp analysis of model-based reinforcement learning with self-play. In *International Conference on Machine Learning*, pages 7001–7010. PMLR, 2021.

- Dipendra Misra, Mikael Henaff, Akshay Krishnamurthy, and John Langford. Kinematic state abstraction and provably efficient rich-observation reinforcement learning. In *International conference on machine learning*, pages 6961–6971. PMLR, 2020.
- Aditya Modi, Jinglin Chen, Akshay Krishnamurthy, Nan Jiang, and Alekh Agarwal. Model-free representation learning and exploration in low-rank mdps. *arXiv preprint arXiv:2102.07035*, 2021.
- Matej Moravčík, Martin Schmid, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. Deepstack: Expert-level artificial intelligence in heads-up no-limit poker. *Science*, 356(6337):508–513, 2017.
- Christos H Papadimitriou and Tim Roughgarden. Computing correlated equilibria in multi-player games. *Journal of the ACM (JACM)*, 55(3):1–29, 2008.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- Shuang Qiu, Lingxiao Wang, Chenjia Bai, Zhuoran Yang, and Zhaoran Wang. Contrastive ucbl: Provably efficient contrastive self-supervised learning in online reinforcement learning. In *International Conference on Machine Learning*, pages 18168–18210. PMLR, 2022.
- Tongzheng Ren, Tianjun Zhang, Csaba Szepesvári, and Bo Dai. A free lunch from the noise: Provable and practical exploration for representation learning. *arXiv preprint arXiv:2111.11485*, 2021.
- Tim Roughgarden. Algorithmic game theory. *Communications of the ACM*, 53(7):78–86, 2010.
- Shai Shalev-Shwartz, Shaked Shammah, and Amnon Shashua. Safe, multi-agent, reinforcement learning for autonomous driving. *arXiv preprint arXiv:1610.03295*, 2016.
- Lloyd S Shapley. Stochastic games. *Proceedings of the national academy of sciences*, 39(10):1095–1100, 1953.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- Alexander L Strehl, Carlos Diuk, and Michael L Littman. Efficient structure learning in factored-state mdps. In *AAAI*, volume 7, pages 645–650, 2007.
- Masatoshi Uehara, Xuezhou Zhang, and Wen Sun. Representation learning for online and offline rl in low-rank mdps. *arXiv preprint arXiv:2110.04652*, 2021.
- Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- Qiaomin Xie, Yudong Chen, Zhaoran Wang, and Zhuoran Yang. Learning zero-sum simultaneous-move markov games using function approximation and correlated equilibrium. In *Conference on learning theory*, pages 3674–3682. PMLR, 2020.
- Lin Yang and Mengdi Wang. Sample-optimal parametric q-learning using linearly additive features. In *International Conference on Machine Learning*, pages 6995–7004. PMLR, 2019.
- Andrea Zanette, David Brandfonbrener, Emma Brunskill, Matteo Pirota, and Alessandro Lazaric. Frequentist regret bounds for randomized least-squares value iteration. In *International Conference on Artificial Intelligence and Statistics*, pages 1954–1964. PMLR, 2020.
- Andrea Zanette, Ching-An Cheng, and Alekh Agarwal. Cautiously optimistic policy optimization and exploration with linear function approximation. In *Conference on Learning Theory*, pages 4473–4525. PMLR, 2021.

Kaiqing Zhang, Sham Kakade, Tamer Basar, and Lin Yang. Model-based multi-agent rl in zero-sum markov games with near-optimal sample complexity. *Advances in Neural Information Processing Systems*, 33:1166–1178, 2020.

Xuezhou Zhang, Yuda Song, Masatoshi Uehara, Mengdi Wang, Wen Sun, and Alekh Agarwal. Efficient reinforcement learning in block mdps: A model-free representation learning approach. *arXiv preprint arXiv:2202.00063*, 2022.

A ADDITIONAL NOTATIONS

Given a (possibly not normalized) transition probability $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H] \rightarrow [0, 1]$ and a policy $\pi : \mathcal{S} \times [H] \rightarrow \Delta_{\mathcal{A}}$, we define the density function of the state-action pair (s, \mathbf{a}) at step h under transition P and π by

$$d_{P,1}^{\pi}(s, \mathbf{a}) := d_1(s)\pi_1(\mathbf{a}|s), \quad d_{P,h+1}^{\pi}(s, \mathbf{a}) := \sum_{\tilde{s} \in \mathcal{S}, \tilde{\mathbf{a}} \in \mathcal{A}} d_{P,h}^{\pi}(\tilde{s}, \tilde{\mathbf{a}})P_h(s|\tilde{s}, \tilde{\mathbf{a}})\pi_{h+1}(\mathbf{a}|s), \forall h \geq 1.$$

We abuse the notations a bit and denote $d_{P,h}^{\pi}(s)$ as the marginalized state distribution, i.e., $d_{P,h}^{\pi}(s) = \sum_{\mathbf{a} \in \mathcal{A}} d_{P,h}^{\pi}(s, \mathbf{a})$. For any $n \in [N], h \in [H]$, define

$$\begin{aligned} \rho_h^{(n)}(s, \mathbf{a}) &= \frac{1}{n} \sum_{i=1}^n d_{P^*,h}^{\pi^{(i)}}(s)u_{\mathcal{A}}(\mathbf{a}), \\ \tilde{\rho}_h^{(n)}(s, \mathbf{a}) &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\tilde{s} \sim d_{P^*,h-1}^{\pi^{(i)}}, \tilde{\mathbf{a}} \sim U(\mathcal{A})} [P^*(s|\tilde{s}, \tilde{\mathbf{a}})u_{\mathcal{A}}(\mathbf{a})], \\ \gamma_h^{(n)}(s, \mathbf{a}) &= \frac{1}{n} \sum_{i=1}^n d_{P^*,h}^{\pi^{(i)}}(s, \mathbf{a}). \end{aligned}$$

When we use the expectation $\mathbb{E}_{(s,\mathbf{a}) \sim \rho} [f(s, \mathbf{a})]$ (or $\mathbb{E}_{s \sim \rho} [f(s)]$) for some (possibly not normalized) distribution ρ and function f , we simply mean $\sum_{s \in \mathcal{S}, \mathbf{a} \in \mathcal{A}} \rho(s, \mathbf{a})f(s, \mathbf{a})$ (or $\sum_{s \in \mathcal{S}} \rho(s)f(s)$) so that the expectation can be naturally extended to the unnormalized distributions. For an iteration n , a distribution ρ and a feature ϕ , we denote the expected feature covariance as

$$\Sigma_{n,\rho,\phi} = n\mathbb{E}_{(s,\mathbf{a}) \sim \rho} [\phi(s, \mathbf{a})\phi(s, \mathbf{a})^{\top}] + \lambda I_d.$$

Meanwhile, define the empirical covariance by

$$\hat{\Sigma}_{h,\phi}^{(n)} := \sum_{(s,\mathbf{a}) \in \mathcal{D}_h^{(n)}} \phi(s, \mathbf{a})\phi(s, \mathbf{a})^{\top} + \lambda I_d.$$

B ANALYSIS OF THE MODEL-BASED METHOD

B.1 HIGH PROBABILITY EVENTS

We define the following event

$$\begin{aligned} \mathcal{E}_1 : \forall n \in [N], h \in [H], \rho \in \left\{ \rho_h^{(n)}, \tilde{\rho}_h^{(n)} \right\}, \quad \mathbb{E}_{(s,\mathbf{a}) \sim \rho} \left[\left\| \hat{P}_h^{(n)}(\cdot|s, \mathbf{a}) - P_h^*(\cdot|s, \mathbf{a}) \right\|_1^2 \right] \leq \zeta^{(n)}, \\ \mathcal{E}_2 : \forall n \in [N], h \in [H], \phi_h \in \Phi_h, s \in \mathcal{S}, \mathbf{a} \in \mathcal{A}, \quad \left\| \phi_h(s, \mathbf{a}) \right\|_{\left(\hat{\Sigma}_{h,\phi_h}^{(n)} \right)^{-1}} = \Theta \left(\left\| \phi_h(s, \mathbf{a}) \right\|_{\Sigma_{n,\rho_h^{(n)},\phi_h}^{-1}} \right) \\ \mathcal{E} := \mathcal{E}_1 \cap \mathcal{E}_2. \end{aligned}$$

To prove \mathcal{E} holds with a high probability, we first introduce the following MLE guarantee, whose original version can be found in (Agarwal et al., 2020b):

Lemma B.1 (MLE guarantee). *For a fixed episode n and any step h , with probability $1 - \delta$,*

$$\mathbb{E}_{(s,\mathbf{a}) \sim \{0.5\rho_h^{(n)} + 0.5\tilde{\rho}_h^{(n)}\}} \left[\left\| \hat{P}_h^{(n)}(\cdot|s, \mathbf{a}) - P_h^*(\cdot|s, \mathbf{a}) \right\|_1^2 \right] \lesssim \frac{1}{n} \log \frac{|\mathcal{M}|}{\delta}.$$

As a straightforward corollary, with probability $1 - \delta$,

$$\forall n \in \mathbb{N}^+, \forall h \in [H], \quad \mathbb{E}_{(s,\mathbf{a}) \sim \{0.5\rho_h^{(n)} + 0.5\tilde{\rho}_h^{(n)}\}} \left[\left\| \hat{P}_h^{(n)}(\cdot|s, \mathbf{a}) - P_h^*(\cdot|s, \mathbf{a}) \right\|_1^2 \right] \lesssim \frac{1}{n} \log \frac{nH|\mathcal{M}|}{\delta}. \quad (6)$$

Proof. See Agarwal et al. (Agarwal et al., 2020b) (Theorem 21). \square

Based on Lemma B.1 and Lemma E.1 in Appendix E, we directly get the following guarantee:

Lemma B.2. When $\hat{P}_h^{(n)}$ is computed using Alg. 2, if we set

$$\lambda = \Theta \left(d \log \frac{NH|\Phi|}{\delta} \right), \quad \zeta^{(n)} = \Theta \left(\frac{1}{n} \log \frac{|\mathcal{M}|HN}{\delta} \right),$$

then \mathcal{E} holds with probability at least $1 - \delta$.

B.2 STATISTICAL GUARANTEES

Lemma B.3 (One-step back inequality for the learned model). *Suppose the event \mathcal{E} holds. Consider a set of functions $\{g_h\}_{h=1}^H$ that satisfies $g_h \in \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}_+$, s.t. $\|g_h\|_\infty \leq B$. For any given policy π , we have*

$$\begin{aligned} & \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^\pi} [g_h(s, \mathbf{a})] \\ & \leq \begin{cases} \sqrt{A \mathbb{E}_{(s,\mathbf{a}) \sim \rho_1^{(n)}} [g_1^2(s, \mathbf{a})]}, & h = 1 \\ \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)},h-1}^\pi} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n,\rho_{h-1}^{(n)},\hat{\phi}_{h-1}^{(n)}}^{-1}} \sqrt{n A \mathbb{E}_{(s,\mathbf{a}) \sim \hat{\rho}_h^{(n)}} [g_h^2(s, \mathbf{a})] + B^2 \lambda d + B^2 n \zeta^{(n)}}, B \right\} \right], & h \geq 2 \end{cases} \end{aligned}$$

Recall $\Sigma_{n,\rho_h^{(n)},\hat{\phi}_h^{(n)}} = n \mathbb{E}_{(s,\mathbf{a}) \sim \rho_h^{(n)}} [\hat{\phi}_h^{(n)}(s, \mathbf{a}) \hat{\phi}_h^{(n)}(s, \mathbf{a})^\top] + \lambda I_d$.

Proof. For step $h = 1$, we have

$$\begin{aligned} \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},1}^\pi} [g_1(s, \mathbf{a})] &= \mathbb{E}_{s \sim d_1, \mathbf{a} \sim \pi_1(s)} [g_1(s, \mathbf{a})] \\ &\leq \sqrt{\max_{(s,\mathbf{a})} \frac{d_1(s) \pi_1(\mathbf{a}|s)}{\rho_1^{(n)}(s, \mathbf{a})} \mathbb{E}_{(s',\mathbf{a}') \sim \rho_1^{(n)}} [g_1^2(s', \mathbf{a}')] } \\ &= \sqrt{\max_{(s,\mathbf{a})} \frac{d_1(s) \pi_1(\mathbf{a}|s)}{d_1(s) u_{\mathcal{A}}(\mathbf{a})} \mathbb{E}_{(s',\mathbf{a}') \sim \rho_1^{(n)}} [g_1^2(s', \mathbf{a}')] } \\ &\leq \sqrt{A \mathbb{E}_{(s,\mathbf{a}) \sim \rho_1^{(n)}} [g_1^2(s, \mathbf{a})]}. \end{aligned}$$

For step $h = 2, \dots, H - 1$, we observe the following one-step-back decomposition:

$$\begin{aligned} & \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^\pi} [g_h(s, \mathbf{a})] \\ &= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)},h-1}^\pi, s \sim \hat{P}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim \pi_h(s)} [g_h(s, \mathbf{a})] \\ &= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)},h-1}^\pi} \left[\hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}})^\top \int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right] \\ &= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)},h-1}^\pi} \left[\min \left\{ \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}})^\top \int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds, B \right\} \right] \\ &\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)},h-1}^\pi} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n,\rho_{h-1}^{(n)},\hat{\phi}_{h-1}^{(n)}}^{-1}} \left\| \int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right\|_{\Sigma_{n,\rho_{h-1}^{(n)},\hat{\phi}_{h-1}^{(n)}}}, B \right\} \right]. \end{aligned}$$

where we use the fact that g_h is bounded by B . Then,

$$\left\| \int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right\|_{\Sigma_{n,\rho_{h-1}^{(n)},\hat{\phi}_{h-1}^{(n)}}}^2$$

$$\begin{aligned}
&\leq \left(\int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right)^\top \left(n \mathbb{E}_{(s, \mathbf{a}) \sim \rho_{h-1}^{(n)}} [\hat{\phi}_{h-1}^{(n)}(s, \mathbf{a}) \hat{\phi}_{h-1}^{(n)}(s, \mathbf{a})^\top] + \lambda I_d \right) \left(\int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right) \\
&\leq n \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \rho_{h-1}^{(n)}} \left[\left(\int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right)^2 \right] + B^2 \lambda d \\
&\quad \left(\left\| \sum_{\mathbf{a} \in \mathcal{A}} \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) \right\|_\infty \leq B \text{ and by assumption } \left\| \hat{w}_{h-1}^{(n)}(s) \right\|_2 \leq \sqrt{d} \right) \\
&= n \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \rho_{h-1}^{(n)}} \left[\left(\mathbb{E}_{s \sim P_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim \pi_h(s)} [g_h(s, \mathbf{a})] \right)^2 \right] + B^2 \lambda d \\
&\leq n \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \rho_{h-1}^{(n)}} \left[\left(\mathbb{E}_{s \sim P_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim \pi_h(s)} [g_h(s, \mathbf{a})] \right)^2 \right] + B^2 \lambda d + n B^2 \xi^{(n)} \quad (\text{Event } \mathcal{E}) \\
&\leq n \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \rho_{h-1}^{(n)}, s \sim P_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim \pi_h(s)} [g_h^2(s, \mathbf{a})] + B^2 \lambda d + B^2 n \xi^{(n)}. \quad (\text{Jensen}) \\
&\leq n A \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \rho_{h-1}^{(n)}, s \sim P_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim U(\mathcal{A})} [g_h^2(s, \mathbf{a})] + B^2 \lambda d + B^2 n \zeta^{(n)} \quad (\text{Importance sampling}) \\
&\leq n A \mathbb{E}_{(s, \mathbf{a}) \sim \tilde{\rho}_h^{(n)}} [g_h^2(s, \mathbf{a})] + B^2 \lambda d + B^2 n \zeta^{(n)}. \quad (\text{Definition of } \tilde{\rho}_h^{(n)})
\end{aligned}$$

Combing the above results together, we get

$$\begin{aligned}
&\mathbb{E}_{(s, \mathbf{a}) \sim d_{P^{(n)}, h}^\pi} [g_h(s, \mathbf{a})] \\
&\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\tilde{P}^{(n)}, h-1}^\pi} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}^{-1}} \left\| \int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}}, B \right\} \right] \\
&\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\tilde{P}^{(n)}, h-1}^\pi} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}^{-1}} \sqrt{n A \mathbb{E}_{(s, \mathbf{a}) \sim \tilde{\rho}_h^{(n)}} [g_h^2(s, \mathbf{a})] + B^2 \lambda d + B^2 n \zeta^{(n)}}, B \right\} \right],
\end{aligned}$$

which has finished the proof. \square

Lemma B.4 (One-step back inequality for the true model). *Consider a set of functions $\{g_h\}_{h=1}^H$ that satisfies $g_h \in \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}_+$, s.t. $\|g_h\|_\infty \leq B$. Then for any given policy π , we have*

$$\begin{aligned}
&\mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^\pi} [g_h(s, \mathbf{a})] \\
&\leq \begin{cases} \sqrt{A \mathbb{E}_{(s, \mathbf{a}) \sim \rho_1^{(n)}} [g_1^2(s, \mathbf{a})]}, & h = 1 \\ \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h-1}^\pi} \left[\left\| \phi_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_{h-1}^{(n)}, \phi_{h-1}^*}^{-1}} \sqrt{n A \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} [g_h^2(s, \mathbf{a})] + B^2 \lambda d} \right], & h \geq 2 \end{cases}
\end{aligned}$$

Recall $\Sigma_{n, \gamma_h^{(n)}, \phi_h^*} = n \mathbb{E}_{(s, \mathbf{a}) \sim \gamma_h^{(n)}} [\phi_h^*(s, \mathbf{a}) \phi_h^*(s, \mathbf{a})^\top] + \lambda I_d$.

Proof. For step $h = 1$, we have

$$\begin{aligned}
\mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, 1}^\pi} [g_1(s, \mathbf{a})] &= \mathbb{E}_{s \sim d_1, \mathbf{a} \sim \pi_1(s)} [g_1(s, \mathbf{a})] \\
&\leq \sqrt{\max_{(s, \mathbf{a})} \frac{d_1(s) \pi_1(\mathbf{a}|s)}{\rho_1^{(n)}(s, \mathbf{a})} \mathbb{E}_{(s', \mathbf{a}') \sim \rho_1^{(n)}} [g_1^2(s', \mathbf{a}')] } \\
&= \sqrt{\max_{(s, \mathbf{a})} \frac{d_1(s) \pi_1(\mathbf{a}|s)}{d_1(s) u_{\mathcal{A}}(\mathbf{a})} \mathbb{E}_{(s', \mathbf{a}') \sim \rho_1^{(n)}} [g_1^2(s', \mathbf{a}')] } \\
&\leq \sqrt{A \mathbb{E}_{(s, \mathbf{a}) \sim \rho_1^{(n)}} [g_1^2(s, \mathbf{a})]}.
\end{aligned}$$

For step $h = 2, \dots, H - 1$, we observe the following one-step-back decomposition:

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^\pi} [g_h(s, \mathbf{a})]$$

$$\begin{aligned}
&= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h-1}^\pi, s \sim P_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim \pi_h(s)} [g_h(s, \mathbf{a})] \\
&= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h-1}^\pi} \left[\phi_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}})^\top \int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} w_{h-1}^*(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right] \\
&\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h-1}^\pi} \left[\left\| \phi_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}} \left\| \int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} w_{h-1}^*(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right\|_{\Sigma} \right].
\end{aligned}$$

Then,

$$\begin{aligned}
&\left\| \int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} w_{h-1}^*(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right\|_{\Sigma}^2 \\
&\leq \left(\int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} w_{h-1}^*(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right)^\top \left(n \mathbb{E}_{(s, \mathbf{a}) \sim \gamma_{h-1}^{(n)}} [\phi_{h-1}^*(s, \mathbf{a}) \phi_{h-1}^*(s, \mathbf{a})^\top] + \lambda I_d \right) \left(\int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} w_{h-1}^*(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right) \\
&\leq n \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \gamma_{h-1}^{(n)}} \left[\left(\int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} w_{h-1}^*(s) \phi_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right)^2 \right] + B^2 \lambda d \\
&\quad \text{(Use the assumption } \left\| \sum_{\mathbf{a} \in \mathcal{A}} \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) \right\|_\infty \leq B \text{ and } \|\phi_{h-1}^*(s)\|_2 \leq \sqrt{d}.) \\
&= n \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \gamma_{h-1}^{(n)}} \left[\left(\mathbb{E}_{s \sim P_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim \pi_h(s)} [g_h(s, \mathbf{a})] \right)^2 \right] + B^2 \lambda d \\
&\leq n \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \gamma_{h-1}^{(n)}, s \sim P_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim \pi_h(s)} [g_h^2(s, \mathbf{a})] + B^2 \lambda d \quad \text{(Jensen)} \\
&\leq n A \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \gamma_{h-1}^{(n)}, s \sim P_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim U(\mathcal{A})} [g_h^2(s, \mathbf{a})] + B^2 \lambda d \quad \text{(Importance sampling)} \\
&\leq n A \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} [g_h^2(s, \mathbf{a})] + B^2 \lambda d, \quad \text{(Definition of } \rho_h^{(n)})
\end{aligned}$$

Combing the above results together, we get

$$\begin{aligned}
&\mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^\pi} [g_h(s, \mathbf{a})] \\
&= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h-1}^\pi, s \sim P_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim \pi_h(s)} [g_h(s, \mathbf{a})] \\
&\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h-1}^\pi} \left[\left\| \phi_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}} \left\| \int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} w_{h-1}^*(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right\|_{\Sigma} \right] \\
&\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h-1}^\pi} \left[\left\| \phi_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}} \sqrt{n A \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} [g_h^2(s, \mathbf{a})] + B^2 \lambda d} \right],
\end{aligned}$$

which has finished the proof. \square

Lemma B.5 (Optimism for NE and CCE). *Consider an episode $n \in [N]$ and set $\alpha^{(n)} = \Theta \left(H \sqrt{n A \zeta^{(n)}} + d \lambda \right)$. When the event \mathcal{E} holds and the policy $\pi^{(n)}$ is computed by solving NE or CCE, we have*

$$\bar{v}_i^{(n)}(s) - v_i^{\dagger, \pi^{(n)}}(s) \geq -H \sqrt{A \zeta^{(n)}}, \quad \forall n \in [N], i \in [M].$$

Proof. Define $\tilde{\mu}_{h,i}^{(n)}(\cdot|s) := \arg \max_{\mu} \left(\mathbb{D}_{\mu, \pi_{h,-i}^{(n)}} Q_{h,i}^{\dagger, \pi_{h,-i}^{(n)}} \right) (s)$ as the best response policy for player i at step h , and let $\tilde{\pi}_h^{(n)} = \tilde{\mu}_{h,i}^{(n)} \times \pi_{h,-i}^{(n)}$. Let $f_h^{(n)}(s, \mathbf{a}) = \left\| \hat{P}_h^{(n)}(\cdot|s, \mathbf{a}) - P_h^*(\cdot|s, \mathbf{a}) \right\|_1$, then according to the event \mathcal{E} , we have

$$\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}, \quad \mathbb{E}_{(s, \mathbf{a}) \sim \tilde{\rho}_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}, \quad \forall n \in [N], h \in [H],$$

$$\|\phi_h(s, \mathbf{a})\|_{\left(\hat{\Sigma}_{h, \phi_h}^{(n)}\right)^{-1}} = \Theta \left(\left\| \phi_h(s, \mathbf{a}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \phi_h}^{-1}} \right), \quad \forall n \in [N], h \in [H], \phi_h \in \Phi_h.$$

A direct conclusion of the event \mathcal{E} is we can find an absolute constant c , such that

$$\begin{aligned} \beta_h^{(n)}(s, \mathbf{a}) &= \min \left\{ \alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\left(\hat{\Sigma}_{h, \hat{\phi}_h}^{(n)}\right)^{-1}}, H \right\} \\ &\geq \min \left\{ c\alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \hat{\phi}_h}^{-1}}, H \right\}, \quad \forall n \in [N], h \in [H]. \end{aligned}$$

Next, we prove by induction that

$$\mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\bar{V}_{h,i}^{(n)}(s) - V_{h,i}^{\dagger, \pi^{(n)}}(s) \right] \geq \sum_{h'=h}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h'}^{\pi^{(n)}}} \left[\hat{\beta}_{h'}^{(n)}(s, \mathbf{a}) - H \min\{f_{h'}^{(n)}(s, \mathbf{a}), 1\} \right], \quad \forall h \in [H]. \quad (7)$$

First, notice that $\forall h \in [H]$,

$$\begin{aligned} \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\bar{V}_{h,i}^{(n)}(s) - V_{h,i}^{\dagger, \pi^{(n)}}(s) \right] &= \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\left(\mathbb{D}_{\pi_h^{(n)}} \bar{Q}_{h,i}^{(n)} \right)(s) - \left(\mathbb{D}_{\pi_h^{(n)}} Q_{h,i}^{\dagger, \pi^{(n)}} \right)(s) \right] \\ &\geq \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\left(\mathbb{D}_{\pi_h^{(n)}} \bar{Q}_{h,i}^{(n)} \right)(s) - \left(\mathbb{D}_{\pi_h^{(n)}} Q_{h,i}^{\dagger, \pi^{(n)}} \right)(s) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\bar{Q}_{h,i}^{(n)}(s, \mathbf{a}) - Q_{h,i}^{\dagger, \pi^{(n)}}(s, \mathbf{a}) \right], \end{aligned}$$

where the inequality uses the fact that $\pi_h^{(n)}$ is the NE (or CCE) solution for $\left\{ \bar{Q}_{h,i}^{(n)} \right\}_{i=1}^M$. Now we are ready to prove equation 7:

- When $h = H$, we have

$$\begin{aligned} \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, H}^{\pi^{(n)}}} \left[\bar{V}_{H,i}^{(n)}(s) - V_{H,i}^{\dagger, \pi^{(n)}}(s) \right] &\geq \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\pi^{(n)}}} \left[\bar{Q}_{H,i}^{(n)}(s, \mathbf{a}) - Q_{H,i}^{\dagger, \pi^{(n)}}(s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\pi^{(n)}}} \left[\hat{\beta}_H^{(n)}(s, \mathbf{a}) \right] \\ &\geq \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\pi^{(n)}}} \left[\hat{\beta}_H^{(n)}(s, \mathbf{a}) - H \min\{f_H^{(n)}(s, \mathbf{a}), 1\} \right]. \end{aligned}$$

- Suppose the statement is true for step $h+1$, then for step h , we have

$$\begin{aligned} &\mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\bar{V}_{h,i}^{(n)}(s) - V_{h,i}^{\dagger, \pi^{(n)}}(s) \right] \\ &\geq \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\bar{Q}_{h,i}^{(n)}(s, \mathbf{a}) - Q_{h,i}^{\dagger, \pi^{(n)}}(s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \bar{V}_{h+1,i}^{(n)} \right)(s, \mathbf{a}) - \left(P_h^* V_{h+1,i}^{\dagger, \pi^{(n)}} \right)(s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \left(\bar{V}_{h+1,i}^{(n)} - V_{h+1,i}^{\dagger, \pi^{(n)}} \right) \right)(s, \mathbf{a}) + \left(\left(\hat{P}_h^{(n)} - P_h^* \right) V_{h+1,i}^{\dagger, \pi^{(n)}} \right)(s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\left(\hat{P}_h^{(n)} - P_h^* \right) V_{h+1,i}^{\dagger, \pi^{(n)}} \right)(s, \mathbf{a}) \right] + \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h+1}^{\pi^{(n)}}} \left[\bar{V}_{h+1,i}^{(n)}(s) - V_{h+1,i}^{\dagger, \pi^{(n)}}(s) \right] \\ &\geq \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - H \min\{f_h^{(n)}(s, \mathbf{a}), 1\} \right] + \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h+1}^{\pi^{(n)}}} \left[\bar{V}_{h+1,i}^{(n)}(s) - V_{h+1,i}^{\dagger, \pi^{(n)}}(s) \right] \end{aligned}$$

$$\geq \sum_{h'=h}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h'}}^{\pi^{(n)}} \left[\hat{\beta}_{h'}^{(n)}(s, \mathbf{a}) - H \min \left\{ f_{h'}^{(n)}(s, \mathbf{a}), 1 \right\} \right],$$

where we use the fact

$$\begin{aligned} \left| \left(\hat{P}_h^{(n)} - P_h^* \right) V_{h+1, i}^{\dagger, \pi^{(n)}}(s, \mathbf{a}) \right| &\leq \min \left\{ H, \left\| \hat{P}_h^{(n)}(\cdot | s, \mathbf{a}) - P_h^*(\cdot | s, \mathbf{a}) \right\|_1 \left\| V_{h+1, i}^{\dagger, \pi^{(n)}} \right\|_\infty \right\} \\ &\leq H \min \left\{ 1, \left\| \hat{P}_h^{(n)}(\cdot | s, \mathbf{a}) - P_h^*(\cdot | s, \mathbf{a}) \right\|_1 \right\} \\ &= H \min \left\{ 1, f_{h'}^{(n)}(s, \mathbf{a}) \right\} \end{aligned}$$

and the last row uses the induction assumption.

Therefore, we have proved equation 7. We then apply $h = 1$ to equation 7, and get

$$\begin{aligned} \mathbb{E}_{s \sim d_1} \left[\bar{V}_{1, i}^{(n)}(s) - V_{1, i}^{\dagger, \pi^{(n)}}(s) \right] &= \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, 1}}^{\pi^{(n)}} \left[\bar{V}_{1, i}^{(n)}(s) - V_{1, i}^{\dagger, \pi^{(n)}}(s) \right] \\ &\geq \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}}^{\pi^{(n)}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - H \min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right] \\ &= \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}}^{\pi^{(n)}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] - H \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}}^{\pi^{(n)}} \left[\min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right]. \end{aligned}$$

Next we are going to bound the second term, let $g_h(s, \mathbf{a}) = \min \{ f_h^{(n)}(s, \mathbf{a}), 1 \}$ and apply Lemma B.3 to g_h , we have for $h = 1$,

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, 1}}^{\pi^{(n)}} \left[\min \left\{ f_1^{(n)}(s, \mathbf{a}), 1 \right\} \right] \leq \sqrt{A \mathbb{E}_{(s, \mathbf{a}) \sim \rho_1^{(n)}} \left[\left(f_1^{(n)}(s, \mathbf{a}) \right)^2 \right]} \leq \sqrt{A \zeta^{(n)}}.$$

And $\forall h \geq 2$, we have

$$\begin{aligned} &\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}}^{\pi^{(n)}} \left[\min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right] \\ &\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}}^{\pi^{(n)}} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}} \sqrt{n A \mathbb{E}_{(s, \mathbf{a}) \sim \tilde{\rho}_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right]} + d\lambda + n\zeta^{(n)}, 1 \right\} \right] \\ &\lesssim \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}}^{\pi^{(n)}} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}} \sqrt{n A \zeta^{(n)} + d\lambda + n\zeta^{(n)}}, 1 \right\} \right]. \end{aligned}$$

Note that we here use the fact $\min \{ f_h^{(n)}(s, \mathbf{a}), 1 \} \leq 1$, $\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}$ and

$$\mathbb{E}_{(s, \mathbf{a}) \sim \tilde{\rho}_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}. \text{ Then according to our choice of } \alpha^{(n)}, \text{ we get}$$

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}}^{\pi^{(n)}} \left[\min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right] \leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}}^{\pi^{(n)}} \left[\min \left\{ \frac{c\alpha^{(n)}}{H} \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}}, 1 \right\} \right].$$

Combining all things together,

$$\begin{aligned} \bar{v}_i^{(n)} - v_i^{\dagger, \pi^{(n)}} &= \mathbb{E}_{s \sim d_1} \left[\bar{V}_{1, i}^{(n)}(s) - V_{1, i}^{\dagger, \pi^{(n)}}(s) \right] \\ &\geq \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}}^{\pi^{(n)}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] - H \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}}^{\pi^{(n)}} \left[\min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right] \end{aligned}$$

$$\begin{aligned} &\geq \sum_{h=1}^{H-1} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h}^{\tilde{\pi}^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - \min \left\{ c\alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \hat{\phi}_h^{(n)}}^{-1}}, H \right\} \right] - H\sqrt{A\zeta^{(n)}} \\ &\geq -H\sqrt{A\zeta^{(n)}}, \end{aligned}$$

which proves the inequality. \square

Lemma B.6 (Optimism for CE). *Consider an episode $n \in [N]$ and set $\alpha^{(n)} = \Theta\left(H\sqrt{nA\zeta^{(n)}} + d\lambda\right)$. When the event \mathcal{E} holds, we have*

$$\bar{v}_i^{(n)}(s) - \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi^{(n)}}(s) \geq -H\sqrt{A\zeta^{(n)}}, \quad \forall n \in [N], i \in [M].$$

Proof. Denote $\tilde{\omega}_{h,i}^{(n)} = \arg \max_{\omega_h \in \Omega_{h,i}} \left(\mathbb{D}_{\omega_h \circ \pi_h^{(n)}} \max_{\omega \in \Omega_i} Q_{h,i}^{\omega \circ \pi^{(n)}} \right)(s)$ and let $\tilde{\pi}_h^{(n)} = \tilde{\omega}_{h,i} \circ \pi_h^{(n)}$.

Let $f_h^{(n)}(s, \mathbf{a}) = \left\| \hat{P}_h^{(n)}(\cdot | s, \mathbf{a}) - P_h^*(\cdot | s, \mathbf{a}) \right\|_1$, then according to the event \mathcal{E} , we have

$$\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}, \quad \mathbb{E}_{(s, \mathbf{a}) \sim \tilde{\rho}_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}, \quad \forall n \in [N], h \in [H],$$

$$\|\phi_h(s, \mathbf{a})\|_{\left(\hat{\Sigma}_{h, \phi_h}^{(n)}\right)^{-1}} = \Theta \left(\left\| \phi_h(s, \mathbf{a}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \phi_h}^{-1}} \right), \quad \forall n \in [N], h \in [H], \phi_h \in \Phi_h.$$

A direct conclusion of the event \mathcal{E} is we can find an absolute constant c , such that

$$\begin{aligned} \beta_h^{(n)}(s, \mathbf{a}) &= \min \left\{ \alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\left(\hat{\Sigma}_{h, \hat{\phi}_h}^{(n)}\right)^{-1}}, H \right\} \\ &\geq \min \left\{ c\alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \hat{\phi}_h}^{-1}}, H \right\}, \quad \forall n \in [N], h \in [H]. \end{aligned}$$

Next, we prove by induction that

$$\mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\tilde{\pi}^{(n)}}} \left[\bar{V}_{h,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{h,i}^{\omega \circ \pi^{(n)}}(s) \right] \geq \sum_{h'=h}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h'}^{\tilde{\pi}^{(n)}}} \left[\hat{\beta}_{h'}^{(n)}(s, \mathbf{a}) - H \min \left\{ f_{h'}^{(n)}(s, \mathbf{a}), 1 \right\} \right], \quad \forall h \in [H]. \quad (8)$$

First, notice that $\forall h \in [H]$,

$$\begin{aligned} \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\tilde{\pi}^{(n)}}} \left[\bar{V}_{h,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{h,i}^{\omega \circ \pi^{(n)}}(s) \right] &= \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\tilde{\pi}^{(n)}}} \left[\left(\mathbb{D}_{\tilde{\pi}_h^{(n)}} \bar{Q}_{h,i}^{(n)} \right)(s) - \left(\mathbb{D}_{\tilde{\pi}_h^{(n)}} \max_{\omega \in \Omega_i} Q_{h,i}^{\omega \circ \pi^{(n)}} \right)(s) \right] \\ &\geq \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\tilde{\pi}^{(n)}}} \left[\left(\mathbb{D}_{\tilde{\pi}_h^{(n)}} \bar{Q}_{h,i}^{(n)} \right)(s) - \left(\mathbb{D}_{\tilde{\pi}_h^{(n)}} \max_{\omega \in \Omega_i} Q_{h,i}^{\omega \circ \pi^{(n)}} \right)(s) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\tilde{\pi}^{(n)}}} \left[\bar{Q}_{h,i}^{(n)}(s, \mathbf{a}) - \max_{\omega \in \Omega_i} Q_{h,i}^{\omega \circ \pi^{(n)}}(s, \mathbf{a}) \right]. \end{aligned}$$

where the inequality uses the fact that $\pi_h^{(n)}$ is the CE solution for $\left\{ \bar{Q}_{h,i}^{(n)} \right\}_{i=1}^M$. Now we are ready to prove equation 8:

- When $h = H$, we have

$$\begin{aligned} \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, H}^{\tilde{\pi}^{(n)}}} \left[\bar{V}_{H,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{H,i}^{\omega \circ \pi^{(n)}}(s) \right] &\geq \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\tilde{\pi}^{(n)}}} \left[\bar{Q}_{H,i}^{(n)}(s, \mathbf{a}) - \max_{\omega \in \Omega_i} Q_{H,i}^{\omega \circ \pi^{(n)}}(s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\tilde{\pi}^{(n)}}} \left[\hat{\beta}_H^{(n)}(s, \mathbf{a}) \right] \\ &\geq \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\tilde{\pi}^{(n)}}} \left[\hat{\beta}_H^{(n)}(s, \mathbf{a}) - H \min \left\{ f_H^{(n)}(s, \mathbf{a}), 1 \right\} \right]. \end{aligned}$$

- Suppose the statement is true for $h + 1$, then for step h , we have

$$\begin{aligned}
& \mathbb{E}_{s \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\overline{V}_{h,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{h,i}^{\omega \circ \pi^{(n)}}(s) \right] \\
& \geq \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\overline{Q}_{h,i}^{(n)}(s, \mathbf{a}) - \max_{\omega \in \Omega_i} Q_{h,i}^{\omega \circ \pi^{(n)}}(s, \mathbf{a}) \right] \\
& = \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \overline{V}_{h+1,i}^{(n)} \right)(s, \mathbf{a}) - \left(P_h^* \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}} \right)(s, \mathbf{a}) \right] \\
& = \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \left(\overline{V}_{h+1,i}^{(n)} - \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}} \right) \right)(s, \mathbf{a}) + \left(\left(\hat{P}_h^{(n)} - P_h^* \right) \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}} \right)(s, \mathbf{a}) \right] \\
& = \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\left(\hat{P}_h^{(n)} - P_h^* \right) \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}} \right)(s, \mathbf{a}) \right] \\
& \quad + \mathbb{E}_{s \sim d_{\hat{P}^{(n)},h+1}^{\pi^{(n)}}} \left[\overline{V}_{h+1,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}}(s) \right] \\
& \geq \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - H \min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right] + \mathbb{E}_{s \sim d_{\hat{P}^{(n)},h+1}^{\pi^{(n)}}} \left[\overline{V}_{h+1,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}}(s) \right] \\
& \geq \sum_{h'=h}^H \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h'}^{\pi^{(n)}}} \left[\hat{\beta}_{h'}^{(n)}(s, \mathbf{a}) - H \min \left\{ f_{h'}^{(n)}(s, \mathbf{a}), 1 \right\} \right],
\end{aligned}$$

where we use the fact

$$\begin{aligned}
\left| \left(\hat{P}_h^{(n)} - P_h^* \right) \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}} \right| (s, \mathbf{a}) & \leq \min \left\{ H, \left\| \hat{P}_h^{(n)}(\cdot | s, \mathbf{a}) - P_h^*(\cdot | s, \mathbf{a}) \right\|_1 \left\| \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}} \right\|_\infty \right\} \\
& \leq H \min \left\{ 1, \left\| \hat{P}_h^{(n)}(\cdot | s, \mathbf{a}) - P_h^*(\cdot | s, \mathbf{a}) \right\|_1 \right\} \\
& = H \min \left\{ 1, f_{h'}^{(n)}(s, \mathbf{a}) \right\}
\end{aligned}$$

and the last row uses the induction assumption.

Therefore, we have proved equation 8. We then apply $h = 1$ to equation 8, and get

$$\begin{aligned}
\mathbb{E}_{s \sim d_1} \left[\overline{V}_{1,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{1,i}^{\omega \circ \pi^{(n)}}(s) \right] & = \mathbb{E}_{s \sim d_{\hat{P}^{(n)},1}^{\pi^{(n)}}} \left[\overline{V}_{1,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{1,i}^{\omega \circ \pi^{(n)}}(s) \right] \\
& \geq \sum_{h=1}^H \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - H \min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right] \\
& = \sum_{h=1}^H \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] - H \sum_{h=1}^H \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right].
\end{aligned}$$

Next we are going to bound the second term, let $g_h(s, \mathbf{a}) = \min \{ f_h^{(n)}(s, \mathbf{a}), 1 \}$ and apply Lemma B.3 to g_h , we have for $h = 1$,

$$\mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},1}^{\pi^{(n)}}} \left[\min \left\{ f_1^{(n)}(s, \mathbf{a}), 1 \right\} \right] \leq \sqrt{A \mathbb{E}_{(s,\mathbf{a}) \sim \rho_1^{(n)}} \left[\left(f_1^{(n)}(s, \mathbf{a}) \right)^2 \right]} \leq \sqrt{A \zeta^{(n)}}.$$

And $\forall h \geq 2$, we have

$$\begin{aligned}
& \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right] \\
& \leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)},h-1}^{\pi^{(n)}}} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}} \sqrt{n A \mathbb{E}_{(s,\mathbf{a}) \sim \tilde{\rho}_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right]} + d\lambda + n\zeta^{(n)}, 1 \right\} \right]
\end{aligned}$$

$$\lesssim \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\pi^{(n)}}} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}^{-1}} \sqrt{nA\zeta^{(n)} + d\lambda + n\zeta^{(n)}}, 1 \right\} \right].$$

Note that we here use the fact $\min\{f_h^{(n)}(s, \mathbf{a}), 1\} \leq 1$, $\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}$ and $\mathbb{E}_{(s, \mathbf{a}) \sim \tilde{\rho}_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}$. Then according to our choice of $\alpha^{(n)}$, we get

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[f_h^{(n)}(s, \mathbf{a}) \right] \leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\pi^{(n)}}} \left[\min \left\{ \frac{c\alpha^{(n)}}{H} \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}^{-1}}, 1 \right\} \right].$$

Combining all things together,

$$\begin{aligned} \bar{v}_i^{(n)} - \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi^{(n)}} &= \mathbb{E}_{s \sim d_1} \left[\bar{V}_{1,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{1,i}^{\omega \circ \pi^{(n)}}(s) \right] \\ &\geq \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] - H \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right] \\ &\geq \sum_{h=1}^{H-1} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - \min \left\{ c\alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \hat{\phi}_h^{(n)}}^{-1}}, H \right\} \right] - H\sqrt{A\zeta^{(n)}} \\ &\geq -H\sqrt{A\zeta^{(n)}}, \end{aligned}$$

which proves the inequality. \square

Lemma B.7 (Pessimism). Consider an episode $n \in [N]$ and set $\alpha^{(n)} = \Theta \left(H\sqrt{nA\zeta^{(n)} + d\lambda} \right)$. When the event \mathcal{E} holds, we have

$$\underline{v}_i^{(n)}(s) - v_i^{\pi^{(n)}}(s) \leq H\sqrt{A\zeta^{(n)}}, \quad \forall n \in [N], i \in [M].$$

Proof. Let $f_h^{(n)}(s, \mathbf{a}) = \left\| \hat{P}_h^{(n)}(\cdot | s, \mathbf{a}) - P_h^*(\cdot | s, \mathbf{a}) \right\|_1$, then according to the event \mathcal{E} , we have

$$\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}, \quad \mathbb{E}_{(s, \mathbf{a}) \sim \tilde{\rho}_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}, \quad \forall n \in [N], h \in [H],$$

$$\left\| \phi_h(s, \mathbf{a}) \right\|_{\left(\hat{\Sigma}_{h, \phi_h}^{(n)} \right)^{-1}} = \Theta \left(\left\| \phi_h(s, \mathbf{a}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \phi_h}^{-1}} \right), \quad \forall n \in [N], h \in [H], \phi_h \in \Phi_h.$$

A direct conclusion of the event \mathcal{E} is we can find an absolute constant c , such that

$$\begin{aligned} \beta_h^{(n)}(s, \mathbf{a}) &= \min \left\{ \alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\left(\hat{\Sigma}_{h, \hat{\phi}_h}^{(n)} \right)^{-1}}, H \right\} \\ &\geq \min \left\{ c\alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \hat{\phi}_h}^{-1}}, H \right\}, \quad \forall n \in [N], h \in [H]. \end{aligned}$$

Again, we prove the following inequality by induction:

$$\mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[V_{h,i}^{(n)}(s) - V_{h,i}^{\pi^{(n)}}(s) \right] \leq \sum_{h'=h}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h'}^{\pi^{(n)}}} \left[-\hat{\beta}_{h'}^{(n)}(s, \mathbf{a}) + H \min \left\{ f_{h'}^{(n)}(s, \mathbf{a}), 1 \right\} \right], \quad \forall h \in [H]. \quad (9)$$

- When $h = H$, we have

$$\begin{aligned} \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, H}^{\pi^{(n)}}} \left[\underline{V}_{H,i}^{(n)}(s) - V_{H,i}^{\pi^{(n)}}(s) \right] &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\pi^{(n)}}} \left[\underline{Q}_{H,i}^{(n)}(s, \mathbf{a}) - Q_{H,i}^{\pi^{(n)}}(s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\pi^{(n)}}} \left[-\hat{\beta}_H^{(n)}(s, \mathbf{a}) \right] \\ &\leq \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\pi^{(n)}}} \left[-\hat{\beta}_H^{(n)}(s, \mathbf{a}) + H \min \left\{ f_H^{(n)}(s, \mathbf{a}), 1 \right\} \right] \end{aligned}$$

- Suppose the statement is true for $h + 1$, then for step h , we have

$$\begin{aligned} &\mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\underline{V}_{h,i}^{(n)}(s) - V_{h,i}^{\pi^{(n)}}(s) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\underline{Q}_{h,i}^{(n)}(s, \mathbf{a}) - Q_{h,i}^{\pi^{(n)}}(s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[-\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \underline{V}_{h+1,i}^{(n)} \right) (s, \mathbf{a}) - \left(P_h^* V_{h+1,i}^{\pi^{(n)}} \right) (s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[-\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \left(\underline{V}_{h+1,i}^{(n)} - V_{h+1,i}^{\pi^{(n)}} \right) \right) (s, \mathbf{a}) + \left(\left(\hat{P}_h^{(n)} - P_h^* \right) V_{h+1,i}^{\pi^{(n)}} \right) (s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[-\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\left(\hat{P}_h^{(n)} - P_h^* \right) V_{h+1,i}^{\pi^{(n)}} \right) (s, \mathbf{a}) \right] + \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h+1}^{\pi^{(n)}}} \left[\left(\underline{V}_{h+1,i}^{(n)} - V_{h+1,i}^{\pi^{(n)}} \right) (s) \right] \\ &\leq \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[-\hat{\beta}_h^{(n)}(s, \mathbf{a}) + H \min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right] + \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h+1}^{\pi^{(n)}}} \left[\left(\underline{V}_{h+1,i}^{(n)} - V_{h+1,i}^{\pi^{(n)}} \right) (s) \right] \\ &\leq \sum_{h'=h}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h'}^{\pi^{(n)}}} \left[-\hat{\beta}_{h'}^{(n)}(s, \mathbf{a}) + H \min \left\{ f_{h'}^{(n)}(s, \mathbf{a}), 1 \right\} \right]. \end{aligned}$$

where we use the fact

$$\begin{aligned} \left| \left(\hat{P}_h^{(n)} - P_h^* \right) V_{h+1,i}^{\pi^{(n)}} \right| (s, \mathbf{a}) &\leq \min \left\{ H, \left\| \hat{P}_h^{(n)}(\cdot | s, \mathbf{a}) - P_h^*(\cdot | s, \mathbf{a}) \right\|_1 \left\| V_{h+1,i}^{\pi^{(n)}} \right\|_\infty \right\} \\ &\leq H \min \left\{ 1, \left\| \hat{P}_h^{(n)}(\cdot | s, \mathbf{a}) - P_h^*(\cdot | s, \mathbf{a}) \right\|_1 \right\} \\ &= H \min \left\{ 1, f_{h'}^{(n)}(s, \mathbf{a}) \right\} \end{aligned}$$

and the last row uses the induction assumption.

The remaining steps are exactly the same as the proof in Lemma B.5 or Lemma B.6, we may prove

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, 1}^{\pi^{(n)}}} \left[\min \left\{ f_1^{(n)}(s, \mathbf{a}), 1 \right\} \right] \leq \sqrt{A\zeta^{(n)}},$$

and

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[f_h^{(n)}(s, \mathbf{a}) \right] \leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\pi^{(n)}}} \left[\min \left\{ \frac{c\alpha^{(n)}}{H} \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}^{-1}}, 1 \right\}, 1 \right], \quad \forall h \geq 2.$$

Combining all things together, we get

$$\begin{aligned} \underline{v}_i^{(n)} - v_i^{\pi^{(n)}} &= \mathbb{E}_{s \sim d_1} \left[\underline{V}_{1,i}^{(n)}(s) - V_{1,i}^{\pi^{(n)}}(s) \right] \\ &\leq \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[-\hat{\beta}_h^{(n)}(s, \mathbf{a}) + H \min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right] \\ &\leq \sum_{h=1}^{H-1} \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[-\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \min \left\{ c\alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \hat{\phi}_h^{(n)}}^{-1}}, H \right\} \right] + H \sqrt{A\zeta^{(n)}} \\ &\leq H \sqrt{A\zeta^{(n)}}, \end{aligned}$$

which has finished the proof. \square

Lemma B.8. For the model-based algorithm, when we pick $\lambda = \Theta\left(d \log \frac{NH|\Phi|}{\delta}\right)$, $\zeta^{(n)} = \Theta\left(\frac{1}{n} \log \frac{|\mathcal{M}|HN}{\delta}\right)$ and $\alpha^{(n)} = \Theta\left(H\sqrt{nA\zeta^{(n)} + d\lambda}\right)$, with probability $1 - \delta$, we have

$$\sum_{n=1}^N \Delta^{(n)} \lesssim H^3 d^2 AN^{\frac{1}{2}} \log \frac{|\mathcal{M}|HN}{\delta}.$$

Proof. With our choice of λ and $\zeta^{(n)}$, according to Lemma B.2, we know \mathcal{E} holds with probability $1 - \delta$. Furthermore, we have

$$\alpha^{(n)} = \Theta\left(H\sqrt{A \log \frac{|\mathcal{M}|HN}{\delta} + d^2 \log \frac{NH|\Phi|}{\delta}}\right) = O\left(dH\sqrt{A \log \frac{|\mathcal{M}|HN}{\delta}}\right)$$

Let $f_h^{(n)}(s, \mathbf{a}) = \left\| \hat{P}_h^{(n)}(\cdot | s, \mathbf{a}) - P_h^*(\cdot | s, \mathbf{a}) \right\|_1$. According to the definition of the event \mathcal{E} , we have

$$\mathbb{E}_{s \sim f_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}, \|\phi_h(s, \mathbf{a})\|_{(\hat{\Sigma}_{h, \phi_h}^{(n)})^{-1}} = \Theta\left(\|\phi_h(s, \mathbf{a})\|_{\Sigma_{n, \rho_h^{(n)}, \phi_h}^{-1}}\right), \quad \forall n \in [N], h \in [H], \phi_h \in \Phi_h. \quad (10)$$

By definition, we have

$$\Delta^{(n)} = \max_{i \in [M]} \left\{ \bar{v}_i^{(n)} - \underline{v}_i^{(n)} \right\} + 2H\sqrt{A\zeta^{(n)}}.$$

For each fixed $i \in [M]$, $h \in [H]$ and $n \in [N]$, we have

$$\begin{aligned} & \mathbb{E}_{s \sim d_{P^*, h}^{\pi^{(n)}}} \left[\bar{V}_{h, i}^{(n)}(s) - \underline{V}_{h, i}^{(n)}(s) \right] \\ &= \mathbb{E}_{s \sim d_{P^*, h}^{\pi^{(n)}}} \left[\left(\mathbb{D}_{\pi_h^{(n)}} \bar{Q}_{h, i}^{(n)} \right)(s) - \left(\mathbb{D}_{\pi_h^{(n)}} \underline{Q}_{h, i}^{(n)} \right)(s) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[\bar{Q}_{h, i}^{(n)}(s, \mathbf{a}) - \underline{Q}_{h, i}^{(n)}(s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[2\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \left(\bar{V}_{h+1, i}^{(n)} - \underline{V}_{h+1, i}^{(n)} \right) \right)(s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[2\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\left(\hat{P}_h^{(n)} - P_h^* \right) \left(\bar{V}_{h+1, i}^{(n)} - \underline{V}_{h+1, i}^{(n)} \right) \right)(s, \mathbf{a}) \right] + \mathbb{E}_{s \sim d_{P^*, h+1}^{\pi^{(n)}}} \left[\bar{V}_{h+1, i}^{(n)}(s) - \underline{V}_{h+1, i}^{(n)}(s) \right] \\ &\leq \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[2\hat{\beta}_h^{(n)}(s, \mathbf{a}) + 2H^2 f_h^{(n)}(s, \mathbf{a}) \right] + \mathbb{E}_{s \sim d_{P^*, h+1}^{\pi^{(n)}}} \left[\bar{V}_{h+1, i}^{(n)}(s) - \underline{V}_{h+1, i}^{(n)}(s) \right]. \end{aligned}$$

Note that we use the fact $\bar{V}_{h+1, i}^{(n)}(s) - \underline{V}_{h+1, i}^{(n)}(s)$ is upper bounded by $2H^2$, which can be proved easily using induction using the fact that $\hat{\beta}_h^{(n)}(s, \mathbf{a}) \leq H$. Applying the above formula recursively to $\mathbb{E}_{s \sim d_{P^*, h+1}^{\pi^{(n)}}} \left[\bar{V}_{h+1, i}^{(n)}(s) - \underline{V}_{h+1, i}^{(n)}(s) \right]$, one gets the following result (or more formally, one can prove by induction, just like what we did in Lemma B.5, Lemma B.6 and Lemma B.7):

$$\mathbb{E}_{s \sim d_{P^*, 1}^{\pi^{(n)}}} \left[\bar{V}_{1, i}^{(n)}(s) - \underline{V}_{1, i}^{(n)}(s) \right] \leq \underbrace{2 \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right]}_{(a)} + \underbrace{2H^2 \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[f_h^{(n)}(s, \mathbf{a}) \right]}_{(b)}. \quad (11)$$

First, we calculate the first term (a) in Inequality equation 11. Following Lemma B.4 and noting the bonus $\hat{\beta}_h^{(n)}$ is $O(H)$, we have

$$\sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right]$$

$$\begin{aligned}
&\lesssim \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi(n)}} \left[\min \left\{ \alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(s, \mathbf{a}) \right\|_{\Sigma^{-1}_{n, \rho_h^{(n)}, \hat{\phi}_h^{(n)}}}, H \right\} \right] && \text{(From equation 10)} \\
&\lesssim \sum_{h=1}^{H-1} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}_{n, \gamma_h^{(n)}, \phi_h^*}} \right] \sqrt{nA (\alpha^{(n)})^2 \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left\| \hat{\phi}_h^{(n)}(s, \mathbf{a}) \right\|_{\Sigma^{-1}_{n, \rho_h^{(n)}, \hat{\phi}_h^{(n)}}}^2 \right]} + H^2 d \lambda \\
&+ \sqrt{A (\alpha^{(n)})^2 \mathbb{E}_{(s, \mathbf{a}) \sim \rho_1^{(n)}} \left[\left\| \hat{\phi}_1^{(n)}(s, \mathbf{a}) \right\|_{\Sigma^{-1}_{n, \rho_1^{(n)}, \hat{\phi}_1^{(n)}}}^2 \right]}.
\end{aligned}$$

Note that we use the fact that $B = H$ when applying Lemma B.4. In addition, we have

$$\begin{aligned}
&n \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left\| \hat{\phi}_h^{(n)}(s, \mathbf{a}) \right\|_{\Sigma^{-1}_{n, \rho_h^{(n)}, \hat{\phi}_h^{(n)}}}^2 \right] \\
&= n \text{Tr} \left(\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\hat{\phi}_h^{(n)}(s, \mathbf{a}) \hat{\phi}_h^{(n)}(s, \mathbf{a})^\top \right] \left(n \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\hat{\phi}_h^{(n)}(s, \mathbf{a}) \hat{\phi}_h^{(n)}(s, \mathbf{a})^\top \right] + \lambda I_d \right)^{-1} \right) \\
&\leq d.
\end{aligned}$$

Then,

$$\sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi(n)}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] \leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}_{n, \gamma_h^{(n)}, \phi_h^*}} \right] \sqrt{dA (\alpha^{(n)})^2 + H^2 d \lambda} + \sqrt{dA (\alpha^{(n)})^2 / n}.$$

Second, we calculate the term (b) in inequality equation 11. Following Lemma B.4 and noting that $f_h^{(n)}(s, \mathbf{a})$ is upper-bounded by 2 (i.e., $B = 2$ in Lemma B.4), we have

$$\begin{aligned}
&\sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi(n)}} [f_h^{(n)}(s, \mathbf{a})] \\
&\leq \sum_{h=1}^{H-1} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}_{n, \gamma_h^{(n)}, \phi_h^*}} \right] \sqrt{nA \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right]} + d \lambda + \sqrt{A \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_1^{(n)}(s, \mathbf{a}) \right)^2 \right]} \\
&\leq \sum_{h=1}^{H-1} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}_{n, \gamma_h^{(n)}, \phi_h^*}} \right] \sqrt{nA \zeta^{(n)} + d \lambda} + \sqrt{A \zeta^{(n)}} \\
&\lesssim \frac{\alpha^{(n)}}{H} \sum_{h=1}^{H-1} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}_{n, \gamma_h^{(n)}, \phi_h^*}} \right] + \sqrt{A \zeta^{(n)}},
\end{aligned}$$

where in the second inequality, we use $\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}$, and in the last line,

recall $\sqrt{nA \zeta^{(n)} + d \lambda} \lesssim \alpha^{(n)} / H$. Then, by combining the above calculation of the term (a) and term (b) in inequality equation 11, we have:

$$\begin{aligned}
\bar{v}_i^{(n)} - v_i^{(n)} &= \mathbb{E}_{s \sim d_{P^*, 1}^{\pi(n)}} \left[\bar{V}_{1, i}^{(n)}(s) - V_{1, i}^{(n)}(s) \right] \\
&\lesssim \sum_{h=1}^{H-1} \left(\mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}_{n, \gamma_h^{(n)}, \phi_h^*}} \right] \sqrt{dA (\alpha^{(n)})^2 + H^2 d \lambda} + \sqrt{\frac{dA (\alpha^{(n)})^2}{n}} \right) \\
&\quad + H^2 \sum_{h=1}^{H-1} \left(\frac{\alpha^{(n)}}{H} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}_{n, \gamma_h^{(n)}, \phi_h^*}} \right] + \sqrt{A \zeta^{(n)}} \right).
\end{aligned}$$

Taking maximum over i on both sides and using the definition of $\Delta^{(n)}$, we get

$$\Delta^{(n)} = \max_{i \in [M]} \left\{ \bar{v}_i^{(n)} - v_i^{(n)} \right\} + 2H \sqrt{A \zeta^{(n)}}$$

$$\begin{aligned} &\lesssim \sum_{h=1}^{H-1} \left(\mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[\left\| \phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_h^{(n)}, \phi_h^*}^{-1}} \right] \sqrt{dA (\alpha^{(n)})^2 + H^2 d\lambda} + \sqrt{\frac{dA (\alpha^{(n)})^2}{n}} \right) \\ &\quad + H^2 \sum_{h=1}^{H-1} \left(\frac{\alpha^{(n)}}{H} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[\left\| \phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_h^{(n)}, \phi_h^*}^{-1}} \right] + \sqrt{A\zeta^{(n)}} \right). \end{aligned}$$

Hereafter, we take the dominating term out. Note that

$$\begin{aligned} \sum_{n=1}^N \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[\left\| \phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_h^{(n)}, \phi_h^*}^{-1}} \right] &\leq \sqrt{N \sum_{n=1}^N \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[\phi_h^*(\tilde{s}, \tilde{\mathbf{a}})^\top \Sigma_{n, \gamma_h^{(n)}, \phi_h^*}^{-1} \phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \right]} \\ &\quad \text{(CS inequality)} \\ &\lesssim \sqrt{N \left(\log \det \left(\sum_{n=1}^N \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi^{(n)}}} [\phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \phi_h^*(\tilde{s}, \tilde{\mathbf{a}})^\top] \right) - \log \det(\lambda I_d) \right)} \quad \text{(Lemma E.2)} \\ &\leq \sqrt{dN \log \left(1 + \frac{N}{d\lambda} \right)}. \end{aligned}$$

(Potential function bound, Lemma E.3 noting $\|\phi_h^*(s, \mathbf{a})\|_2 \leq 1$ for any (s, \mathbf{a}) .)

Finally,

$$\begin{aligned} \sum_{n=1}^N \Delta^{(n)} &\lesssim H \left(\sqrt{dN \log \left(1 + \frac{N}{d} \right)} \sqrt{dA (\alpha^{(N)})^2 + H^2 d\lambda} + \sum_{n=1}^N \sqrt{\frac{dA (\alpha^{(n)})^2}{n}} \right) \\ &\quad + H^3 \left(\frac{1}{H} \sqrt{dN \log \left(1 + \frac{N}{d\lambda} \right)} \alpha^{(N)} + \sum_{n=1}^N \sqrt{A\zeta^{(n)}} \right) \\ &\lesssim H^2 d \sqrt{NA \log \left(1 + \frac{N}{d\lambda} \right)} \alpha^{(N)} \\ &\quad \text{(Some algebra. We take the dominating term out. Note that } \alpha^{(n)} \text{ is increasing in } n) \\ &\lesssim H^3 d^2 AN^{\frac{1}{2}} \log \frac{|\mathcal{M}|HN}{\delta}. \end{aligned}$$

This concludes the proof. \square

Proof of Theorem 4.1

Proof. For any fixed episode n and agent i , by Lemma B.5, Lemma B.6 and Lemma B.7, we have

$$v_i^{\dagger, \pi_i^{(n)}} - v_i^{\pi_i^{(n)}} \left(\text{or } \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi_i^{(n)}} - v_i^{\pi_i^{(n)}} \right) \leq \bar{v}_i^{(n)} - \underline{v}_i^{(n)} + 2H \sqrt{A\zeta^{(n)}} \leq \Delta^{(n)}.$$

Taking maximum over i on both sides, we have

$$\max_{i \in [M]} \left\{ v_i^{\dagger, \pi_i^{(n)}} - v_i^{\pi_i^{(n)}} \right\} \left(\text{or } \max_{i \in [M]} \left\{ \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi_i^{(n)}} - v_i^{\pi_i^{(n)}} \right\} \right) \leq \Delta^{(n)}. \quad (12)$$

From Lemma B.8, with probability $1 - \delta$, we can ensure

$$\sum_{n=1}^N \Delta^{(n)} \lesssim H^3 d^2 AN^{\frac{1}{2}} \log \frac{|\mathcal{M}|HN}{\delta}.$$

Therefore, according to Lemma E.4, when we pick N to be

$$O \left(\frac{H^6 d^4 A^2}{\varepsilon^2} \log^2 \left(\frac{HdA|\mathcal{M}|}{\delta\varepsilon} \right) \right),$$

we have

$$\frac{1}{N} \sum_{n=1}^N \Delta^{(n)} \leq \varepsilon.$$

On the other hand, from equation 12, we have

$$\begin{aligned} & \max_{i \in [M]} \left\{ v_i^{\dagger, \hat{\pi}^{-i}} - v_i^{\hat{\pi}} \right\} \left(\text{or } \max_{i \in [M]} \left\{ \max_{\omega \in \Omega_i} v_i^{\omega \circ \hat{\pi}} - v_i^{\hat{\pi}} \right\} \right) \\ &= \max_{i \in [M]} \left\{ v_i^{\dagger, \pi^{(n^*)}^{-i}} - v_i^{\pi^{(n^*)}} \right\} \left(\text{or } \max_{i \in [M]} \left\{ \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi^{(n^*)}} - v_i^{\pi^{(n^*)}} \right\} \right) \\ &\leq \Delta^{(n^*)} = \min_{n \in [N]} \Delta^{(n)} \leq \frac{1}{N} \sum_{n=1}^N \Delta^{(n)} \leq \varepsilon, \end{aligned}$$

which has finished the proof. \square

C ANALYSIS OF THE MODEL-FREE METHOD

For the model-free method, throughout this section we assume the Markov game is a block Markov game.

C.1 CONSTRUCTION OF \mathcal{N}_h AND \mathcal{F}_h

Let $\mathcal{C}_h = \{\Sigma_h : \Sigma_h = \lambda I_d + \sum_{k=1}^l \phi_h(s_k, \mathbf{a}_k) \phi_h(s_k, \mathbf{a}_k)^\top | \phi_h \in \Phi_h, l \in [L], s_k \in \mathcal{S}, \mathbf{a}_k \in \mathcal{A}, \forall k \in [l]\}$. Fix a variable L , for each $h \in [H]$, define a function class $\tilde{\mathcal{F}}_h \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ by

$$\begin{aligned} \tilde{\mathcal{F}}_h = \{ & f(s, \mathbf{a}) := r_{h,i}(s, \mathbf{a}) + \phi_h(s, \mathbf{a})^\top \theta + \min \left(c \|\phi_h(s, \mathbf{a})\|_{\Sigma_h^{-1}}, H \right) \\ & i \in [M], \phi_h \in \Phi_h, \|\theta\|_2 \leq 2H^2 \sqrt{d}, c \in [0, L], \Sigma_h \in \mathcal{C}_h \} \end{aligned}$$

For a given parameter $\tilde{\varepsilon}$, let \mathcal{N}_h be a $\tilde{\varepsilon}$ -net of $\tilde{\mathcal{F}}_h$ under the $\|\cdot\|_\infty$ metric. Define Π_h as the set of all possible policies produced by equation 2 (or equation 3 or equation 4, according to the problem setting). We then define the discriminator function class \mathcal{F}_h as followings:

$$\begin{aligned} \mathcal{F}_{1,h} &:= \left\{ f(s) := \mathbb{E}_{\mathbf{a} \sim U(\mathcal{A})} \left[\|\phi_h(s, \mathbf{a})^\top \theta - \phi'_h(s, \mathbf{a})^\top \theta'\| \mid \phi_h, \phi'_h \in \Phi_h, \max\{\|\theta\|_2, \|\theta'\|_2\} \leq \sqrt{d} \right], \right. \\ \mathcal{F}_{2,h} &:= \left\{ f(s) := \mathbb{E}_{\mathbf{a} \sim \pi_{h+1}(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \theta \right] \mid i \in [M], \pi_{h+1} \in \Pi_{h+1}, \phi_{h+1} \in \Phi_{h+1}, \|\theta\|_2 \leq \sqrt{d} \right\}, \\ \mathcal{F}_{3,h} &:= \left\{ f(s) := \max_{\tilde{\mu}_{h+1,i}} \mathbb{E}_{\mathbf{a} \sim (\tilde{\mu}_{h+1,i} \times \pi_{h+1,-i})(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \theta \right] \mid \right. \\ & \quad \left. i \in [M], \pi_{h+1} \in \Pi_{h+1}, \phi_{h+1} \in \Phi_{h+1}, \|\theta\|_2 \leq \sqrt{d} \right\}, \quad (\text{For NE and CCE}) \\ \mathcal{F}_{3,h} &:= \left\{ f(s) := \max_{\omega_{h+1,i} \in \Omega_{h+1,i}} \mathbb{E}_{\mathbf{a} \sim (\omega_{h+1,i} \circ \pi_{h+1})(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \theta \right] \mid \right. \\ & \quad \left. i \in [M], \pi_{h+1} \in \Pi_{h+1}, \phi_{h+1} \in \Phi_{h+1}, \|\theta\|_2 \leq \sqrt{d} \right\}, \quad (\text{For CE}) \\ \mathcal{F}_{4,h} &:= \left\{ f(s) := \mathbb{E}_{\mathbf{a} \sim \pi_{h+1}(s)} \left[\frac{\min \left\{ c \|\phi_{h+1}(s, \mathbf{a})\|_{\Sigma_{h+1}^{-1}}, H \right\}}{H^2} + \phi_{h+1}(s, \mathbf{a})^\top \theta \right] \mid \right. \\ & \quad \left. c \in [0, L], \pi_{h+1} \in \Pi_{h+1}, \Sigma_{h+1} \in \mathcal{C}_{h+1}, \phi_{h+1} \in \Phi_{h+1}, \|\theta\|_2 \leq \sqrt{d} \right\}, \end{aligned}$$

$$\begin{aligned} \mathcal{G} &:= \{f : \mathcal{S} \rightarrow [0, 1]\}, \\ \mathcal{F}_h &:= (\mathcal{F}_{1,h} \cup \mathcal{F}_{2,h} \cup \mathcal{F}_{3,h} \cup \mathcal{F}_{4,h}) \cap \mathcal{G}. \end{aligned}$$

C.2 HIGH PROBABILITY EVENTS

We define the following event

$$\begin{aligned}\mathcal{E}_1 &: \forall n \in [N], h \in [H], \rho \in \left\{ \rho_h^{(n)}, \tilde{\rho}_h^{(n)} \right\}, f \in \mathcal{F}_h, \quad \mathbb{E}_\rho \left[\left(\left(\hat{P}_h^{(n)} - P_h^* \right) f \right) (s, \mathbf{a}) \right]^2 \leq \zeta^{(n)}, \\ \mathcal{E}_2 &: \forall n \in [N], h \in [H], \phi_h \in \Phi_h, \quad \|\phi_h(s, \mathbf{a})\|_{\left(\hat{\Sigma}_{h, \phi_h}^{(n)} \right)^{-1}} = \Theta \left(\|\phi_h(s, \mathbf{a})\|_{\Sigma_{n, \rho_h^{(n)}, \phi_h}^{-1}} \right) \\ \mathcal{E} &:= \mathcal{E}_1 \cap \mathcal{E}_2.\end{aligned}$$

Similar to the procedure of the model-based case, we first prove a few lemmas which lead to the conclusion that \mathcal{E} holds with a high probability.

Lemma C.1. *For any $n \in [N], h \in [H]$, we have $\hat{P}_h^{(n)}(s'|s, \mathbf{a}) = \hat{\phi}_h^{(n)}(s, \mathbf{a})^\top \hat{w}_h^{(n)}(s')$ for some $\hat{w}_h^{(n)} : \mathcal{S} \rightarrow \mathbb{R}^d$. For any function $f : \mathcal{S} \rightarrow [0, 1]$ and $n \in [N], h \in [H]$, we have $\left\| \int_{\mathcal{S}} \hat{w}_h^{(n)}(s') f(s') ds' \right\|_2 \leq \sqrt{d}$, and there exist $\theta, \tilde{\theta} \in \mathbb{R}^d$ such that $(P_h^* f)(s, \mathbf{a}) = \phi_h^*(s, \mathbf{a})^\top \theta$, $(\hat{P}_h^{(n)} f)(s, \mathbf{a}) = \hat{\phi}_h^{(n)}(s, \mathbf{a})^\top \tilde{\theta}$ and $\max\{\|\theta\|_2, \|\tilde{\theta}\|_2\} \leq \sqrt{d}$. Furthermore, we have $\|\tilde{\theta}\|_\infty \leq 1$.*

Proof. By definition, we have

$$\begin{aligned}(P_h^* f)(s, \mathbf{a}) &= \int_{\mathcal{S}} P_h^*(s'|s, \mathbf{a}) f(s') ds' \\ &= \phi_h^*(s, \mathbf{a})^\top \int_{\mathcal{S}} w_h^*(s') f(s') ds' \\ &= \phi_h^*(s, \mathbf{a})^\top \theta,\end{aligned}$$

where $\theta = \int_{\mathcal{S}} w_h^*(s') f(s') ds'$. Furthermore, note that $\|f\|_\infty \leq 1$, according to the assumption on w_h^* , we have

$$\left\| \int_{\mathcal{S}} w_h^*(s') f(s') ds' \right\|_2 \leq \sqrt{d},$$

which implies $\|\theta\|_2 \leq \sqrt{d}$. For $(\hat{P}_h^{(n)} f)(s, \mathbf{a})$, let

$$\hat{w}_h^{(n)}(s') := \left(\sum_{(\tilde{s}, \tilde{\mathbf{a}}) \in \mathcal{D}_h^{(n)} \cup \tilde{\mathcal{D}}_h^{(n)}} \phi_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \phi_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}})^\top + \lambda I_d \right)^{-1} \sum_{(\tilde{s}, \tilde{\mathbf{a}}, \tilde{s}') \in \mathcal{D}_h^{(n)} \cup \tilde{\mathcal{D}}_h^{(n)}} \phi_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \mathbf{1}_{\tilde{s}'=s'}.$$

Since $\phi_h^{(n)}(s, \mathbf{a})$ is an one-hot vector, one has $\|\hat{w}_h^{(n)}(s')\|_\infty \leq 1, \forall s' \in \mathcal{S}$. It follows that $\left\| \int_{\mathcal{S}} \hat{w}_h^{(n)}(s') f(s') ds' \right\|_\infty \leq 1$, and therefore, $\left\| \int_{\mathcal{S}} \hat{w}_h^{(n)}(s') f(s') ds' \right\|_2 \leq \sqrt{d}$. By definition, we have

$$\begin{aligned}(\hat{P}_h^{(n)} f)(s, \mathbf{a}) &= \int_{\mathcal{S}} \hat{P}_h^{(n)}(s'|s, \mathbf{a}) f(s') ds' \\ &= \hat{\phi}_h^{(n)}(s, \mathbf{a})^\top \int_{\mathcal{S}} \hat{w}_h^{(n)}(s') f(s') ds' \\ &= \hat{\phi}_h^{(n)}(s, \mathbf{a})^\top \tilde{\theta},\end{aligned}$$

where $\tilde{\theta} = \int_{\mathcal{S}} \hat{w}_h^{(n)}(s') f(s') ds'$. Due to the property we just derived for $\hat{w}_h^{(n)}$, similar to the proof of the true model, we also have $\|\tilde{\theta}\|_2 \leq \sqrt{d}$. Meanwhile, one can easily see that $\|\tilde{\theta}\|_\infty \leq 1$, using the fact $\left\| \int_{\mathcal{S}} \hat{w}_h^{(n)}(s') f(s') ds' \right\|_\infty \leq 1$. \square

Lemma C.2 (Covering Number of $\tilde{\mathcal{F}}_h$). *When Φ_h is the set of one-hot vectors and $\lambda \geq 1$, it's possible to construct the $\tilde{\varepsilon}$ -net \mathcal{N}_h such that $|\mathcal{N}_h| \leq M \left(\frac{12H^2L^2d}{\tilde{\varepsilon}} \right)^{3d} |\Phi|, \forall h \in [H]$. Furthermore, we have $|\Pi_h| \leq |\mathcal{N}_h|^M \leq M^M \left(\frac{12H^2L^2d}{\tilde{\varepsilon}} \right)^{3Md} |\Phi|^M$.*

Proof. Recall that

$$\begin{aligned} \tilde{\mathcal{F}}_h &= \left\{ f(s, \mathbf{a}) := r_{h,i}(s, \mathbf{a}) + \phi_h(s, \mathbf{a})^\top \theta + \min\{c\|\phi_h(s, \mathbf{a})\|_{\Sigma_h^{-1}}, H\} \right. \\ &\quad \left. i \in [M], \phi_h \in \Phi_h, \|\theta\|_2 \leq 2H^2\sqrt{d}, c \in [0, L], \Sigma \in \mathcal{C}_h \right\}. \end{aligned}$$

Note that when Φ_h is the set of one-hot vectors, Σ_h will be a diagonal matrix. In this case, $\tilde{\mathcal{F}}_h$ is the subset of the following function class:

$$\begin{aligned} \tilde{\mathcal{F}}'_h &:= \left\{ f(s, \mathbf{a}) := r_{h,i}(s, \mathbf{a}) + \min\{c\phi_h(s, \mathbf{a})^\top \theta', H\} + \phi_h(s, \mathbf{a})^\top \theta \right. \\ &\quad \left. i \in [M], \phi_h \in \Phi_h, 0 \leq c \leq L, \max\{\|\theta\|_2, \|\theta'\|_2\} \leq 2H^2\sqrt{d} \right\}. \end{aligned}$$

Let Θ be an ℓ_2 -cover of the set $\{\theta \in \mathbb{R}^d : \|\theta\|_2 \leq 2H^2\sqrt{d}\}$ at scale $\tilde{\varepsilon}$. Then we know $|\Theta| \leq \left(\frac{4H^2\sqrt{d}}{\tilde{\varepsilon}} \right)^d$. Let \mathcal{W} be an ℓ_∞ -cover of the set $[0, L]$ at scale $\tilde{\varepsilon}' := \frac{\tilde{\varepsilon}}{2H^2\sqrt{d}}$, we have $|\mathcal{W}| \leq \frac{2H^2L\sqrt{d}}{\tilde{\varepsilon}}$. Define the covering set by

$$\bar{\mathcal{F}}_h := \left\{ \bar{f}(s, \mathbf{a}) := r_{h,i}(s, \mathbf{a}) + \min\{\tilde{c}\phi_h(s, \mathbf{a})^\top \tilde{\theta}', H\} + \phi_h(s, \mathbf{a})^\top \tilde{\theta} \mid i \in [M], \phi_h \in \Phi_h, \tilde{c} \in \mathcal{W}, \tilde{\theta}, \tilde{\theta}' \in \Theta \right\}.$$

Then, for any $f \in \tilde{\mathcal{F}}_h$, by definition, suppose f takes the following form:

$$f(s, \mathbf{a}) := r_{h,i}(s, \mathbf{a}) + \min\{c\phi_h(s, \mathbf{a})^\top \theta', H\} + \phi_h(s, \mathbf{a})^\top \theta, \quad 0 \leq c \leq L, \max\{\|\theta\|_2, \|\theta'\|_2\} \leq 2H^2\sqrt{d}.$$

Then we can find $\tilde{\theta}, \tilde{\theta}' \in \Theta, \tilde{c} \in \mathcal{W}$ such that $\|\theta - \tilde{\theta}\|_2 \leq \tilde{\varepsilon}, \|\theta' - \tilde{\theta}'\|_2 \leq \tilde{\varepsilon}$ and $|c - \tilde{c}| \leq \tilde{\varepsilon}'$. Let

$$\bar{f}(s, \mathbf{a}) := r_{h,i}(s, \mathbf{a}) + \min\{\tilde{c}\phi_h(s, \mathbf{a})^\top \tilde{\theta}', H\} + \phi_h(s, \mathbf{a})^\top \tilde{\theta},$$

then we have

$$\begin{aligned} &|f(s, \mathbf{a}) - \bar{f}(s, \mathbf{a})| \\ &\leq \|\phi_h(s, \mathbf{a})\|_2 \|\theta - \tilde{\theta}\|_2 + \|\phi_h(s, \mathbf{a})\|_2 \|\tilde{c}\tilde{\theta}' - c\theta'\|_2 \\ &\leq \tilde{\varepsilon} + |\tilde{c} - c| \|\tilde{\theta}'\|_2 + c \|\theta' - \tilde{\theta}'\|_2 \\ &\leq \tilde{\varepsilon} + 2H^2\sqrt{d}\tilde{\varepsilon}' + L\tilde{\varepsilon} \\ &\leq 3L\tilde{\varepsilon}, \end{aligned}$$

which implies $\bar{\mathcal{F}}_h$ is a $3L\tilde{\varepsilon}$ -covering of $\tilde{\mathcal{F}}'_h$ (therefore, is a $3L\tilde{\varepsilon}$ -covering of $\tilde{\mathcal{F}}_h$), and we have

$$|\bar{\mathcal{F}}_h| \leq M \left(\frac{4H^2Ld}{\tilde{\varepsilon}} \right)^{3d} |\Phi|.$$

Replacing $\tilde{\varepsilon}$ by $\frac{\tilde{\varepsilon}}{3L}$, we get an $\tilde{\varepsilon}$ -covering of $\tilde{\mathcal{F}}_h$ whose size is no larger than $M \left(\frac{12H^2L^2d}{\tilde{\varepsilon}} \right)^{3d} |\Phi|$. For Π_h , since each policy is determined by M members from \mathcal{N}_h , we have $|\Pi_h| \leq |\mathcal{N}_h|^M$, which has finished the proof. \square

Lemma C.3 (Covering Number of \mathcal{F}_h). *When Φ_h is the set of one-hot vectors and $\lambda \geq 1$. The γ -covering number of \mathcal{F}_h is at most $4M|\Pi_{h+1}| \left(\frac{6L^2d}{\gamma} \right)^{3d} |\Phi|^2$.*

Proof. We cover $\mathcal{F}_{1,h}, \mathcal{F}_{2,h}, \mathcal{F}_{3,h}, \mathcal{F}_{4,h}$ separately. For $\mathcal{F}_{1,h}$, let Θ be an ℓ_2 -cover of the set $\{\theta \in \mathbb{R}^d : \|\theta\|_2 \leq \sqrt{d}\}$ at scale γ . Then we know $|\Theta| \leq \left(\frac{2\sqrt{d}}{\gamma} \right)^d$. Define the covering set of $\mathcal{F}_{1,h}$ as

$$\tilde{\mathcal{F}}_{1,h} := \left\{ \tilde{f}(s) := \mathbb{E}_{\mathbf{a} \sim U(\mathcal{A})} \left[\left| \phi_h(s, \mathbf{a})^\top \tilde{\theta} - \phi'_h(s, \mathbf{a})^\top \tilde{\theta}' \right| \right] \mid \phi_h, \phi'_h \in \Phi_h, \tilde{\theta}, \tilde{\theta}' \in \Theta \right\}.$$

For any $f \in \mathcal{F}_{1,h}$, suppose

$$f(s) = \mathbb{E}_{\mathbf{a} \sim U(\mathcal{A})} [\phi_h(s, \mathbf{a})^\top \theta - \phi'_h(s, \mathbf{a})^\top \theta'], \quad \phi_h, \phi'_h \in \Phi_h, \max\{\|\theta\|_2, \|\theta'\|_2\} \leq \sqrt{d},$$

Then we can find $\tilde{\theta}, \tilde{\theta}' \in \Theta$ such that $\|\theta - \tilde{\theta}\|_2 \leq \gamma, \|\theta' - \tilde{\theta}'\|_2 \leq \gamma$. Let

$$\tilde{f}(s) := \mathbb{E}_{\mathbf{a} \sim U(\mathcal{A})} [\phi_h(s, \mathbf{a})^\top \tilde{\theta} - \phi'_h(s, \mathbf{a})^\top \tilde{\theta}'].$$

Then we have

$$\begin{aligned} |f(s) - \tilde{f}(s)| &\leq \frac{1}{A} \sum_{\mathbf{a} \in \mathcal{A}} \|\phi_h(s, \mathbf{a})\|_2 \|\theta - \tilde{\theta}\|_2 + \frac{1}{A} \sum_{\mathbf{a} \in \mathcal{A}} \|\phi'_h(s, \mathbf{a})\|_2 \|\theta' - \tilde{\theta}'\|_2 \\ &\leq 2\gamma, \end{aligned}$$

which implies $\tilde{\mathcal{F}}_{1,h}$ is a 2γ covering of $\mathcal{F}_{1,h}$. Furthermore, we have

$$|\tilde{\mathcal{F}}_{1,h}| \leq \left(\frac{2d}{\gamma}\right)^{2d} |\Phi|^2.$$

For $\mathcal{F}_{2,h}$ and $\mathcal{F}_{3,h}$, we construct

$$\tilde{\mathcal{F}}_{2,h} := \left\{ \tilde{f}(s) := \mathbb{E}_{\mathbf{a} \sim \pi_{h+1}(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \tilde{\theta} \right] \mid i \in [M], \phi_{h+1} \in \Phi_{h+1}, \tilde{\theta} \in \Theta, \pi_{h+1} \in \Pi_{h+1} \right\}.$$

Similar to the proof of $\mathcal{F}_{1,h}$, we may verify $\tilde{\mathcal{F}}_{2,h}$ is a γ -covering of $\mathcal{F}_{2,h}$, and

$$|\tilde{\mathcal{F}}_{2,h}| \leq M |\Pi_{h+1}| \left(\frac{2d}{\gamma}\right)^d |\Phi|.$$

For $\mathcal{F}_{3,h}$, we only prove the case of NE or CCE, the case of CE can be proved in a similar manner. We construct

$$\begin{aligned} \tilde{\mathcal{F}}_{3,h} := &\left\{ \tilde{f}(s) := \max_{\tilde{\mu}_{h+1,i}} \mathbb{E}_{\mathbf{a} \sim (\tilde{\mu}_{h+1,i} \times \pi_{h+1,-i})(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \tilde{\theta} \right] \mid \right. \\ &\left. i \in [M], \phi_{h+1} \in \Phi_{h+1}, \tilde{\theta} \in \Theta, \pi_{h+1} \in \Pi_{h+1} \right\}. \end{aligned}$$

For any $f \in \mathcal{F}_{3,h}$, suppose

$$f(s) = \max_{\tilde{\mu}_{h+1,i}} \mathbb{E}_{\mathbf{a} \sim (\tilde{\mu}_{h+1,i} \times \pi_{h+1,-i})(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \theta \right], \quad i \in [M], \pi_{h+1} \in \Pi_{h+1}, \phi_{h+1} \in \Phi_{h+1}, \|\theta\|_2 \leq \sqrt{d}.$$

Then we can find $\tilde{\theta} \in \Theta$ such that $\|\theta - \tilde{\theta}\|_2 \leq \gamma$. Let

$$\tilde{f}(s) = \max_{\tilde{\mu}_{h+1,i}} \mathbb{E}_{\mathbf{a} \sim (\tilde{\mu}_{h+1,i} \times \pi_{h+1,-i})(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \tilde{\theta} \right],$$

we have

$$\begin{aligned} f(s) - \tilde{f}(s) &= \max_{\tilde{\mu}_{h+1,i}} \mathbb{E}_{\mathbf{a} \sim (\tilde{\mu}_{h+1,i} \times \pi_{h+1,-i})(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \theta \right] \\ &\quad - \max_{\tilde{\mu}_{h+1,i}} \mathbb{E}_{\mathbf{a} \sim (\tilde{\mu}_{h+1,i} \times \pi_{h+1,-i})(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \tilde{\theta} \right] \\ &\leq \max_{\tilde{\mu}_{h+1,i}} \left(\mathbb{E}_{\mathbf{a} \sim (\tilde{\mu}_{h+1,i} \times \pi_{h+1,-i})(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \theta \right] \right. \\ &\quad \left. - \mathbb{E}_{\mathbf{a} \sim (\tilde{\mu}_{h+1,i} \times \pi_{h+1,-i})(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \tilde{\theta} \right] \right) \\ &= \max_{\tilde{\mu}_{h+1,i}} \left(\mathbb{E}_{\mathbf{a} \sim (\tilde{\mu}_{h+1,i} \times \pi_{h+1,-i})(s)} \left[\phi_{h+1}(s, \mathbf{a})^\top \theta - \phi_{h+1}(s, \mathbf{a})^\top \tilde{\theta} \right] \right) \\ &\leq \|\theta - \tilde{\theta}\|_2 \end{aligned}$$

$$\leq \gamma,$$

and

$$\begin{aligned} \tilde{f}(s) - f(s) &= \max_{\tilde{\mu}_{h+1,i}} \mathbb{E}_{\mathbf{a} \sim (\tilde{\mu}_{h+1,i} \times \pi_{h+1,-i})(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \tilde{\theta} \right] \\ &\quad - \max_{\tilde{\mu}_{h+1,i}} \mathbb{E}_{\mathbf{a} \sim (\tilde{\mu}_{h+1,i} \times \pi_{h+1,-i})(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \theta \right] \\ &\leq \max_{\tilde{\mu}_{h+1,i}} \left(\mathbb{E}_{\mathbf{a} \sim (\tilde{\mu}_{h+1,i} \times \pi_{h+1,-i})(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \tilde{\theta} \right] \right. \\ &\quad \left. - \mathbb{E}_{\mathbf{a} \sim (\tilde{\mu}_{h+1,i} \times \pi_{h+1,-i})(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \phi_{h+1}(s, \mathbf{a})^\top \theta \right] \right) \\ &= \max_{\tilde{\mu}_{h+1,i}} \left(\mathbb{E}_{\mathbf{a} \sim (\tilde{\mu}_{h+1,i} \times \pi_{h+1,-i})(s)} \left[\phi_{h+1}(s, \mathbf{a})^\top \tilde{\theta} - \phi_{h+1}(s, \mathbf{a})^\top \theta \right] \right) \\ &\leq \|\theta - \tilde{\theta}\|_2 \\ &\leq \gamma, \end{aligned}$$

which implies

$$|\tilde{f}(s) - f(s)| \leq \gamma.$$

Therefore, we conclude $\tilde{\mathcal{F}}_{3,h}$ is a γ -covering of $\mathcal{F}_{3,h}$, and

$$|\tilde{\mathcal{F}}_{3,h}| \leq M |\Pi_{h+1}| \left(\frac{2d}{\gamma} \right)^d |\Phi|.$$

For $\mathcal{F}_{4,h}$, note that when Φ_h is the set of one-hot vectors, Σ_h will be a diagonal matrix. In this case, $\mathcal{F}_{4,h}$ is the subset of the following function class:

$$\begin{aligned} \mathcal{F}'_{4,h} := &\left\{ f(s) := \mathbb{E}_{\mathbf{a} \sim \pi_{h+1}(s)} \left[\frac{\min\{c\phi_{h+1}(s, \mathbf{a})^\top \theta', H\}}{H^2} + \phi_{h+1}(s, \mathbf{a})^\top \theta \right] \right. \\ &\left. 0 \leq c \leq L, \pi_{h+1} \in \Pi_{h+1}, \max\{\|\theta\|_2, \|\theta'\|_2\} \leq \sqrt{d}, \phi_{h+1} \in \Phi_{h+1} \right\}. \end{aligned}$$

In this case, let \mathcal{W} be an ℓ_∞ cover of the set $[0, L]$ at scale $\tilde{\gamma} := \frac{\gamma}{\sqrt{d}}$, we have $|\mathcal{W}| \leq \frac{L\sqrt{d}}{\tilde{\gamma}}$. Let

$$\begin{aligned} \tilde{\mathcal{F}}_{4,h} := &\left\{ \tilde{f}(s) := \mathbb{E}_{\mathbf{a} \sim \pi_{h+1}(s)} \left[\frac{\min\{\tilde{c}\phi_{h+1}(s, \mathbf{a})^\top \tilde{\theta}', H\}}{H^2} + \phi_{h+1}(s, \mathbf{a})^\top \tilde{\theta} \right] \right. \\ &\left. \tilde{c} \in \mathcal{W}, \pi_{h+1} \in \Pi_{h+1}, \tilde{\theta}, \tilde{\theta}' \in \Theta, \phi_{h+1} \in \Phi_{h+1} \right\}. \end{aligned}$$

Then, for any $f \in \mathcal{F}_{4,h}$, suppose

$$\begin{aligned} f(s) &:= \mathbb{E}_{\mathbf{a} \sim \pi_{h+1}(s)} \left[\frac{\min\{c\phi_{h+1}(s, \mathbf{a})^\top \theta', H\}}{H^2} + \phi_{h+1}(s, \mathbf{a})^\top \theta \right], \\ &0 \leq c \leq L, \pi_{h+1} \in \Pi_{h+1}, \max\{\|\theta\|_2, \|\theta'\|_2\} \leq \sqrt{d}, \phi_{h+1} \in \Phi_{h+1}. \end{aligned}$$

Then we can find $\tilde{\theta}, \tilde{\theta}' \in \Theta, \tilde{c} \in \mathcal{W}$ such that $\|\theta - \tilde{\theta}\|_2 \leq \gamma, \|\theta' - \tilde{\theta}'\|_2 \leq \gamma$ and $|c - \tilde{c}| \leq \tilde{\gamma}$. Let

$$\tilde{f}(s) := \mathbb{E}_{\mathbf{a} \sim \pi_{h+1}(s)} \left[\frac{\min\{\tilde{c}\phi_{h+1}(s, \mathbf{a})^\top \tilde{\theta}', H\}}{H^2} + \phi_{h+1}(s, \mathbf{a})^\top \tilde{\theta} \right],$$

then we have

$$\begin{aligned} &|f(s) - \tilde{f}(s)| \\ &\leq \mathbb{E}_{\mathbf{a} \sim \pi_{h+1}(s)} \left[\|\phi_{h+1}(s, \mathbf{a})\|_2 \|\theta - \tilde{\theta}\|_2 \right] + \frac{1}{H^2} \mathbb{E}_{\mathbf{a} \sim \pi_{h+1}(s)} \left[\|\phi_{h+1}(s, \mathbf{a})\|_2 \|\tilde{c}\tilde{\theta}' - c\theta'\|_2 \right] \end{aligned}$$

$$\begin{aligned}
&\leq \gamma + \frac{1}{H^2} \left(|\tilde{c} - c| \|\tilde{\theta}'\|_2 + c \|\theta' - \tilde{\theta}'\|_2 \right) \\
&\leq \gamma + \frac{\sqrt{d}}{H^2} \tilde{\gamma} + \frac{L}{H^2} \gamma \\
&\leq 3L\gamma,
\end{aligned}$$

which implies $\tilde{\mathcal{F}}_{4,h}$ is a $3L\gamma$ -covering of $\mathcal{F}_{4,h}$, and we have

$$|\tilde{\mathcal{F}}_{4,h}| \leq |\Pi_{h+1}| \left(\frac{2Ld}{\gamma} \right)^{3d} |\Phi|.$$

In summary, we know $\tilde{\mathcal{F}}_h := \tilde{\mathcal{F}}_{1,h} \cup \tilde{\mathcal{F}}_{2,h} \cup \tilde{\mathcal{F}}_{3,h} \cup \tilde{\mathcal{F}}_{4,h}$ is a $3L\gamma$ -covering of \mathcal{F}_h . And

$$|\mathcal{F}_h| \leq 4M |\Pi_{h+1}| \left(\frac{2Ld}{\gamma} \right)^{3d} |\Phi|^2.$$

Replacing γ by $\frac{\gamma}{3L}$, we get an γ -covering of \mathcal{F}_h whose size is no larger than $4M |\Pi_{h+1}| \left(\frac{6L^2d}{\gamma} \right)^{3d} |\Phi|^2$, which has finished the proof. \square

Below we omit the superscript n and subscript h when clear from the context. Denote

$$\mathcal{L}_{\lambda, \mathcal{D}}(\phi, \theta, f) = \frac{1}{|\mathcal{D}|} \sum_{(s, \mathbf{a}, s') \in \mathcal{D}} (\phi(s, \mathbf{a})^\top \theta - f(s'))^2 + \frac{\lambda}{|\mathcal{D}|} \|\theta\|_2^2 \quad (13)$$

$$\mathcal{L}_{\mathcal{D}}(\phi, \theta, f) = \frac{1}{|\mathcal{D}|} \sum_{(s, \mathbf{a}, s') \in \mathcal{D}} (\phi(s, \mathbf{a})^\top \theta - f(s'))^2 \quad (14)$$

$$\mathcal{L}_{\rho}(\phi, \theta, f) = \mathbb{E}_{(s, \mathbf{a}) \sim \rho, s' \sim P^*(s, \mathbf{a})} [(\phi(s, \mathbf{a})^\top \theta - f(s'))^2]. \quad (15)$$

Lemma C.4 (Uniform Convergence for Square Loss). *Let there be a dataset $\mathcal{D} := \{(s_i, \mathbf{a}_i, s'_i)\}_{i=1}^n$ collected in n episodes. Denote that the data generating distribution in iteration i by d_i , and $\rho = \frac{1}{n} \sum_{i=1}^n d_i$. Note that d_i can depend on the randomness in episodes $1, \dots, i-1$. For a finite feature class Φ and a discriminator class $\mathcal{F} : \mathcal{S} \rightarrow [0, 1]$ with γ -covering number $\|\mathcal{F}\|_{\gamma}$, we will show that, with probability at least $1 - \delta$:*

$$\begin{aligned}
&|[\mathcal{L}_{\rho}(\phi, \theta, f) - \mathcal{L}_{\rho}(\phi^*, \theta_f^*, f)] - [\mathcal{L}_{\mathcal{D}}(\phi, \theta, f) - \mathcal{L}_{\mathcal{D}}(\phi^*, \theta_f^*, f)]| \\
&\leq \frac{1}{2} [\mathcal{L}_{\rho}(\phi, \theta, f) - \mathcal{L}_{\rho}(\phi^*, \theta_f^*, f)] + \frac{64 \log\left(\frac{2(4n)^d \cdot |\Phi| \cdot \|\mathcal{F}\|_{1/2n}}{\delta}\right)}{n}
\end{aligned}$$

for all $\phi \in \Phi$, $\|\theta\|_{\infty} \leq 1$ and $f \in \mathcal{F}$, where recall that ϕ^* is the true feature and θ_f^* is defined as $\mathbb{E}_{s' \sim P^*(s, \mathbf{a})}[f(s')] = \langle \phi^*(s, \mathbf{a}), \theta_f^* \rangle$.

Proof. To start, we focus on a given $f \in \mathcal{F}$. We first give a high probability bound on the following deviation term:

$$|\mathcal{L}_{\rho}(\phi, \theta, f) - \mathcal{L}_{\rho}(\phi^*, \theta_f^*, f) - (\mathcal{L}_{\mathcal{D}}(\phi, \theta, f) - \mathcal{L}_{\mathcal{D}}(\phi^*, \theta_f^*, f))|.$$

Denote $g(s_i, \mathbf{a}_i) = \phi(s_i, \mathbf{a}_i)^\top \theta$ and $g^*(s_i, \mathbf{a}_i) = \phi^*(s_i, \mathbf{a}_i)^\top \theta_f^*$. At episode i , let \mathcal{F}_{i-1} be the σ -field generated by all the random variables over the first $i-1$ episodes, for the random variable $Y_i := (g(s_i, \mathbf{a}_i) - f(s'_i))^2 - (g^*(s_i, \mathbf{a}_i) - f(s'_i))^2$, we have

$$\begin{aligned}
\mathbb{E}[Y_i | \mathcal{F}_{i-1}] &= \mathbb{E} \left[(g(s_i, \mathbf{a}_i) - f(s'_i))^2 - (g^*(s_i, \mathbf{a}_i) - f(s'_i))^2 \right] \\
&= \mathbb{E} [(g(s_i, \mathbf{a}_i) + g^*(s_i, \mathbf{a}_i) - 2f(s'_i)) (g(s_i, \mathbf{a}_i) - g^*(s_i, \mathbf{a}_i))] \\
&= \mathbb{E} [(g(s_i, \mathbf{a}_i) - g^*(s_i, \mathbf{a}_i))^2].
\end{aligned}$$

Here the conditional expectation is taken according to the distribution $d_i | \mathcal{F}_{i-1}$. The last equality is due to the fact that

$$\mathbb{E} [(g^*(s_i, \mathbf{a}_i) - f(s'_i)) (g(s_i, \mathbf{a}_i) - g^*(s_i, \mathbf{a}_i))]$$

$$\begin{aligned}
&= \mathbb{E}_{s_i, \mathbf{a}_i} \left[\mathbb{E}_{s'_i} \left[(g^*(s_i, \mathbf{a}_i) - f(s'_i)) (g(s_i, \mathbf{a}_i) - g^*(s_i, \mathbf{a}_i)) \mid s_i, \mathbf{a}_i \right] \right] \\
&= 0.
\end{aligned}$$

Next, for the conditional variance of the random variable, we have:

$$\begin{aligned}
\mathbb{V}[Y_i | \mathcal{F}_{i-1}] &\leq \mathbb{E}[Y_i^2 | \mathcal{F}_{i-1}] = \mathbb{E} \left[(g(s_i, \mathbf{a}_i) + g^*(s_i, \mathbf{a}_i) - 2f(s'_i))^2 (g(s_i, \mathbf{a}_i) - g^*(s_i, \mathbf{a}_i))^2 \mid \mathcal{F}_{i-1} \right] \\
&\leq 16 \mathbb{E} \left[(g(s_i, \mathbf{a}_i) - g^*(s_i, \mathbf{a}_i))^2 \mid \mathcal{F}_{i-1} \right] \\
&\leq 16 \mathbb{E}[Y_i | \mathcal{F}_{i-1}].
\end{aligned}$$

Noticing $Y_i \in [-4, 4]$. Applying Lemma 1 in (Foster and Rakhlin, 2020), we get with probability at least $1 - \delta'$, we can bound the deviation term above as:

$$\begin{aligned}
&\left| \mathcal{L}_\rho(\phi, \theta, f) - \mathcal{L}_\rho(\phi^*, \theta_f^*, f) - (\mathcal{L}_\mathcal{D}(\phi, \theta, f) - \mathcal{L}_\mathcal{D}(\phi^*, \theta_f^*, f)) \right| \\
&\leq \sqrt{\frac{2 \sum_{i=1}^n \mathbb{V}[Y_i | \mathcal{F}_{i-1}] \log \frac{2}{\delta'}}{n^2}} + \frac{16 \log \frac{2}{\delta'}}{3n} \\
&\leq \sqrt{\frac{32 \sum_{i=1}^n \mathbb{E}[Y_i | \mathcal{F}_{i-1}] \log \frac{2}{\delta'}}{n^2}} + \frac{16 \log \frac{2}{\delta'}}{3n},
\end{aligned}$$

Further, consider a finite point-wise cover of the function class $\mathcal{G} := \{g(s, \mathbf{a}) = \phi(s, \mathbf{a})^\top \theta : \phi \in \Phi, \|\theta\|_\infty \leq 1\}$. Note that, with a ℓ_∞ -cover $\bar{\mathcal{W}}$ of $\mathcal{W} = \{\|\theta\|_\infty \leq 1\}$ at scale γ , we have for all (s, \mathbf{a}) and $\phi \in \Phi$, there exists $\bar{\theta} \in \bar{\mathcal{W}}$, $|\langle \phi(s, \mathbf{a}), \theta - \bar{\theta} \rangle| \leq \gamma$, and we have $|\mathcal{W}| = \left(\frac{2}{\gamma}\right)^d$. Let $\tilde{\mathcal{F}}$ be a γ -covering set of \mathcal{F} . For any $f \in \mathcal{F}$, there exists $\bar{f} \in \tilde{\mathcal{F}}$ such that $\|f - \bar{f}\|_\infty \leq \gamma$. Then, applying a union bound over elements in $\Phi \times \bar{\mathcal{W}} \times \tilde{\mathcal{F}}$, with probability $1 - |\Phi| |\bar{\mathcal{W}}| |\tilde{\mathcal{F}}| \delta'$, for all $\theta \in \mathcal{W}$, $f \in \mathcal{F}$, we have:

$$\begin{aligned}
&\left| \mathcal{L}_\rho(\phi, \theta, f) - \mathcal{L}_\rho(\phi^*, \theta_f^*, f) - (\mathcal{L}_\mathcal{D}(\phi, \theta, f) - \mathcal{L}_\mathcal{D}(\phi^*, \theta_f^*, f)) \right| \\
&\leq \left| \mathcal{L}_\rho(\phi, \bar{\theta}, \bar{f}) - \mathcal{L}_\rho(\phi^*, \theta_{\bar{f}}^*, \bar{f}) - (\mathcal{L}_\mathcal{D}(\phi, \bar{\theta}, \bar{f}) - \mathcal{L}_\mathcal{D}(\phi^*, \theta_{\bar{f}}^*, \bar{f})) \right| + 16\gamma \\
&\leq \sqrt{\frac{32 \sum_{i=1}^n \mathbb{E}[\bar{Y}_i | \mathcal{F}_{i-1}] \log \frac{2}{\delta'}}{n^2}} + \frac{16 \log \frac{2}{\delta'}}{3n} + 16\gamma \\
&\leq \frac{1}{2n} \sum_{i=1}^n \mathbb{E}[\bar{Y}_i | \mathcal{F}_{i-1}] + \frac{16 \log \frac{2}{\delta'}}{n} + \frac{16 \log \frac{2}{\delta'}}{3n} + 16\gamma \\
&\leq \frac{1}{2n} \sum_{i=1}^n \mathbb{E}[Y_i | \mathcal{F}_{i-1}] + \frac{16 \log \frac{2}{\delta'}}{n} + \frac{16 \log \frac{2}{\delta'}}{3n} + 32\gamma \\
&\leq \frac{1}{2} (\mathcal{L}_\rho(\phi, \theta, f) - \mathcal{L}_\rho(\phi^*, \theta_f^*, f)) + \frac{32 \log \frac{2}{\delta'}}{n} + 32\gamma \\
&\leq \frac{1}{2} (\mathcal{L}_\rho(\phi, \theta, f) - \mathcal{L}_\rho(\phi^*, \theta_f^*, f)) + \frac{64 \log \frac{2}{\delta'}}{n} \quad (\text{setting } \gamma = 1/n)
\end{aligned}$$

where $\bar{Y}_i := (\phi(s_i, \mathbf{a}_i)^\top \bar{\theta} - \bar{f}(s'_i))^2 - (\phi(s_i, \mathbf{a}_i)^\top \theta_{\bar{f}}^* - \bar{f}(s'_i))^2$. Finally, setting $\delta = \delta' / (|\Phi| |\bar{\mathcal{W}}| |\tilde{\mathcal{F}}|)$, we get $\log \frac{2}{\delta'} \leq \log \frac{2(4n)^d |\Phi| |\tilde{\mathcal{F}}|}{\delta}$. This completes the proof. \square

Lemma C.5 (Deviation Bounds for Alg. 3). *Let $\varepsilon' = \frac{128 \log(\frac{2(4n)^d |\Phi| \|\mathcal{F}\|_{1/2n}}{\delta})}{n}$. If Alg. 3 is called with a dataset \mathcal{D} of size n , then with probability at least $1 - \delta$, for any $f \in \mathcal{F} \subset [0, 1]^S$, we have*

$$\mathbb{E}_\rho \left[\left(\hat{\phi}(s, \mathbf{a})^\top \hat{\theta}_f - \phi^*(s, \mathbf{a})^\top \theta_f^* \right)^2 \right] \leq \varepsilon' + \frac{2\lambda d}{n}.$$

Proof. We begin by using the result in Lemma C.4 such that, with probability at least $1 - \delta$, for all $\|\theta\|_\infty \leq 1$, $\phi \in \Phi$ and $f \in \mathcal{F}$, we have

$$\left| [\mathcal{L}_\rho(\phi, \theta, f) - \mathcal{L}_\rho(\phi^*, \theta_f^*, f)] - [\mathcal{L}_\mathcal{D}(\phi, \theta, f) - \mathcal{L}_\mathcal{D}(\phi^*, \theta_f^*, f)] \right| \leq \frac{1}{2} [\mathcal{L}_\rho(\phi, \theta, f) - \mathcal{L}_\rho(\phi^*, \theta_f^*, f)] + \varepsilon'/2.$$

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

$$\begin{aligned}
& \mathbb{E}_\rho \left[\left(\hat{\phi}(s, \mathbf{a})^\top \hat{\theta}_f - \phi^*(s, \mathbf{a})^\top \theta_f^* \right)^2 \right] \\
&= \mathcal{L}_\rho(\hat{\phi}, \hat{\theta}_f, f) - \mathcal{L}_\rho(\phi^*, \theta_f^*, f) \quad (\text{since } \mathbb{E}_{s' \sim P^*(s, \mathbf{a})} [f(s')] = \phi^*(s, \mathbf{a})^\top \theta_f^*) \\
&\leq 2 \left(\mathcal{L}_\mathcal{D}(\hat{\phi}, \hat{\theta}_f, f) - \mathcal{L}_\mathcal{D}(\phi^*, \theta_f^*, f) \right) + \varepsilon' \\
&\quad (\text{Lemma C.4, and } \|\hat{\theta}_f\|_\infty \leq 1 \text{ according to the proof in Lemma C.1}) \\
&\leq 2 \left(\mathcal{L}_{\lambda, \mathcal{D}}(\hat{\phi}, \hat{\theta}_f, f) - \mathcal{L}_{\lambda, \mathcal{D}}(\phi^*, \theta_f^*, f) + \frac{\lambda}{n} \|\theta_f^*\|_2^2 \right) + \varepsilon' \\
&\leq \varepsilon' + \frac{2\lambda d}{n}, \quad (\text{by the optimality of } \hat{\phi}, \hat{\theta}_f \text{ under } \mathcal{L}_{\lambda, \mathcal{D}}(\cdot, \cdot, f))
\end{aligned}$$

which means the inequality in the lemma statement holds. Here, we use $\|\theta_f^*\|_2^2 \leq d$. \square

Lemma C.6. When $\hat{P}_h^{(n)}$ is computed using Alg. 3 and the Markov games is a block Markov game, if we set

$$\lambda = \Theta \left(d \log \frac{NH|\Phi|}{\delta} \right), \quad \zeta^{(n)} = \Theta \left(\frac{d^2 M \log \frac{dNHML|\Phi|}{\delta \tilde{\varepsilon}}}{n} \right).$$

then \mathcal{E} holds with probability at least $1 - \delta$.

Proof. Combining Lemma C.5 and Lemma C.3, we have that

$$\max_{f \in \tilde{\mathcal{F}}_h} \mathbb{E}_\rho \left[\left(\hat{\phi}(s, \mathbf{a})^\top \hat{\theta}_f - \phi^*(s, \mathbf{a})^\top \theta_f^* \right)^2 \right] \leq \varepsilon' + \frac{2\lambda d}{n} \leq \zeta^{(n)} := \Theta \left(d^2 M \frac{\log \left(\frac{dNHML|\Phi|}{\delta \tilde{\varepsilon}} \right)}{n} \right),$$

which shows \mathcal{E}_1 holds with a high probability. Combining this result with Lemma E.1, we have proved Lemma C.6. \square

C.3 STATISTICAL GUARANTEES

To ensure the algorithm is well-defined, we first prove the following lemma which implies the optimistic Q-value estimators always belong to the function class $\tilde{\mathcal{F}}_h$.

Lemma C.7. When $\alpha^{(n)} \leq L$, we have $\bar{Q}_{h,i}^{(n)} \in \tilde{\mathcal{F}}_h, \forall h \in [H], i \in [M], n \in [N]$.

Proof. Because $\hat{\beta}_h^{(n)}$ is upper bounded by H , by induction one can easily get $\bar{V}_{h+1,i}^{(n)} \leq 2H^2$. Then according to the result of Lemma C.1, we know $(\hat{P}_h^{(n)} \bar{V}_{h+1,i}^{(n)})(s, a) = \phi_h^{(n)}(s, \mathbf{a})^\top \theta$ with $\|\theta\|_2 \leq 2H^2 \sqrt{d}$. We conclude $\bar{Q}_{h,i}^{(n)} \in \tilde{\mathcal{F}}_h$. \square

We will show later that our choice of $\alpha^{(n)}$ and L always satisfies the condition $\alpha^{(n)} \leq L$.

Lemma C.8. We have

- For NE and CCE,

$$\max_{\pi_{h,i}} \left(\mathbb{D}_{\pi_{h,i}, \pi_{h,-i}^{(n)}} \bar{Q}_{h,i}^{(n)} \right) (s) \leq \left(\mathbb{D}_{\pi_h^{(n)}} \bar{Q}_{h,i}^{(n)} \right) (s) + 2\tilde{\varepsilon};$$

- For CE,

$$\max_{\omega_{h,i} \in \Omega_{h,i}} \left(\mathbb{D}_{\omega_{h,i}, \pi_h^{(n)}} \bar{Q}_{h,i}^{(n)} \right) (s) \leq \left(\mathbb{D}_{\pi_h^{(n)}} \bar{Q}_{h,i}^{(n)} \right) (s) + 2\tilde{\varepsilon}.$$

Proof. We only prove the case of NE and CCE, the case of CE can be proved similarly. Let $\tilde{Q}_{h,i}^{(n)}$ be the nearest neighbour of $\bar{Q}_{h,i}^{(n)}$ in \mathcal{N}_h , we have

$$\begin{aligned} \max_{\pi_{h,i}} \left(\mathbb{D}_{\pi_{h,i}, \pi_{h,-i}^{(n)}} \bar{Q}_{h,i}^{(n)} \right) (s) &\leq \max_{\pi_{h,i}} \left(\mathbb{D}_{\pi_{h,i}, \pi_{h,-i}^{(n)}} \tilde{Q}_{h,i}^{(n)} \right) (s) + \tilde{\varepsilon} \\ &\leq \left(\mathbb{D}_{\pi_h^{(n)}} \tilde{Q}_{h,i}^{(n)} \right) (s) + \tilde{\varepsilon} && \text{(Definition of } \pi_h^{(n)}) \\ &\leq \left(\mathbb{D}_{\pi_h^{(n)}} \bar{Q}_{h,i}^{(n)} \right) (s) + 2\tilde{\varepsilon}, \end{aligned}$$

which has finished the proof. \square

Lemma C.9 (One-step back inequality for the learned model). *Suppose the event \mathcal{E} holds. Consider a set of functions $\{g_h\}_{h=1}^H$ that satisfies $g_h \in \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}_+$, s.t. $\|g_h\|_\infty \leq B$. For a given policy π , suppose $\mathbb{E}_{\mathbf{a} \sim U(\mathcal{A})} [g_h(\cdot, \mathbf{a})] \in \mathcal{F}_{1,h}$, then we have*

$$\begin{aligned} &\left| \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^\pi} [g_h(s, \mathbf{a})] \right| \\ &\leq \begin{cases} \sqrt{A \mathbb{E}_{(s,\mathbf{a}) \sim \rho_1^{(n)}} [g_1^2(s, \mathbf{a})]}, & h = 1 \\ \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)},h-1}^\pi} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n,\rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}^{-1}} \sqrt{nA^2 \mathbb{E}_{(s,\mathbf{a}) \sim \hat{\rho}_h^{(n)}} [g_h^2(s, \mathbf{a})] + B^2 \lambda d + nA^2 \zeta^{(n)}}, B \right\} \right], & h \geq 2 \end{cases} \end{aligned}$$

Recall $\Sigma_{n,\rho_h^{(n)}, \hat{\phi}_h^{(n)}} = n \mathbb{E}_{(s,\mathbf{a}) \sim \rho_h^{(n)}} \left[\hat{\phi}_h^{(n)}(s, \mathbf{a}) \hat{\phi}_h^{(n)}(s, \mathbf{a})^\top \right] + \lambda I_d$.

Proof. For step $h = 1$, we have

$$\begin{aligned} \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},1}^\pi} [g_1(s, \mathbf{a})] &= \mathbb{E}_{s \sim d_1, \mathbf{a} \sim \pi_1(s)} [g_1(s, \mathbf{a})] \\ &\leq \sqrt{\max_{(s,\mathbf{a})} \frac{d_1(s) \pi_1(\mathbf{a}|s)}{\rho_1^{(n)}(s, \mathbf{a})} \mathbb{E}_{(s',\mathbf{a}') \sim \rho_1^{(n)}} [g_1^2(s', \mathbf{a}')] } \\ &= \sqrt{\max_{(s,\mathbf{a})} \frac{d_1(s) \pi_1(\mathbf{a}|s)}{d_1(s) u_{\mathcal{A}}(\mathbf{a})} \mathbb{E}_{(s',\mathbf{a}') \sim \rho_1^{(n)}} [g_1^2(s', \mathbf{a}')] } \\ &\leq \sqrt{A \mathbb{E}_{(s,\mathbf{a}) \sim \rho_1^{(n)}} [g_1^2(s, \mathbf{a})]}. \end{aligned}$$

For step $h = 2, \dots, H-1$, we observe the following one-step-back decomposition:

$$\begin{aligned} &\mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^\pi} [g_h(s, \mathbf{a})] \\ &= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)},h-1}^\pi, s \sim \hat{P}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim \pi_h(s)} [g_h(s, \mathbf{a})] \\ &= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)},h-1}^\pi} \left[\hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}})^\top \int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right] \\ &= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)},h-1}^\pi} \left[\min \left\{ \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}})^\top \int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds, B \right\} \right] \\ &\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)},h-1}^\pi} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n,\rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}^{-1}} \left\| \int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right\|_{\Sigma_{n,\rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}}, B \right\} \right]. \end{aligned}$$

where we use the fact that g_h is bounded by B . Then,

$$\left\| \int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right\|_{\Sigma_{n,\rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}}^2$$

$$\begin{aligned}
&\leq \left(\int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right)^\top \left(n \mathbb{E}_{(s, \mathbf{a}) \sim \rho_{h-1}^{(n)}} \left[\hat{\phi}_{h-1}^{(n)}(s, \mathbf{a}) \hat{\phi}_{h-1}^{(n)}(s, \mathbf{a})^\top \right] + \lambda I_d \right) \left(\int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right) \\
&\leq n \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \rho_{h-1}^{(n)}} \left[\left(\int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right)^2 \right] + B^2 \lambda d \\
&\quad \left(\left\| \sum_{\mathbf{a} \in \mathcal{A}} \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) \right\|_\infty \leq B \text{ and by Lemma C.1 } \left\| \int_{\mathcal{S}} \hat{w}_{h-1}^{(n)}(s) l(s) ds \right\|_2 \leq \sqrt{d} \text{ for any } l : \mathcal{S} \rightarrow [0, 1] \right) \\
&= n \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \rho_{h-1}^{(n)}} \left[\left(\mathbb{E}_{s \sim \hat{P}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim \pi_h(s)} [g_h(s, \mathbf{a})] \right)^2 \right] + B^2 \lambda d \\
&\leq n A^2 \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \rho_{h-1}^{(n)}} \left[\left(\mathbb{E}_{s \sim \hat{P}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim U(\mathcal{A})} [g_h(s, \mathbf{a})] \right)^2 \right] + B^2 \lambda d \quad (\text{Importance sampling}) \\
&\leq n A^2 \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \rho_{h-1}^{(n)}} \left[\left(\mathbb{E}_{s \sim P_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim U(\mathcal{A})} [g_h(s, \mathbf{a})] \right)^2 \right] + B^2 \lambda d + n A^2 \zeta^{(n)} \quad (\text{Assumption on } g_h) \\
&\leq n A^2 \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \rho_{h-1}^{(n)}, s \sim P_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a} \sim U(\mathcal{A})} [g_h^2(s, \mathbf{a})] + B^2 \lambda d + n A^2 \zeta^{(n)}. \quad (\text{Jensen}) \\
&\leq n A^2 \mathbb{E}_{(s, \mathbf{a}) \sim \hat{\rho}_h^{(n)}} [g_h^2(s, \mathbf{a})] + B^2 \lambda d + n A^2 \zeta^{(n)}. \quad (\text{Definition of } \hat{\rho}_h^{(n)})
\end{aligned}$$

Combing the above results together, we get

$$\begin{aligned}
&\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}_{h-1}, h}^\pi} [g_h(s, \mathbf{a})] \\
&\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}_{h-1}, h-1}^\pi} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}^{-1}} \left\| \int_{\mathcal{S}} \sum_{\mathbf{a} \in \mathcal{A}} \hat{w}_{h-1}^{(n)}(s) \pi_h(\mathbf{a}|s) g_h(s, \mathbf{a}) ds \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}}, B \right\} \right] \\
&\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}_{h-1}, h-1}^\pi} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}^{-1}} \sqrt{n A^2 \mathbb{E}_{(s, \mathbf{a}) \sim \hat{\rho}_h^{(n)}} [g_h^2(s, \mathbf{a})] + B^2 \lambda d + n A^2 \zeta^{(n)}}, B \right\} \right],
\end{aligned}$$

which has finished the proof. \square

The following lemma is an exact copy of Lemma B.4, and here we state it again just for completeness.

Lemma C.10 (One-step back inequality for the true model). *Consider a set of functions $\{g_h\}_{h=1}^H$ that satisfies $g_h \in \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}_+$, s.t. $\|g_h\|_\infty \leq B$. Then for any given policy π , we have*

$$\begin{aligned}
&\left| \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}_{h-1}, h}^\pi} [g_h(s, \mathbf{a})] \right| \\
&\leq \begin{cases} \sqrt{A \mathbb{E}_{(s, \mathbf{a}) \sim \rho_1^{(n)}} [g_1^2(s, \mathbf{a})]}, & h = 1 \\ \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}_{h-1}, h-1}^\pi} \left[\left\| \phi_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_{h-1}^{(n)}, \phi_{h-1}^*}^{-1}} \sqrt{n A \mathbb{E}_{(s, \mathbf{a}) \sim \hat{\rho}_h^{(n)}} [g_h^2(s, \mathbf{a})] + B^2 \lambda d}, \right] & h \geq 2 \end{cases}
\end{aligned}$$

Recall $\Sigma_{n, \gamma_h^{(n)}, \phi_h^*} = n \mathbb{E}_{(s, \mathbf{a}) \sim \gamma_h^{(n)}} [\phi_h^*(s, \mathbf{a}) \phi_h^*(s, \mathbf{a})^\top] + \lambda I_d$.

Lemma C.11 (Optimism for NE and CCE). *Consider an episode $n \in [N]$ and set $\alpha^{(n)} = \Theta \left(H \sqrt{n A^2 \zeta^{(n)} + d \lambda} \right)$. When the event \mathcal{E} holds and the policy $\pi^{(n)}$ is computed by solving NE or CCE, we have*

$$\bar{v}_i^{(n)}(s) - v_i^{\dagger, \pi^{(n)}}(s) \geq -H \sqrt{A \zeta^{(n)}} - 2H \tilde{\varepsilon}, \quad \forall n \in [N], i \in [M].$$

Proof. Denote $\tilde{\mu}_{h,i}^{(n)}(\cdot|s) := \arg \max_{\mu} \left(\mathbb{D}_{\mu, \pi^{(n)}} Q_{h,i}^{\dagger, \pi^{(n)}} \right)(s)$ and let $\tilde{\pi}_h^{(n)} = \tilde{\mu}_{h,i}^{(n)} \times \pi_{h,-i}^{(n)}$. Let $f_h^{(n)}(s, \mathbf{a}) = \left| \frac{1}{H} \left(\hat{P}_h^{(n)} - P_h^* \right) V_{h+1,i}^{\dagger, \pi^{(n)}} \right|(s, \mathbf{a})$, note that by definition, we have $\frac{1}{H} V_{h+1,i}^{\dagger, \pi^{(n)}}(s)$ is

bounded by 1, and

$$\begin{aligned} \frac{1}{H} V_{h+1,i}^{\dagger,\pi_{-i}^{(n)}}(s) &= \mathbb{E}_{\mathbf{a} \sim \bar{\pi}_h^{(n)}(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \frac{1}{H} \left(P_{h+1}^* V_{h+2,i}^{\dagger,\pi_{-i}^{(n)}} \right) (s, \mathbf{a}) \right] \\ &= \max_{\mu_{h+1,i}} \mathbb{E}_{\mathbf{a} \sim (\mu_{h+1,i} \times \pi_{h+1,-i}^{(n)})(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \frac{1}{H} \left(P_{h+1}^* V_{h+2,i}^{\dagger,\pi_{-i}^{(n)}} \right) (s, \mathbf{a}) \right] \in \mathcal{F}_{3,h}. \end{aligned}$$

where we use the result of Lemma C.1 and get $\frac{1}{H} \left(P_{h+1}^* V_{h+2,i}^{\dagger,\pi_{-i}^{(n)}} \right) (s, \mathbf{a})$ is a linear function in ϕ_{h+1}^* and the 2-norm of the weight is upper bounded by \sqrt{d} . Then according to the event \mathcal{E} , we have

$$\begin{aligned} \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] &\leq \zeta^{(n)}, \quad \mathbb{E}_{(s, \mathbf{a}) \sim \hat{\rho}_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}, \quad \forall n \in [N], h \in [H] \\ \|\phi_h(s, \mathbf{a})\|_{\left(\hat{\Sigma}_{h, \phi_h}^{(n)} \right)^{-1}} &= \Theta \left(\|\phi_h(s, \mathbf{a})\|_{\Sigma_{n, \rho_h^{(n)}, \phi_h}^{-1}} \right), \quad \forall n \in [N], h \in [H], \phi_h \in \Phi_h. \end{aligned}$$

A direct conclusion of the event \mathcal{E} is we can find an absolute constant c , such that

$$\begin{aligned} \beta_h^{(n)}(s, \mathbf{a}) &= \min \left\{ \alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\left(\Sigma_{h, \hat{\phi}_h}^{(n)} \right)^{-1}}, H \right\} \\ &\geq \min \left\{ c\alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \hat{\phi}_h}^{-1}}, H \right\}, \quad \forall n \in [N], h \in [H]. \end{aligned}$$

Next, we prove by induction that

$$\begin{aligned} &\mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\bar{\pi}^{(n)}}} \left[\bar{V}_{h,i}^{(n)}(s) - V_{h,i}^{\dagger,\pi_{-i}^{(n)}}(s) \right] \\ &\geq \sum_{h'=h}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h'}^{\bar{\pi}^{(n)}}} \left[\hat{\beta}_{h'}^{(n)}(s, \mathbf{a}) - H f_{h'}^{(n)}(s, \mathbf{a}) \right] - 2(H-h+1)\tilde{\varepsilon}, \quad \forall h \in [H]. \end{aligned} \quad (16)$$

First, notice that $\forall h \in [H]$,

$$\begin{aligned} \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\bar{\pi}^{(n)}}} \left[\bar{V}_{h,i}^{(n)}(s) - V_{h,i}^{\dagger,\pi_{-i}^{(n)}}(s) \right] &= \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\bar{\pi}^{(n)}}} \left[\left(\mathbb{D}_{\bar{\pi}_h^{(n)}} \bar{Q}_{h,i}^{(n)} \right) (s) - \left(\mathbb{D}_{\bar{\pi}_h^{(n)}} Q_{h,i}^{\dagger,\pi_{-i}^{(n)}} \right) (s) \right] \\ &\geq \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\bar{\pi}^{(n)}}} \left[\left(\mathbb{D}_{\bar{\pi}_h^{(n)}} \bar{Q}_{h,i}^{(n)} \right) (s) - \left(\mathbb{D}_{\bar{\pi}_h^{(n)}} Q_{h,i}^{\dagger,\pi_{-i}^{(n)}} \right) (s) \right] - 2\tilde{\varepsilon} \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\bar{\pi}^{(n)}}} \left[\bar{Q}_{h,i}^{(n)}(s, \mathbf{a}) - Q_{h,i}^{\dagger,\pi_{-i}^{(n)}}(s, \mathbf{a}) \right] - 2\tilde{\varepsilon}, \end{aligned}$$

where the inequality uses the result of Lemma C.8. Now we are ready to prove equation 16,

- When $h = H$, we have

$$\begin{aligned} \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, H}^{\bar{\pi}^{(n)}}} \left[\bar{V}_{H,i}^{(n)}(s) - V_{H,i}^{\dagger,\pi_{-i}^{(n)}}(s) \right] &\geq \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\bar{\pi}^{(n)}}} \left[\bar{Q}_{H,i}^{(n)}(s, \mathbf{a}) - Q_{H,i}^{\dagger,\pi_{-i}^{(n)}}(s, \mathbf{a}) \right] - 2\tilde{\varepsilon} \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\bar{\pi}^{(n)}}} \left[\hat{\beta}_H^{(n)}(s, \mathbf{a}) \right] - 2\tilde{\varepsilon} \\ &\geq \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\bar{\pi}^{(n)}}} \left[\hat{\beta}_H^{(n)}(s, \mathbf{a}) - H f_H^{(n)}(s, \mathbf{a}) \right] - 2\tilde{\varepsilon}. \end{aligned}$$

- Suppose the statement is true for $h+1$, then for step h , we have

$$\mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\bar{\pi}^{(n)}}} \left[\bar{V}_{h,i}^{(n)}(s) - V_{h,i}^{\dagger,\pi_{-i}^{(n)}}(s) \right]$$

$$\begin{aligned}
&\geq \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\overline{Q}_{h,i}^{(n)}(s, \mathbf{a}) - Q_{h,i}^{\dagger, \pi^{(n)}}(s, \mathbf{a}) \right] - 2\tilde{\varepsilon} \\
&= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \overline{V}_{h+1,i}^{(n)} \right)(s, \mathbf{a}) - \left(P_h^* V_{h+1,i}^{\dagger, \pi^{(n)}} \right)(s, \mathbf{a}) \right] - 2\tilde{\varepsilon} \\
&= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \left(\overline{V}_{h+1,i}^{(n)} - V_{h+1,i}^{\dagger, \pi^{(n)}} \right) \right)(s, \mathbf{a}) + \left(\left(\hat{P}_h^{(n)} - P_h^* \right) V_{h+1,i}^{\dagger, \pi^{(n)}} \right)(s, \mathbf{a}) \right] - 2\tilde{\varepsilon} \\
&= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\left(\hat{P}_h^{(n)} - P_h^* \right) V_{h+1,i}^{\dagger, \pi^{(n)}} \right)(s, \mathbf{a}) \right] + \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h+1}^{\pi^{(n)}}} \left[\overline{V}_{h+1,i}^{(n)}(s) - V_{h+1,i}^{\dagger, \pi^{(n)}}(s) \right] - 2\tilde{\varepsilon} \\
&\geq \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - H f_h^{(n)}(s, \mathbf{a}) \right] + \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h+1}^{\pi^{(n)}}} \left[\overline{V}_{h+1,i}^{(n)}(s) - V_{h+1,i}^{\dagger, \pi^{(n)}}(s) \right] - 2\tilde{\varepsilon} \\
&\geq \sum_{h'=h}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h'}^{\pi^{(n)}}} \left[\hat{\beta}_{h'}^{(n)}(s, \mathbf{a}) - H f_{h'}^{(n)}(s, \mathbf{a}) \right] - 2(H-h+1)\tilde{\varepsilon},
\end{aligned}$$

where the last row uses the induction assumption.

Therefore, we have proved equation 16. We then apply $h = 1$ to equation 16, and get

$$\begin{aligned}
&\mathbb{E}_{s \sim d_1} \left[\overline{V}_{1,i}^{(n)}(s) - V_{1,i}^{\dagger, \pi^{(n)}}(s) \right] \\
&= \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, 1}^{\pi^{(n)}}} \left[\overline{V}_{1,i}^{(n)}(s) - V_{1,i}^{\dagger, \pi^{(n)}}(s) \right] \\
&\geq \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - H f_h^{(n)}(s, \mathbf{a}) \right] - 2H\tilde{\varepsilon} \\
&= \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] - H \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[f_h^{(n)}(s, \mathbf{a}) \right] - 2H\tilde{\varepsilon}.
\end{aligned}$$

For the second term, since $\frac{1}{H} \hat{P}_h^{(n)} V_{h+1,i}^{\dagger, \pi^{(n)}}$ is linear in $\hat{\phi}_h^{(n)}$ and $\frac{1}{H} P_h^* V_{h+1,i}^{\dagger, \pi^{(n)}}$ is linear in ϕ_h^* , and according to the result of Lemma C.1, the 2-norm of their weights are both upper bounded by \sqrt{d} . Therefore, we have $\mathbb{E}_{\mathbf{a} \sim U(\mathcal{A})} \left[f_h^{(n)}(\cdot, \mathbf{a}) \right] \in \mathcal{F}_{1,h}$. By Lemma C.9, we have for $h = 1$,

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, 1}^{\pi^{(n)}}} \left[f_1^{(n)}(s, \mathbf{a}) \right] \leq \sqrt{A \mathbb{E}_{(s, \mathbf{a}) \sim \rho_1^{(n)}} \left[\left(f_1^{(n)}(s, \mathbf{a}) \right)^2 \right]} \leq \sqrt{A \zeta^{(n)}}.$$

And $\forall h \geq 2$, we have

$$\begin{aligned}
&\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[f_h^{(n)}(s, \mathbf{a}) \right] \\
&\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\pi^{(n)}}} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1} \begin{smallmatrix} (n) \\ n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)} \end{smallmatrix}} \sqrt{nA^2 \mathbb{E}_{(s, \mathbf{a}) \sim \tilde{\rho}_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] + d\lambda + nA^2 \zeta^{(n)}, 1} \right\} \right] \\
&\lesssim \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\pi^{(n)}}} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1} \begin{smallmatrix} (n) \\ n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)} \end{smallmatrix}} \sqrt{nA^2 \zeta^{(n)} + d\lambda}, 1 \right\} \right].
\end{aligned}$$

Note that we here use $f_h^{(n)}(s, \mathbf{a}) \leq 1$, $\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}$ and

$\mathbb{E}_{(s, \mathbf{a}) \sim \tilde{\rho}_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}$. Then according to our choice of $\alpha^{(n)}$, we get

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[f_h^{(n)}(s, \mathbf{a}) \right] \leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\pi^{(n)}}} \left[\min \left\{ \frac{c\alpha^{(n)}}{H} \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1} \begin{smallmatrix} (n) \\ n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)} \end{smallmatrix}}, 1 \right\} \right].$$

Combining all things together,

$$\begin{aligned}
\bar{v}_i^{(n)} - v_i^{\dagger, \pi_{-i}^{(n)}} &= \mathbb{E}_{s \sim d_1} \left[\bar{V}_{1,i}^{(n)}(s) - V_{1,i}^{\dagger, \pi_{-i}^{(n)}}(s) \right] \\
&\geq \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - H \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}} \left[f_h^{(n)}(s, \mathbf{a}) \right] - 2H\tilde{\varepsilon} \right] \\
&\geq \sum_{h=1}^{H-1} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - \min \left\{ c\alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \hat{\phi}_h^{(n)}}^{-1}}, H \right\} \right] - H\sqrt{A\zeta^{(n)}} - 2H\tilde{\varepsilon} \\
&= -H\sqrt{A\zeta^{(n)}} - 2H\tilde{\varepsilon},
\end{aligned}$$

which proves the inequality. \square

Lemma C.12 (Optimism for CE). *Consider an episode $n \in [N]$ and set $\alpha^{(n)} = \Theta \left(H\sqrt{nA^2\zeta^{(n)}} + d\lambda \right)$. When the event \mathcal{E} holds, we have*

$$\bar{v}_i^{(n)}(s) - \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi^{(n)}}(s) \geq -H\sqrt{A\zeta^{(n)}} - 2H\tilde{\varepsilon}, \quad \forall n \in [N], i \in [M].$$

Proof. Denote $\tilde{\omega}_{h,i}^{(n)} = \arg \max_{\omega_h \in \Omega_{h,i}} \left(\mathbb{D}_{\omega_h \circ \pi_h^{(n)}} \max_{\omega \in \Omega_i} Q_{h,i}^{\omega \circ \pi^{(n)}} \right)(s)$ and let $\tilde{\pi}_h^{(n)} = \tilde{\omega}_{h,i} \circ \pi_h^{(n)}$. Let $f_h^{(n)}(s, \mathbf{a}) = \left| \frac{1}{H} \left(\hat{P}_h^{(n)} - P_h^* \right) \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}} \right| (s, \mathbf{a})$, note that by definition, we have $\frac{1}{H} \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}}(s)$ is bounded by 1, and

$$\frac{1}{H} \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}}(s) = \max_{\omega_{h+1,i} \in \Omega_{h+1,i}} \mathbb{E}_{\mathbf{a} \sim (\omega_{h+1,i} \circ \pi_h)(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \frac{1}{H} \left(P_{h+1}^* \max_{\omega \in \Omega_i} V_{h+2,i}^{\omega \circ \pi^{(n)}} \right)(s, \mathbf{a}) \right] \in \mathcal{F}_{3,h}.$$

where we use the result of Lemma C.1 and get $\frac{1}{H} \left(P_{h+1}^* \max_{\omega \in \Omega_i} V_{h+2,i}^{\omega \circ \pi^{(n)}} \right)(s, \mathbf{a})$ is a linear function in ϕ_h^* and the 2-norm of the weight is upper bounded by \sqrt{d} . Then according to the event \mathcal{E} , we have

$$\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}, \quad \mathbb{E}_{(s, \mathbf{a}) \sim \hat{\rho}_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}, \quad \forall n \in [N], h \in [H]$$

$$\|\phi_h(s, \mathbf{a})\|_{\left(\hat{\Sigma}_{h, \hat{\phi}_h}^{(n)} \right)^{-1}} = \Theta \left(\|\phi_h(s, \mathbf{a})\|_{\Sigma_{n, \rho_h^{(n)}, \phi_h}^{-1}} \right), \quad \forall n \in [N], h \in [H], \phi_h \in \Phi_h.$$

A direct conclusion of the event \mathcal{E} is we can find an absolute constant c , such that

$$\begin{aligned}
\beta_h^{(n)}(s, \mathbf{a}) &= \min \left\{ \alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\left(\Sigma_{h, \hat{\phi}_h}^{(n)} \right)^{-1}}, H \right\} \\
&\geq \min \left\{ c\alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \hat{\phi}_h}^{-1}}, H \right\}, \quad \forall n \in [N], h \in [H].
\end{aligned}$$

Next, we prove by induction that

$$\begin{aligned}
&\mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}} \left[\bar{V}_{h,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{h,i}^{\omega \circ \pi^{(n)}}(s) \right] \\
&\geq \sum_{h'=h}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h'}} \left[\hat{\beta}_{h'}^{(n)}(s, \mathbf{a}) - H f_{h'}^{(n)}(s, \mathbf{a}) \right] - 2(H-h+1)\tilde{\varepsilon}, \quad \forall h \in [H]. \quad (17)
\end{aligned}$$

First, notice that $\forall h \in [H]$,

$$\mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}} \left[\bar{V}_{h,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{h,i}^{\omega \circ \pi^{(n)}}(s) \right] = \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}} \left[\left(\mathbb{D}_{\pi_h^{(n)}} \bar{Q}_{h,i}^{(n)} \right)(s) - \left(\mathbb{D}_{\pi_h^{(n)}} \max_{\omega \in \Omega_i} Q_{h,i}^{\omega \circ \pi^{(n)}} \right)(s) \right]$$

$$\begin{aligned}
&\geq \mathbb{E}_{s \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\left(\mathbb{D}_{\hat{\pi}_h^{(n)}} \overline{Q}_{h,i}^{(n)} \right) (s) - \left(\mathbb{D}_{\hat{\pi}_h^{(n)}} \max_{\omega \in \Omega_i} Q_{h,i}^{\omega \circ \pi^{(n)}} \right) (s) \right] - 2\tilde{\varepsilon} \\
&= \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\overline{Q}_{h,i}^{(n)}(s, \mathbf{a}) - \max_{\omega \in \Omega_i} Q_{h,i}^{\omega \circ \pi^{(n)}}(s, \mathbf{a}) \right] - 2\tilde{\varepsilon}.
\end{aligned}$$

where the inequality uses the result of Lemma C.8. Now we are ready to prove equation 17,

- When $h = H$, we have

$$\begin{aligned}
\mathbb{E}_{s \sim d_{\hat{P}^{(n)},H}^{\pi^{(n)}}} \left[\overline{V}_{H,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{H,i}^{\omega \circ \pi^{(n)}}(s) \right] &\geq \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},H}^{\pi^{(n)}}} \left[\overline{Q}_{H,i}^{(n)}(s, \mathbf{a}) - \max_{\omega \in \Omega_i} Q_{H,i}^{\omega \circ \pi^{(n)}}(s, \mathbf{a}) \right] \\
&= \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},H}^{\pi^{(n)}}} \left[\hat{\beta}_H^{(n)}(s, \mathbf{a}) \right] \\
&\geq \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},H}^{\pi^{(n)}}} \left[\hat{\beta}_H^{(n)}(s, \mathbf{a}) - H f_H^{(n)}(s, \mathbf{a}) \right] - 2\tilde{\varepsilon}.
\end{aligned}$$

- Suppose the statement is true for $h + 1$, then for step h , we have

$$\begin{aligned}
&\mathbb{E}_{s \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\overline{V}_{h,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{h,i}^{\omega \circ \pi^{(n)}}(s) \right] \\
&\geq \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\overline{Q}_{h,i}^{(n)}(s, \mathbf{a}) - \max_{\omega \in \Omega_i} Q_{h,i}^{\omega \circ \pi^{(n)}}(s, \mathbf{a}) \right] - 2\tilde{\varepsilon} \\
&= \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \overline{V}_{h+1,i}^{(n)} \right) (s, \mathbf{a}) - \left(P_h^* \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}} \right) (s, \mathbf{a}) \right] - 2\tilde{\varepsilon} \\
&= \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \left(\overline{V}_{h+1,i}^{(n)} - \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}} \right) \right) (s, \mathbf{a}) - \left(\left(\hat{P}_h^{(n)} - P_h^* \right) \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}} \right) (s, \mathbf{a}) \right] \\
&\quad - 2\tilde{\varepsilon} \\
&= \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - \left(\left(\hat{P}_h^{(n)} - P_h^* \right) \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}} \right) (s, \mathbf{a}) \right] \\
&\quad + \mathbb{E}_{s \sim d_{\hat{P}^{(n)},h+1}^{\pi^{(n)}}} \left[\overline{V}_{h+1,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}}(s) \right] - 2\tilde{\varepsilon} \\
&\geq \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - H f_h^{(n)}(s, \mathbf{a}) \right] + \mathbb{E}_{s \sim d_{\hat{P}^{(n)},h+1}^{\pi^{(n)}}} \left[\overline{V}_{h+1,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}}(s) \right] - 2\tilde{\varepsilon} \\
&\geq \sum_{h'=h}^H \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h'}^{\pi^{(n)}}} \left[\hat{\beta}_{h'}^{(n)}(s, \mathbf{a}) - H f_{h'}^{(n)}(s, \mathbf{a}) \right] - 2(H-h+1)\tilde{\varepsilon},
\end{aligned}$$

where the last row uses the induction assumption.

Therefore, we have proved equation 17. We then apply $h = 1$ to equation 17, and get

$$\begin{aligned}
&\mathbb{E}_{s \sim d_1} \left[\overline{V}_{1,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{1,i}^{\omega \circ \pi^{(n)}}(s) \right] \\
&= \mathbb{E}_{s \sim d_{\hat{P}^{(n)},1}^{\pi^{(n)}}} \left[\overline{V}_{1,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{1,i}^{\omega \circ \pi^{(n)}}(s) \right] \\
&\geq \sum_{h=1}^H \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - H f_h^{(n)}(s, \mathbf{a}) \right] - 2H\tilde{\varepsilon} \\
&= \sum_{h=1}^H \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] - H \sum_{h=1}^H \mathbb{E}_{(s,\mathbf{a}) \sim d_{\hat{P}^{(n)},h}^{\pi^{(n)}}} \left[f_h^{(n)}(s, \mathbf{a}) \right] - 2H\tilde{\varepsilon}.
\end{aligned}$$

For the second term, since $\frac{1}{H} \hat{P}_h^{(n)} \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}}$ is linear in $\hat{\phi}_h^{(n)}$ and $\frac{1}{H} P_h^* \max_{\omega \in \Omega_i} V_{h+1,i}^{\omega \circ \pi^{(n)}}$ is linear in ϕ_h^* , and according to the result of Lemma C.1, the 2-norm of their weights are both upper

bounded by \sqrt{d} . Therefore, we have $\mathbb{E}_{\mathbf{a} \sim U(\mathcal{A})} [f_h^{(n)}(\cdot, \mathbf{a})] \in \mathcal{F}_{1,h}$. By Lemma C.9, we have for $h = 1$,

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, 1}^{\pi^{(n)}}} [f_1^{(n)}(s, \mathbf{a})] \leq \sqrt{A \mathbb{E}_{(s, \mathbf{a}) \sim \rho_1^{(n)}} \left[\left(f_1^{(n)}(s, \mathbf{a}) \right)^2 \right]} \leq \sqrt{A \zeta^{(n)}}.$$

And $\forall h \geq 2$, we have

$$\begin{aligned} & \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} [f_h^{(n)}(s, \mathbf{a})] \\ & \leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\pi^{(n)}}} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}^{-1}} \sqrt{n A^2 \mathbb{E}_{(s, \mathbf{a}) \sim \hat{\rho}_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] + d\lambda + n A^2 \zeta^{(n)}}, 1 \right\} \right] \\ & \lesssim \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\pi^{(n)}}} \left[\min \left\{ \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}^{-1}} \sqrt{n A^2 \zeta^{(n)} + d\lambda}, 1 \right\} \right]. \end{aligned}$$

Note that we here use $f_h^{(n)}(s, \mathbf{a}) \leq 1$, $\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}$ and

$\mathbb{E}_{(s, \mathbf{a}) \sim \hat{\rho}_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}$. Then according to our choice of $\alpha^{(n)}$, we get

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} [f_h^{(n)}(s, \mathbf{a})] \leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\pi^{(n)}}} \left[\min \left\{ \frac{c\alpha^{(n)}}{H} \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \hat{\phi}_{h-1}^{(n)}}^{-1}}, 1 \right\} \right].$$

Combining all things together,

$$\begin{aligned} & \bar{v}_i^{(n)} - \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi^{(n)}} \\ & = \mathbb{E}_{s \sim d_1} \left[\bar{V}_{1,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{1,i}^{\omega \circ \pi^{(n)}}(s) \right] \\ & \geq \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} [\hat{\beta}_h^{(n)}(s, \mathbf{a})] - H \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} [f_h^{(n)}(s, \mathbf{a})] - 2H\tilde{\varepsilon} \\ & \geq \sum_{h=1}^{H-1} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - \min \left\{ c\alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \hat{\phi}_h^{(n)}}^{-1}}, H \right\} \right] - H\sqrt{A\zeta^{(n)}} - 2H\tilde{\varepsilon} \\ & = -H\sqrt{A\zeta^{(n)}} - 2H\tilde{\varepsilon}, \end{aligned}$$

which proves the inequality. \square

Lemma C.13 (pessimism). *Consider an episode $n \in [N]$ and set $\alpha^{(n)} = \Theta \left(H\sqrt{nA^2\zeta^{(n)} + d\lambda} \right)$. When the event \mathcal{E} holds, we have*

$$\underline{v}_i^{(n)}(s) - v_i^{\pi^{(n)}}(s) \leq H\sqrt{A\zeta^{(n)}}, \quad \forall n \in [N], i \in [M].$$

Proof. Let $\tilde{f}_h^{(n)}(s, \mathbf{a}) = \left| \frac{1}{H} \left(\hat{P}_h^{(n)} - P_h^* \right) V_{h+1,i}^{\pi^{(n)}}(s, \mathbf{a}) \right|$, note that by definition, we have $\frac{1}{H} V_{h+1,i}^{\pi^{(n)}}(s)$ is bounded by 1, and

$$\frac{1}{H} V_{h+1,i}^{\pi^{(n)}}(s) = \mathbb{E}_{\mathbf{a} \sim \pi_h^{(n)}(s)} \left[\frac{r_{h+1,i}(s, \mathbf{a})}{H} + \frac{1}{H} \left(P_{h+1}^* V_{h+2,i}^{\pi^{(n)}}(s, \mathbf{a}) \right) \right] \in \mathcal{F}_{2,h}.$$

where we use the result of Lemma C.1 and get $\frac{1}{H} \left(P_{h+1}^* V_{h+2,i}^{\pi^{(n)}}(s, \mathbf{a}) \right)$ is a linear function in ϕ_{h+1}^* and the 2-norm of the weight is upper bounded by \sqrt{d} . Then according to the event \mathcal{E} , we have

$$\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}, \quad \mathbb{E}_{(s, \mathbf{a}) \sim \hat{\rho}_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}, \quad \forall n \in [N], h \in [H]$$

$$\|\phi_h(s, \mathbf{a})\|_{\left(\hat{\Sigma}_{h, \phi_h}^{(n)}\right)^{-1}} = \Theta \left(\|\phi_h(s, \mathbf{a})\|_{\Sigma_{n, \rho_h^{(n)}, \phi_h}^{-1}} \right), \quad \forall n \in [N], h \in [H], \phi_h \in \Phi_h.$$

A direct conclusion of the event \mathcal{E} is we can find an absolute constant c , such that

$$\begin{aligned} \beta_h^{(n)}(s, \mathbf{a}) &= \min \left\{ \alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\left(\hat{\Sigma}_{h, \hat{\phi}_h}^{(n)}\right)^{-1}}, H \right\} \\ &\geq \min \left\{ c\alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \hat{\phi}_h}^{-1}}, H \right\}, \quad \forall n \in [N], h \in [H]. \end{aligned}$$

Again, we prove the following inequality by induction:

$$\mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[V_{h,i}^{(n)}(s) - V_{h,i}^{\pi^{(n)}}(s) \right] \leq \sum_{h'=h}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h'}^{\pi^{(n)}}} \left[-\hat{\beta}_{h'}^{(n)}(s, \mathbf{a}) + H f_{h'}^{(n)}(s, \mathbf{a}) \right], \quad \forall h \in [H]. \quad (18)$$

- When $h = H$, we have

$$\begin{aligned} \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, H}^{\pi^{(n)}}} \left[V_{H,i}^{(n)}(s) - V_{H,i}^{\pi^{(n)}}(s) \right] &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\pi^{(n)}}} \left[\underline{Q}_{H,i}^{(n)}(s, \mathbf{a}) - Q_{H,i}^{\pi^{(n)}}(s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\pi^{(n)}}} \left[-\hat{\beta}_H^{(n)}(s, \mathbf{a}) \right] \\ &\leq \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, H}^{\pi^{(n)}}} \left[-\hat{\beta}_H^{(n)}(s, \mathbf{a}) + H f_H^{(n)}(s, \mathbf{a}) \right] \end{aligned}$$

- Suppose the statement is true for $h+1$, then for step h , we have

$$\begin{aligned} &\mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[V_{h,i}^{(n)}(s) - V_{h,i}^{\pi^{(n)}}(s) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\underline{Q}_{h,i}^{(n)}(s, \mathbf{a}) - Q_{h,i}^{\pi^{(n)}}(s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[-\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \underline{V}_{h+1,i}^{(n)} \right)(s, \mathbf{a}) - \left(P_h^* V_{h+1,i}^{\pi^{(n)}} \right)(s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[-\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\hat{P}_h^{(n)} \left(V_{h+1,i}^{(n)} - V_{h+1,i}^{\pi^{(n)}} \right) \right)(s, \mathbf{a}) + \left(\left(\hat{P}_h^{(n)} - P_h^* \right) V_{h+1,i}^{\pi^{(n)}} \right)(s, \mathbf{a}) \right] \\ &= \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[-\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \left(\left(\hat{P}_h^{(n)} - P_h^* \right) V_{h+1,i}^{\pi^{(n)}} \right)(s, \mathbf{a}) \right] + \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h+1}^{\pi^{(n)}}} \left[\left(V_{h+1,i}^{(n)} - V_{h+1,i}^{\pi^{(n)}} \right)(s) \right] \\ &\leq \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[-\hat{\beta}_h^{(n)}(s, \mathbf{a}) + H f_h^{(n)}(s, \mathbf{a}) \right] + \mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h+1}^{\pi^{(n)}}} \left[\left(V_{h+1,i}^{(n)} - V_{h+1,i}^{\pi^{(n)}} \right)(s) \right] \\ &\leq \sum_{h'=h}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h'}^{\pi^{(n)}}} \left[-\hat{\beta}_{h'}^{(n)}(s, \mathbf{a}) + H f_{h'}^{(n)}(s, \mathbf{a}) \right]. \end{aligned}$$

where the last row uses the induction assumption.

The remaining steps are exactly the same as the proof in Lemma C.11 or Lemma C.12, we may prove

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, 1}^{\pi^{(n)}}} \left[\min \left\{ f_1^{(n)}(s, \mathbf{a}), 1 \right\} \right] \leq \sqrt{A\zeta^{(n)}},$$

and

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[f_h^{(n)}(s, \mathbf{a}) \right] \leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\pi^{(n)}}} \left[\min \left\{ \frac{c\alpha^{(n)}}{H} \left\| \hat{\phi}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}, \hat{\phi}_{h-1}}^{(n)}}^{-1}}, 1 \right\} \right], \quad \forall h \geq 2.$$

Combining all things together, we get

$$\begin{aligned}
\underline{v}_i^{(n)} - \overline{v}_i^{(n)} &= \mathbb{E}_{s \sim d_1} \left[\underline{V}_{1,i}^{(n)}(s) - \overline{V}_{1,i}^{(n)}(s) \right] \\
&\leq \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}} \left[-\hat{\beta}_h^{(n)}(s, \mathbf{a}) + H f_h^{(n)}(s, \mathbf{a}) \right] \\
&\leq \sum_{h=1}^{H-1} \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}} \left[-\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \min \left\{ c\alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \hat{\phi}_h^{(n)}}^{-1}}, H \right\} \right] + H\sqrt{A\zeta^{(n)}} \\
&\leq H\sqrt{A\zeta^{(n)}},
\end{aligned}$$

which has finished the proof. \square

Lemma C.14. *For the model-free algorithm, suppose N is large enough, when we pick $\lambda = \Theta\left(d \log \frac{NH|\Phi|}{\delta}\right)$, $\zeta^{(n)} = \Theta\left(\frac{d^2 M}{n} \log \frac{dNHML|\Phi|}{\varepsilon\delta}\right)$, $L = \Theta(NHAMd)$, $\tilde{\varepsilon} = \frac{1}{2HN}$ and $\alpha^{(n)} = \Theta\left(H\sqrt{nA^2\zeta^{(n)}} + d\lambda\right)$, with probability $1 - \delta$, we have*

$$\sum_{n=1}^N \Delta^{(n)} \lesssim H^3 d^2 A^{\frac{3}{2}} N^{\frac{1}{2}} M^{\frac{1}{2}} \log \frac{dNHAM|\Phi|}{\delta}.$$

Proof. With our choice of λ and $\zeta^{(n)}$, according to Lemma C.6, we know \mathcal{E} holds with probability $1 - \delta$. Furthermore, with a proper choice of the absolute constants, we have

$$\begin{aligned}
\alpha^{(n)} &= \Theta\left(H\sqrt{d^2 A^2 M \log \frac{dNHML|\Phi|}{\delta}} + d^2 \log \frac{NH|\Phi|}{\delta}\right) \\
&\leq O\left(HdA\sqrt{M \log \frac{dNHMA|\Phi|}{\delta}}\right) \\
&\leq O(NHAMd) \leq L.
\end{aligned}$$

Let $f_h^{(n)}(s, \mathbf{a}) = \frac{1}{2H^2} \left| \left(\hat{P}_h^{(n)} - P_h^* \right) \left(\overline{V}_{h+1,i}^{(n)} - \underline{V}_{h+1,i}^{(n)} \right) \right| (s, \mathbf{a})$. We first verify $\frac{1}{2H^2} \left(\overline{V}_{h+1,i}^{(n)} - \underline{V}_{h+1,i}^{(n)} \right) \in \mathcal{F}_{4,h}$. By definition, we have

$$\frac{1}{2H^2} \left(\overline{V}_{h+1,i}^{(n)} - \underline{V}_{h+1,i}^{(n)} \right) = \mathbb{E}_{\mathbf{a} \sim \pi_h^{(n)}(s)} \left[\frac{1}{H^2} \hat{\beta}_{h+1}^{(n)}(s, \mathbf{a}) + \frac{1}{2H^2} P_{h+1}^* \left(\overline{V}_{h+2,i}^{(n)} - \underline{V}_{h+2,i}^{(n)} \right) (s, \mathbf{a}) \right]$$

The first term is equal to $\frac{1}{H^2} \min \left(\alpha^{(n)} \sqrt{\hat{\phi}_h^{(n)}(s, \mathbf{a})^\top \left(\hat{\Sigma}_h^{(n)} \right)^{-1} \hat{\phi}_h^{(n)}(s, \mathbf{a}), H} \right)$, which is exactly the same as that in the definition of $\mathcal{F}_{4,h}$ (note that we use the property $\alpha^{(n)} \leq L, \forall n \in [N]$). For the second term, note that we have $0 \leq \frac{1}{2H^2} \left(\overline{V}_{h,i}^{(n)} - \underline{V}_{h,i}^{(n)} \right) \leq 1, \forall h$. Therefore, by Lemma C.1, $\frac{1}{2H^2} P_{h+1}^* \left(\overline{V}_{h+2,i}^{(n)} - \underline{V}_{h+2,i}^{(n)} \right) (s, \mathbf{a})$ is a linear function in ϕ_{h+1}^* whose weight's 2-norm is upper bounded by \sqrt{d} . Combing the above arguments, we conclude $\frac{1}{2H^2} \left(\overline{V}_{h+1,i}^{(n)} - \underline{V}_{h+1,i}^{(n)} \right) \in \mathcal{F}_{4,h}$. According to the definition of the event \mathcal{E} , we have

$$\mathbb{E}_{s \sim \rho_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}, \quad \left\| \phi_h(s, \mathbf{a}) \right\|_{\left(\hat{\Sigma}_h^{(n)} \right)^{-1}} = \Theta \left(\left\| \phi_h(s, \mathbf{a}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \phi_h}^{-1}} \right), \quad \forall n \in [N], \phi_h \in \Phi_h, h \in [H]. \quad (19)$$

By definition, we have

$$\Delta^{(n)} = \max_{i \in [M]} \left\{ \overline{v}_i^{(n)} - \underline{v}_i^{(n)} \right\} + 2H\sqrt{A\zeta^{(n)}}.$$

For each fixed $i \in [M]$, $h \in [H]$ and $n \in [N]$, we have

$$\begin{aligned}
& \mathbb{E}_{s \sim d_{P^*,h}^{\pi(n)}} \left[\bar{V}_{h,i}^{(n)}(s) - \underline{V}_{h,i}^{(n)}(s) \right] \\
&= \mathbb{E}_{s \sim d_{P^*,h}^{\pi(n)}} \left[\left(\mathbb{D}_{\pi_h^{(n)}} \bar{Q}_{h,i}^{(n)} \right) (s) - \left(\mathbb{D}_{\pi_h^{(n)}} Q_{h,i}^{(n)} \right) (s) \right] \\
&= \mathbb{E}_{(s,\mathbf{a}) \sim d_{P^*,h}^{\pi(n)}} \left[\bar{Q}_{h,i}^{(n)}(s, \mathbf{a}) - Q_{h,i}^{(n)}(s, \mathbf{a}) \right] \\
&= \mathbb{E}_{(s,\mathbf{a}) \sim d_{P^*,h}^{\pi(n)}} \left[2\hat{\beta}_h^{(n)} + \left(\hat{P}_h^{(n)} \left(\bar{V}_{h+1,i}^{(n)} - \underline{V}_{h+1,i}^{(n)} \right) \right) (s, \mathbf{a}) \right] \\
&= \mathbb{E}_{(s,\mathbf{a}) \sim d_{P^*,h}^{\pi(n)}} \left[2\hat{\beta}_h^{(n)} + \left(\left(\hat{P}_h^{(n)} - P_h^* \right) \left(\bar{V}_{h+1,i}^{(n)} - \underline{V}_{h+1,i}^{(n)} \right) \right) (s, \mathbf{a}) \right] + \mathbb{E}_{s \sim d_{P^*,h+1}^{\pi(n)}} \left[\bar{V}_{h+1,i}^{(n)}(s) - \underline{V}_{h+1,i}^{(n)}(s) \right] \\
&\leq \mathbb{E}_{(s,\mathbf{a}) \sim d_{P^*,h}^{\pi(n)}} \left[2\hat{\beta}_h^{(n)}(s, \mathbf{a}) + 2H^2 f_h^{(n)}(s, \mathbf{a}) \right] + \mathbb{E}_{s \sim d_{P^*,h+1}^{\pi(n)}} \left[\bar{V}_{h+1,i}^{(n)}(s) - \underline{V}_{h+1,i}^{(n)}(s) \right] \\
&\leq \dots \\
&\leq 2 \sum_{h'=h}^H \mathbb{E}_{(s,\mathbf{a}) \sim d_{P^*,h'}^{\pi(n)}} \left[\hat{\beta}_{h'}^{(n)}(s, \mathbf{a}) + H^2 f_{h'}^{(n)}(s, \mathbf{a}) \right],
\end{aligned}$$

where the last inequality is calculated using induction. In particular,

$$\mathbb{E}_{s \sim d_{P^*,1}^{\pi(n)}} \left[\bar{V}_{1,i}^{(n)}(s) - \underline{V}_{1,i}^{(n)}(s) \right] \leq 2 \underbrace{\sum_{h=1}^H \mathbb{E}_{(s,\mathbf{a}) \sim d_{P^*,h}^{\pi(n)}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right]}_{(a)} + 2H^2 \underbrace{\sum_{h=1}^H \mathbb{E}_{(s,\mathbf{a}) \sim d_{P^*,h}^{\pi(n)}} \left[f_h^{(n)}(s, \mathbf{a}) \right]}_{(b)}. \quad (20)$$

First, we calculate the first term (a) in Inequality equation 20. Following Lemma C.10 and noting the bonus $\hat{\beta}_h^{(n)}$ is $O(H)$, we have

$$\begin{aligned}
& \sum_{h=1}^H \mathbb{E}_{(s,\mathbf{a}) \sim d_{P^*,h}^{\pi(n)}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] \\
&\lesssim \sum_{h=1}^H \mathbb{E}_{(s,\mathbf{a}) \sim d_{P^*,h}^{\pi(n)}} \left[\min \left(\alpha^{(n)} \left\| \hat{\phi}_h^{(n)}(s, \mathbf{a}) \right\|_{\Sigma_{n,\rho_h^{(n)},\hat{\phi}_h^{(n)}}^{-1}}, H \right) \right] \quad (\text{From equation 19}) \\
&\lesssim \sum_{h=1}^{H-1} \mathbb{E}_{(\tilde{s},\tilde{\mathbf{a}}) \sim d_{P^*,h}^{\pi(n)}} \left[\left\| \phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n,\gamma_h^{(n)},\phi_h^*}^{-1}} \right] \sqrt{nA (\alpha^{(n)})^2 \mathbb{E}_{(s,\mathbf{a}) \sim \rho_h^{(n)}} \left[\left\| \hat{\phi}_h^{(n)}(s, \mathbf{a}) \right\|_{\Sigma_{n,\rho_h^{(n)},\hat{\phi}_h^{(n)}}^{-1}}^2 \right]} + H^2 d \lambda \\
&+ \sqrt{A (\alpha^{(n)})^2 \mathbb{E}_{(s,\mathbf{a}) \sim \rho_1^{(n)}} \left[\left\| \hat{\phi}_1^{(n)}(s, \mathbf{a}) \right\|_{\Sigma_{n,\rho_1^{(n)},\hat{\phi}_1^{(n)}}^{-1}}^2 \right]}.
\end{aligned}$$

Note that we use the fact that $B = H$ when applying Lemma D.3. In addition, we have

$$\begin{aligned}
& n \mathbb{E}_{(s,\mathbf{a}) \sim \rho_h^{(n)}} \left[\left\| \hat{\phi}_h^{(n)}(s, \mathbf{a}) \right\|_{\Sigma_{n,\rho_h^{(n)},\hat{\phi}_h^{(n)}}^{-1}}^2 \right] \\
&= n \text{Tr} \left(\mathbb{E}_{(s,\mathbf{a}) \sim \rho_h^{(n)}} \left[\hat{\phi}_h^{(n)}(s, \mathbf{a}) \hat{\phi}_h^{(n)}(s, \mathbf{a})^\top \right] \left(n \mathbb{E}_{(s,\mathbf{a}) \sim \rho_h^{(n)}} \left[\hat{\phi}_h^{(n)}(s, \mathbf{a}) \hat{\phi}_h^{(n)}(s, \mathbf{a})^\top \right] + \lambda I_d \right)^{-1} \right) \\
&\leq d.
\end{aligned}$$

Then,

$$\sum_{h=1}^H \mathbb{E}_{(s,\mathbf{a}) \sim d_{P^*,h}^{\pi(n)}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] \leq \mathbb{E}_{(\tilde{s},\tilde{\mathbf{a}}) \sim d_{P^*,h}^{\pi(n)}} \left[\left\| \phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n,\gamma_h^{(n)},\phi_h^*}^{-1}} \right] \sqrt{dA (\alpha^{(n)})^2 + H^2 d \lambda} + \sqrt{dA (\alpha^{(n)})^2 / n}.$$

Second, we calculate the term (b) in inequality equation 23. Following Lemma D.3 and noting $(f_h^{(n)}(s, \mathbf{a}))^2$ is upper-bounded by 1 (i.e., $B = 1$ in Lemma D.3), we have

$$\begin{aligned}
& \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi(n)}} [f_h^{(n)}(s, \mathbf{a})] \\
& \leq \sum_{h=1}^{H-1} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\|\phi_h^*(\tilde{s}, \tilde{\mathbf{a}})\|_{\Sigma_{n, \gamma_h^{(n)}, \phi_h^*}^{-1}} \right] \sqrt{nA \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[(f_h^{(n)}(s, \mathbf{a}))^2 \right] + d\lambda} + \sqrt{A \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[(f_1^{(n)}(s, \mathbf{a}))^2 \right]} \\
& \leq \sum_{h=1}^{H-1} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\|\phi_h^*(\tilde{s}, \tilde{\mathbf{a}})\|_{\Sigma_{n, \gamma_h^{(n)}, \phi_h^*}^{-1}} \right] \sqrt{nA\zeta^{(n)} + d\lambda} + \sqrt{A\zeta^{(n)}} \\
& \lesssim \frac{\alpha^{(n)}}{H} \sum_{h=1}^{H-1} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\|\phi_h^*(\tilde{s}, \tilde{\mathbf{a}})\|_{\Sigma_{n, \gamma_h^{(n)}, \phi_h^*}^{-1}} \right] + \sqrt{A\zeta^{(n)}},
\end{aligned}$$

where in the second inequality, we use $\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[(f_h^{(n)}(s, \mathbf{a}))^2 \right] \leq \zeta^{(n)}$, and in the last line,

recall $\sqrt{nA\zeta^{(n)} + d\lambda} \lesssim \alpha^{(n)}/H$. Then, by combining the above calculation of the term (a) and term (b) in inequality equation 23, we have:

$$\begin{aligned}
\bar{v}_i^{(n)} - v_i^{(n)} &= \mathbb{E}_{s \sim d_{P^*, 1}^{\pi(n)}} \left[\bar{V}_{1, i}^{(n)}(s) - V_{1, i}^{(n)}(s) \right] \\
&\lesssim \sum_{h=1}^{H-1} \left(\mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\|\phi_h^*(\tilde{s}, \tilde{\mathbf{a}})\|_{\Sigma_{n, \gamma_h^{(n)}, \phi_h^*}^{-1}} \right] \sqrt{dA (\alpha^{(n)})^2 + H^2 d\lambda} + \sqrt{\frac{dA (\alpha^{(n)})^2}{n}} \right) \\
&\quad + H^2 \sum_{h=1}^{H-1} \left(\frac{\alpha^{(n)}}{H} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\|\phi_h^*(\tilde{s}, \tilde{\mathbf{a}})\|_{\Sigma_{n, \gamma_h^{(n)}, \phi_h^*}^{-1}} \right] + \sqrt{A\zeta^{(n)}} \right).
\end{aligned}$$

Taking maximum over i on both sides and use the definition of $\Delta^{(n)}$, we get

$$\begin{aligned}
\Delta^{(n)} &= \max_{i \in [M]} \left\{ \bar{v}_i^{(n)} - v_i^{(n)} \right\} + 2H \sqrt{A\zeta^{(n)}} \\
&\lesssim \sum_{h=1}^{H-1} \left(\mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\|\phi_h^*(\tilde{s}, \tilde{\mathbf{a}})\|_{\Sigma_{n, \gamma_h^{(n)}, \phi_h^*}^{-1}} \right] \sqrt{dA (\alpha^{(n)})^2 + H^2 d\lambda} + \sqrt{\frac{dA (\alpha^{(n)})^2}{n}} \right) \\
&\quad + H^2 \sum_{h=1}^{H-1} \left(\frac{\alpha^{(n)}}{H} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\|\phi_h^*(\tilde{s}, \tilde{\mathbf{a}})\|_{\Sigma_{n, \gamma_h^{(n)}, \phi_h^*}^{-1}} \right] + \sqrt{A\zeta^{(n)}} \right).
\end{aligned}$$

Hereafter, we take the dominating term out. Note that

$$\begin{aligned}
\sum_{n=1}^N \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\|\phi_h^*(\tilde{s}, \tilde{\mathbf{a}})\|_{\Sigma_{n, \gamma_h^{(n)}, \phi_h^*}^{-1}} \right] &\leq \sqrt{N \sum_{n=1}^N \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\phi_h^*(\tilde{s}, \tilde{\mathbf{a}})^\top \Sigma_{n, \gamma_h^{(n)}, \phi_h^*}^{-1} \phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \right]} \\
&\quad \text{(CS inequality)} \\
&\lesssim \sqrt{N \left(\log \det \left(\sum_{n=1}^N \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\phi_h^*(\tilde{s}, \tilde{\mathbf{a}}) \phi_h^*(\tilde{s}, \tilde{\mathbf{a}})^\top \right] \right) - \log \det(\lambda I_d) \right)} \quad \text{(Lemma E.2)} \\
&\leq \sqrt{dN \log \left(1 + \frac{N}{d\lambda} \right)}.
\end{aligned}$$

(Potential function bound, Lemma E.3 noting $\|\phi_h^*(s, \mathbf{a})\|_2 \leq 1$ for any (s, \mathbf{a}) .)

Finally,

$$\sum_{n=1}^N \Delta^{(n)} \lesssim H \left(\sqrt{dN \log \left(1 + \frac{N}{d} \right)} \sqrt{dA (\alpha^{(N)})^2 + H^2 d\lambda} + \sum_{n=1}^N \sqrt{\frac{dA (\alpha^{(n)})^2}{n}} \right)$$

$$\begin{aligned}
& + H^3 \left(\frac{1}{H} \sqrt{dN \log \left(1 + \frac{N}{d\lambda} \right)} \alpha^{(N)} + \sum_{n=1}^N \sqrt{A\zeta^{(n)}} \right) + 2HN\tilde{\varepsilon} \\
& \lesssim H^2 d \sqrt{NA \log \left(1 + \frac{N}{d\lambda} \right)} \alpha^{(N)} \\
& \text{(Some algebra. We take the dominating term out. Note that } \alpha^{(n)} \text{ is increasing in } n) \\
& \lesssim H^3 d^2 A^{\frac{3}{2}} N^{\frac{1}{2}} M^{\frac{1}{2}} \log \frac{dNHAM|\Phi|}{\delta}.
\end{aligned}$$

This concludes the proof. \square

Proof of Theorem 4.2

Proof. For any fixed episode n and agent i , by Lemma C.11, Lemma C.12 and Lemma C.13, we have

$$v_i^{\dagger, \pi^{(n)}} - v_i^{\pi^{(n)}} \left(\text{or } \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi^{(n)}} - v_i^{\pi^{(n)}} \right) \leq \bar{v}_i^{(n)} - \underline{v}_i^{(n)} + 2\sqrt{A\zeta^{(n)}} + 2H\tilde{\varepsilon} \leq \Delta^{(n)} + 2H\tilde{\varepsilon}.$$

Taking maximum over i on both sides, we have

$$\max_{i \in [M]} \left\{ v_i^{\dagger, \pi^{(n)}} - v_i^{\pi^{(n)}} \right\} \left(\text{or } \max_{i \in [M]} \left\{ \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi^{(n)}} - v_i^{\pi^{(n)}} \right\} \right) \leq \Delta^{(n)} + 2H\tilde{\varepsilon}. \quad (21)$$

From Lemma C.14, with probability $1 - \delta$, we can ensure

$$\sum_{n=1}^N (\Delta^{(n)} + 2H\tilde{\varepsilon}) \lesssim H^3 d^2 A^{\frac{3}{2}} N^{\frac{1}{2}} M^{\frac{1}{2}} \log \frac{dNHAM|\Phi|}{\delta}.$$

Therefore, according to Lemma E.4, when we pick N to be

$$O \left(\frac{H^6 d^4 A^3 M}{\varepsilon^2} \log^2 \left(\frac{HdAM|\Phi|}{\delta\varepsilon} \right) \right),$$

we have

$$\frac{1}{N} \sum_{n=1}^N (\Delta^{(n)} + 2H\tilde{\varepsilon}) \leq \varepsilon.$$

On the other hand, from equation 21, we have

$$\begin{aligned}
& \max_{i \in [M]} \left\{ v_i^{\dagger, \hat{\pi}^{(n)}} - v_i^{\hat{\pi}^{(n)}} \right\} \left(\text{or } \max_{i \in [M]} \left\{ \max_{\omega \in \Omega_i} v_i^{\omega \circ \hat{\pi}^{(n)}} - v_i^{\hat{\pi}^{(n)}} \right\} \right) \\
& = \max_{i \in [M]} \left\{ v_i^{\dagger, \pi^{(n^*)}} - v_i^{\pi^{(n^*)}} \right\} \left(\text{or } \max_{i \in [M]} \left\{ \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi^{(n^*)}} - v_i^{\pi^{(n^*)}} \right\} \right) \\
& \leq \Delta^{(n^*)} + 2H\tilde{\varepsilon} = \min_{n \in [N]} \Delta^{(n)} + 2H\tilde{\varepsilon} \leq \frac{1}{N} \sum_{n=1}^N (\Delta^{(n)} + 2H\tilde{\varepsilon}) \leq \varepsilon,
\end{aligned}$$

which has finished the proof. \square

D ANALYSIS OF THE FACTORED MARKOV GAMES

D.1 HIGH PROBABILITY EVENTS

Define the set $\bar{\Phi}_{h,i} = \{\bar{\phi}_{h,i}(s, \mathbf{a}) := \bigotimes_{j \in Z_i} \phi_{h,j}(s[Z_j], \mathbf{a}_j) \mid \phi_{h,j} \in \Phi_{h,j}\}$. Let $|\Phi| = \max_{h,j} |\Phi_{h,j}|$ and $|\bar{\Phi}| = \max_{h,i} |\bar{\Phi}_{h,i}|$. Clearly, we have $|\bar{\Phi}| \leq |\Phi|^L$. Define the following event

$$\mathcal{E}_1 : \forall n \in [N], h \in [H], i \in [M], \rho \in \left\{ \rho_h^{(n)}, \tilde{\rho}_h^{(n)} \right\}, \quad \mathbb{E}_\rho \left[\left\| \hat{P}_{h,i}^{(n)}(\cdot \mid s[Z_i], \mathbf{a}_i) - P_{h,i}^*(\cdot \mid s[Z_i], \mathbf{a}_i) \right\|_1^2 \right] \leq \zeta^{(n)},$$

$$\mathcal{E}_2 : \forall n \in [N], h \in [H], i \in [M], \bar{\phi}_{h,i} \in \bar{\Phi}_{h,i}, \quad \|\bar{\phi}_{h,i}(s, \mathbf{a})\|_{\left(\hat{\Sigma}_{h, \bar{\phi}_{h,i}}^{(n)}\right)^{-1}} = \Theta \left(\|\bar{\phi}_{h,i}(s, \mathbf{a})\|_{\Sigma_{n, \rho_h^{(n)}, \bar{\phi}_{h,i}}^{-1}} \right)$$

$$\mathcal{E} := \mathcal{E}_1 \cap \mathcal{E}_2.$$

The following lemma shows that the event \mathcal{E} holds with a high probability with proper choices of the parameters.

Lemma D.1. *When $\hat{P}_h^{(n)}$ is computed using Alg. 2, if we set*

$$\lambda = \Theta \left(Ld^L \log \frac{NHM|\Phi|}{\delta} \right), \quad \zeta^{(n)} = \Theta \left(\frac{1}{n} \log \frac{|\mathcal{M}|HNM}{\delta} \right),$$

then \mathcal{E} holds with probability at least $1 - \delta$.

The proof of Lemma D.1 is follows a similar procedure as that of Lemma B.2, with minor changes on the notations as well as some modifications on the union bound.

D.2 STATISTICAL GUARANTEES

Lemma D.2 (One-step back inequality for the learned model). *Suppose the event \mathcal{E} holds. Consider a set of functions $\{g_h\}_{h=1}^H$ that satisfies $g_h \in \mathcal{S}[Z_i] \times \mathcal{A}_i \rightarrow \mathbb{R}_+$, s.t. $\|g_h\|_\infty \leq B$. For a given policy π , we have*

$$\begin{aligned} & \left| \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}_h^{(n)}, h}^\pi} [g_h(s[Z_i], \mathbf{a}_i)] \right| \\ & \leq \left\{ \begin{array}{l} \sqrt{\tilde{A} \mathbb{E}_{(s, \mathbf{a}) \sim \rho_1^{(n)}} [g_1^2(s[Z_i], \mathbf{a}_i)]}, \quad h = 1 \\ \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}_h^{(n)}, h-1}^\pi} \left[\min \left\{ \tilde{A} \left\| \bar{\phi}_{h-1, i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \bar{\phi}_{h-1, i}^{(n)}}^{-1}} \sqrt{n \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} [g_h^2(s[Z_i], \mathbf{a}_i)] + B^2 \lambda d^L + n B^2 \zeta^{(n)}}, B \right\} \right] \end{array} \right\}, \quad h \geq 2 \end{aligned}$$

$$\text{where } \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) := \bigotimes_{j \in Z_i} \hat{\phi}_{h-1, j}^{(n)}(s[Z_j], \mathbf{a}_j), \quad \text{and } \Sigma_{n, \rho_h^{(n)}, \bar{\phi}_{h,i}^{(n)}} = n \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a})^\top \right] + \lambda I_{d^{Z_i}}.$$

Proof. For step $h = 1$, we have

$$\begin{aligned} \mathbb{E}_{(s, \mathbf{a}_i) \sim d_{\hat{P}_h^{(n)}, 1}^\pi} [g_1(s[Z_i], \mathbf{a}_i)] &= \mathbb{E}_{s \sim d_1, \mathbf{a}_i \sim \pi_1(s)} [g_1(s[Z_i], \mathbf{a}_i)] \\ &\leq \sqrt{\max_{(s, \mathbf{a}_i)} \frac{d_1(s) \pi_1(\mathbf{a}_i | s)}{\rho_1^{(n)}(s, \mathbf{a}_i)} \mathbb{E}_{(s', \mathbf{a}'_i) \sim \rho_1^{(n)}} [g_1^2(s'[Z_i], \mathbf{a}'_i)]} \\ &= \sqrt{\max_{(s, \mathbf{a}_i)} \frac{d_1(s) \pi_1(\mathbf{a}_i | s)}{d_1(s) u_{\mathcal{A}}(\mathbf{a}_i)} \mathbb{E}_{(s', \mathbf{a}'_i) \sim \rho_1^{(n)}} [g_1^2(s'[Z_i], \mathbf{a}'_i)]} \\ &\leq \sqrt{\tilde{A} \mathbb{E}_{(s, \mathbf{a}_i) \sim \rho_1^{(n)}} [g_1^2(s[Z_i], \mathbf{a}_i)]}. \end{aligned}$$

For $h \geq 2$, we observe the following one-step-back decomposition:

$$\begin{aligned} & \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}_i) \sim d_{\hat{P}_h^{(n)}, h}^\pi} [g_h(s[Z_i], \mathbf{a}_i)] \\ &= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}_h^{(n)}, h-1}^\pi, s \sim \hat{P}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a}_i \sim \pi_h(s)} [g_h(s[Z_i], \mathbf{a}_i)] \\ &= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}_h^{(n)}, h-1}^\pi} \left[\int_{\mathcal{S}} \prod_{j=1}^M \left[\hat{\phi}_{h-1, j}^{(n)}(\tilde{s}[Z_j], \tilde{\mathbf{a}}_j)^\top \hat{w}_{h-1, j}^{(n)}(s_j) \right] \sum_{\mathbf{a}_i \in \mathcal{A}_i} \pi_h(\mathbf{a}_i | s) g_h(s[Z_i], \mathbf{a}_i) ds \right] \\ &= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}_h^{(n)}, h-1}^\pi} \left[\min \left\{ \int_{\mathcal{S}} \prod_{j=1}^M \left[\hat{\phi}_{h-1, j}^{(n)}(\tilde{s}[Z_j], \tilde{\mathbf{a}}_j)^\top \hat{w}_{h-1, j}^{(n)}(s_j) \right] \sum_{\mathbf{a}_i \in \mathcal{A}_i} \pi_h(\mathbf{a}_i | s) g_h(s[Z_i], \mathbf{a}_i) ds, B \right\} \right] \end{aligned}$$

$$\begin{aligned}
&\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^\pi} \left[\min \left\{ \tilde{A} \int_{\mathcal{S}} \prod_{j=1}^M [\hat{\phi}_{h-1,j}^{(n)}(\tilde{s}[Z_j], \tilde{\mathbf{a}}_j)^\top \hat{w}_{h-1,j}^{(n)}(s_j)] \frac{1}{|\mathcal{A}_i|} \sum_{\mathbf{a}_i \in \mathcal{A}_i} g_h(s[Z_i], \mathbf{a}_i) ds, B \right\} \right] \\
&= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^\pi} \left[\min \left\{ \tilde{A} \int_{\mathcal{S}[Z_i]} \prod_{j \in Z_i} [\hat{\phi}_{h-1,j}^{(n)}(\tilde{s}[Z_j], \tilde{\mathbf{a}}_j)^\top \hat{w}_{h-1,j}^{(n)}(s_j)] \frac{1}{|\mathcal{A}_i|} \sum_{\mathbf{a}_i \in \mathcal{A}_i} g_h(s[Z_i], \mathbf{a}_i) ds[Z_i], B \right\} \right] \\
&= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^\pi} \left[\min \left\{ \tilde{A} \int_{\mathcal{S}[Z_i]} \left[\bigotimes_{j \in Z_i} \hat{\phi}_{h-1,j}^{(n)}(\tilde{s}[Z_j], \tilde{\mathbf{a}}_j) \right]^\top \left[\bigotimes_{j \in Z_i} \hat{w}_{h-1,j}^{(n)}(s_j) \right] \frac{1}{|\mathcal{A}_i|} \sum_{\mathbf{a}_i \in \mathcal{A}_i} g_h(s[Z_i], \mathbf{a}_i) ds[Z_i], B \right\} \right] \\
&= \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^\pi} \left[\min \left\{ \tilde{A} \int_{\mathcal{S}[Z_i]} \bar{\phi}_{h-1,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}})^\top \left[\bigotimes_{j \in Z_i} \hat{w}_{h-1,j}^{(n)}(s_j) \right] \frac{1}{|\mathcal{A}_i|} \sum_{\mathbf{a}_i \in \mathcal{A}_i} g_h(s[Z_i], \mathbf{a}_i) ds[Z_i], B \right\} \right] \\
&\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^\pi} \left[\min \left\{ \tilde{A} \left\| \bar{\phi}_{h-1,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \bar{\phi}_{h-1,i}^{(n)}}^{-1}} \left\| \int_{\mathcal{S}[Z_i]} \frac{1}{|\mathcal{A}_i|} \sum_{\mathbf{a}_i \in \mathcal{A}_i} \left(\bigotimes_{j \in Z_i} \hat{w}_{h-1,j}^{(n)}(s_j) \right) g_h(s[Z_i], \mathbf{a}_i) ds[Z_i] \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \bar{\phi}_{h-1,i}^{(n)}}}, B \right\} \right].
\end{aligned}$$

Then,

$$\begin{aligned}
&\left\| \int_{\mathcal{S}[Z_i]} \frac{1}{|\mathcal{A}_i|} \sum_{\mathbf{a}_i \in \mathcal{A}_i} \left(\bigotimes_{j \in Z_i} \hat{w}_{h-1,j}^{(n)}(s_j) \right) g_h(s[Z_i], \mathbf{a}_i) ds[Z_i] \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \bar{\phi}_{h-1,i}^{(n)}}}^2 \\
&\leq n \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \rho_{h-1}^{(n)}} \left[\left(\int_{\mathcal{S}[Z_i]} \frac{1}{|\mathcal{A}_i|} \sum_{\mathbf{a}_i \in \mathcal{A}_i} \prod_{j \in Z_i} \left(\hat{w}_{h-1,j}^{(n)}(s_j)^\top \hat{\phi}_{h-1,j}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}_j) \right) g_h(s[Z_i], \mathbf{a}_i) ds[Z_i] \right)^2 \right] + B^2 \lambda d^L \\
&\quad \left(\left\| \frac{1}{|\mathcal{A}_i|} \sum_{\mathbf{a}_i \in \mathcal{A}_i} g_h(s[Z_i], \mathbf{a}_i) \right\|_\infty \leq B \text{ and } \left\| \hat{w}_{h-1,i}^{(n)}(s_i) \right\|_2 \leq \sqrt{d}. \right) \\
&= n \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \rho_{h-1}^{(n)}} \left[\left(\mathbb{E}_{s \sim \hat{P}_{h-1}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a}_i \sim U(\mathcal{A}_i)} [g_h(s[Z_i], \mathbf{a}_i)] \right)^2 \right] + B^2 \lambda d^L \\
&\leq n \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \rho_{h-1}^{(n)}} \left[\left(\mathbb{E}_{s \sim P_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a}_i \sim U(\mathcal{A}_i)} [g_h(s[Z_i], \mathbf{a}_i)] \right)^2 \right] + B^2 \lambda d^L + n B^2 \xi^{(n)} \quad (\text{Event } \mathcal{E}) \\
&\leq n \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim \rho_{h-1}^{(n)}, s \sim P_{h-1}^*(\tilde{s}, \tilde{\mathbf{a}}), \mathbf{a}_i \sim U(\mathcal{A}_i)} [g_h^2(s[Z_i], \mathbf{a}_i)] + B^2 \lambda d^L + B^2 n \xi^{(n)}. \quad (\text{Jensen}) \\
&= n \mathbb{E}_{(s, \mathbf{a}_i) \sim \hat{\rho}_h^{(n)}} [g_h^2(s[Z_i], \mathbf{a}_i)] + B^2 \lambda d^L + B^2 n \zeta^{(n)}. \quad (\text{Definition of } \hat{\rho}_h^{(n)})
\end{aligned}$$

Combing the above results together, we get

$$\begin{aligned}
&\mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}_i) \sim d_{\hat{P}^{(n)}, h}^\pi} [g_h(s[Z_i], \mathbf{a}_i)] \\
&\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^\pi} \left[\min \left\{ \tilde{A} \left\| \bar{\phi}_{h-1,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \bar{\phi}_{h-1,i}^{(n)}}^{-1}} \left\| \int_{\mathcal{S}[Z_i]} \frac{1}{|\mathcal{A}_i|} \sum_{\mathbf{a}_i \in \mathcal{A}_i} \left(\bigotimes_{j \in Z_i} \hat{w}_{h-1,j}^{(n)}(s_j) \right) g_h(s[Z_i], \mathbf{a}_i) ds[Z_i] \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \bar{\phi}_{h-1,i}^{(n)}}}, B \right\} \right] \\
&\leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^\pi} \left[\min \left\{ \tilde{A} \left\| \bar{\phi}_{h-1,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \bar{\phi}_{h-1,i}^{(n)}}^{-1}} \sqrt{n \mathbb{E}_{(s, \mathbf{a}_i) \sim \hat{\rho}_h^{(n)}} [g_h^2(s[Z_i], \mathbf{a}_i)] + B^2 \lambda d^L + B^2 n \zeta^{(n)}}, B \right\} \right],
\end{aligned}$$

which has finished the proof. \square

Lemma D.3 (One-step back inequality for the true model). *Consider a set of functions $\{g_h\}_{h=1}^H$ that satisfies $g_h \in \mathcal{S}[Z_i] \times \mathcal{A}_i \rightarrow \mathbb{R}_+$, s.t. $\|g_h\|_\infty \leq B$. Then for any policy π , we have*

$$\left| \mathbb{E}_{(s, \mathbf{a}_i) \sim d_{P^*, h}^\pi} [g_h(s[Z_i], \mathbf{a}_i)] \right|$$

$$\leq \begin{cases} \sqrt{\tilde{A}\mathbb{E}_{(s,\mathbf{a}_i)\sim\rho_1^{(n)}}[g_1^2(s[Z_i],\mathbf{a}_i)]}, & h=1, \\ \tilde{A}\mathbb{E}_{(\tilde{s},\tilde{\mathbf{a}})\sim d_{P^*,h-1}^\pi} \left[\left\| \bar{\phi}_{h-1,i}^*(\tilde{s},\tilde{\mathbf{a}}) \right\|_{\Sigma_{n,\gamma_{h-1}^{(n)},\bar{\phi}_{h-1,i}^*}^{-1}} \right] \sqrt{n\mathbb{E}_{(s,\mathbf{a})\sim\rho_h^{(n)}}[g_h^2(s[Z_i],\mathbf{a}_i)] + B^2\lambda d^L}, & h \geq 2, \end{cases}$$

where $\bar{\phi}_{h,i}^*(s,\mathbf{a}) := \bigotimes_{j \in Z_i} \phi_{h-1,j}^*(s[Z_j],\mathbf{a}_j)$, and $\Sigma_{n,\gamma_h^{(n)},\bar{\phi}_{h,i}^*} = n\mathbb{E}_{(s,\mathbf{a})\sim\gamma_h^{(n)}} \left[\bar{\phi}_{h,i}^*(s,\mathbf{a})\bar{\phi}_{h,i}^*(s,\mathbf{a})^\top \right] + \lambda I_{d|Z_i|}$.

Proof. For step $h=1$, we have

$$\begin{aligned} \mathbb{E}_{(s,\mathbf{a})\sim d_{P^*,1}^\pi} [g_1(s[Z_i],\mathbf{a}_i)] &= \mathbb{E}_{s \sim d_1, \mathbf{a}_i \sim \pi_1(s)} [g_1(s[Z_i],\mathbf{a}_i)] \\ &\leq \sqrt{\max_{(s,\mathbf{a}_i)} \frac{d_1(s)\pi_1(\mathbf{a}_i|s)}{\rho_1^{(n)}(s,\mathbf{a}_i)} \mathbb{E}_{(s',\mathbf{a}'_i)\sim\rho_1^{(n)}} [g_1^2(s'[Z_i],\mathbf{a}'_i)]} \\ &= \sqrt{\max_{(s,\mathbf{a}_i)} \frac{d_1(s)\pi_1(\mathbf{a}_i|s)}{d_1(s)u_{\mathcal{A}_i}(\mathbf{a}_i)} \mathbb{E}_{(s',\mathbf{a}'_i)\sim\rho_1^{(n)}} [g_1^2(s'[Z_i],\mathbf{a}'_i)]} \\ &\leq \sqrt{\tilde{A}\mathbb{E}_{(s,\mathbf{a}_i)\sim\rho_1^{(n)}} [g_1^2(s[Z_i],\mathbf{a}_i)]}. \end{aligned}$$

For step $h=2, \dots, H-1$, we observe the following one-step-back decomposition:

$$\begin{aligned} &\mathbb{E}_{(\tilde{s},\tilde{\mathbf{a}}_i)\sim d_{P^*,h}^\pi} [g_h(s[Z_i],\mathbf{a}_i)] \\ &= \mathbb{E}_{(\tilde{s},\tilde{\mathbf{a}})\sim d_{P^*,h-1}^\pi, s \sim P_{h-1}^*(\tilde{s},\tilde{\mathbf{a}}), \mathbf{a}_i \sim \pi_h(s)} [g_h(s[Z_i],\mathbf{a}_i)] \\ &= \mathbb{E}_{(\tilde{s},\tilde{\mathbf{a}})\sim d_{P^*,h-1}^\pi} \left[\left(\bigotimes_{j \in Z_i} \phi_{h-1,j}^*(\tilde{s}[Z_j],\tilde{\mathbf{a}}_j) \right)^\top \int_{\mathcal{S}} \sum_{\mathbf{a}_i \in \mathcal{A}_i} \left(\bigotimes_{j \in Z_i} w_{h-1,j}^*(s_j) \right) \pi_h(\mathbf{a}_i|s) g_h(s[Z_i],\mathbf{a}_i) ds \right] \\ &\leq \tilde{A}\mathbb{E}_{(\tilde{s},\tilde{\mathbf{a}})\sim d_{P^*,h-1}^\pi} \left[\left(\bigotimes_{j \in Z_i} \phi_{h-1,j}^*(\tilde{s}[Z_j],\tilde{\mathbf{a}}_j) \right)^\top \int_{\mathcal{S}} \sum_{\mathbf{a}_i \in \mathcal{A}_i} \frac{1}{|\mathcal{A}_i|} \left(\bigotimes_{j \in Z_i} w_{h-1,j}^*(s_j) \right) g_h(s[Z_i],\mathbf{a}_i) ds [Z_i] \right] \\ &\leq \tilde{A}\mathbb{E}_{(\tilde{s},\tilde{\mathbf{a}})\sim d_{P^*,h-1}^\pi} \left[\left\| \bigotimes_{j \in Z_i} \phi_{h-1,j}^*(\tilde{s}[Z_j],\tilde{\mathbf{a}}_j) \right\|_{\Sigma_{n,\gamma_{h-1}^{(n)},\bar{\phi}_{h-1,i}^*}^{-1}} \right] \\ &\quad \cdot \left\| \int_{\mathcal{S}} \sum_{\mathbf{a}_i \in \mathcal{A}_i} \frac{1}{|\mathcal{A}_i|} \left(\bigotimes_{j \in Z_i} w_{h-1,j}^*(s_j) \right) g_h(s[Z_i],\mathbf{a}_i) ds [Z_i] \right\|_{\Sigma_{n,\gamma_{h-1}^{(n)},\bar{\phi}_{h-1,i}^*}}. \end{aligned}$$

Then,

$$\begin{aligned} &\left\| \int_{\mathcal{S}} \sum_{\mathbf{a}_i \in \mathcal{A}_i} \frac{1}{|\mathcal{A}_i|} \left(\bigotimes_{j \in Z_i} w_{h-1,j}^*(s_j) \right) g_h(s[Z_i],\mathbf{a}_i) ds [Z_i] \right\|_{\Sigma_{n,\gamma_{h-1}^{(n)},\bar{\phi}_{h-1,i}^*}}^2 \\ &\leq n\mathbb{E}_{(\tilde{s},\tilde{\mathbf{a}})\sim\gamma_{h-1}^{(n)}} \left[\left(\int_{\mathcal{S}} \sum_{\mathbf{a}_i \in \mathcal{A}_i} \frac{1}{|\mathcal{A}_i|} \left(\bigotimes_{j \in Z_i} w_{h-1,j}^*(s_j) \right) \left(\bigotimes_{j \in Z_i} \phi_{h-1,j}^*(\tilde{s}[Z_j],\tilde{\mathbf{a}}_j) \right) g_h(s[Z_i],\mathbf{a}_i) ds [Z_i] \right)^2 \right] + B^2\lambda d^L \\ &\quad \text{(Use the assumption } \left\| \sum_{\mathbf{a}_i \in \mathcal{A}_i} \frac{1}{|\mathcal{A}_i|} g_h(s[Z_i],\mathbf{a}_i) \right\|_\infty \leq B \text{ and } \|w_{h-1,i}^*(s_i)\|_2 \leq \sqrt{d}.) \\ &= n\mathbb{E}_{(\tilde{s},\tilde{\mathbf{a}})\sim\gamma_{h-1}^{(n)}} \left[\left(\mathbb{E}_{s \sim P_{h-1}^*(\tilde{s},\tilde{\mathbf{a}}), \mathbf{a}_i \sim U(\mathcal{A}_i)} [g_h(s[Z_i],\mathbf{a}_i)] \right)^2 \right] + B^2\lambda d^L \\ &\leq n\mathbb{E}_{(\tilde{s},\tilde{\mathbf{a}})\sim\gamma_{h-1}^{(n)}, s \sim P_{h-1}^*(\tilde{s},\tilde{\mathbf{a}}), \mathbf{a}_i \sim U(\mathcal{A}_i)} [g_h^2(s[Z_i],\mathbf{a}_i)] + B^2\lambda d^L \quad \text{(Jensen)} \end{aligned}$$

$$\leq n\mathbb{E}_{(s,\mathbf{a}_i)\sim\rho_h^{(n)}} [g_h^2(s[Z_i], \mathbf{a}_i)] + B^2\lambda d^L, \quad (\text{Definition of } \rho_h^{(n)})$$

which has finished the proof. \square

Lemma D.4 (One-step back inequality for the true model). *Consider a set of functions $\{g_h\}_{h=1}^H$ that satisfies $g_h \in \mathcal{S}[\cup_{j \in Z_i} Z_j] \times \mathcal{A}[Z_i] \rightarrow \mathbb{R}$, s.t. $\|g_h\|_\infty \leq B$. Then for any policy π , we have*

$$\begin{aligned} & \left| \mathbb{E}_{(s,\mathbf{a})\sim d_{P^*,h}^\pi} [g_h(s[\cup_{j \in Z_i} Z_j], \mathbf{a}[Z_i])] \right| \\ & \leq \begin{cases} \sqrt{\tilde{A}^L \mathbb{E}_{(s,\mathbf{a})\sim\rho_1^{(n)}} [g_1^2(s[\cup_{j \in Z_i} Z_j], \mathbf{a}[Z_i])]}, & h = 1, \\ \tilde{A}^L \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}})\sim d_{P^*,h-1}^\pi} \left[\left\| \tilde{\phi}_{h-1,i}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n,\gamma_{h-1}^{(n)}, \tilde{\phi}_{h-1,i}^*}^{-1}} \right] \sqrt{n\mathbb{E}_{(s,\mathbf{a})\sim\rho_h^{(n)}} [g_h^2(s[\cup_{j \in Z_i} Z_j], \mathbf{a}[Z_i])] + B^2\lambda d^L}, & h \geq 2, \end{cases} \end{aligned}$$

$$\text{where } \tilde{\phi}_{h,i}^*(s, \mathbf{a}) := \bigotimes_{k \in \cup_{j \in Z_i} Z_j} \phi_{h-1,j}^*(s[Z_k], \mathbf{a}_k), \quad \text{and } \Sigma_{n,\gamma_h^{(n)}, \tilde{\phi}_{h,i}^*} = n\mathbb{E}_{(s,\mathbf{a})\sim\gamma_h^{(n)}} [\tilde{\phi}_{h,i}^*(s, \mathbf{a}) \tilde{\phi}_{h,i}^*(s, \mathbf{a})^\top] + \lambda I_{d^{\cup_{j \in Z_i} Z_j}}.$$

Proof. This Lemma can be proved using similar steps as those in the proof of Lemma D.3, noting that in this case the dimension of $\tilde{\phi}_{h,i}^*$ is at most L^2 . \square

Lemma D.5 (Optimism for NE and CCE). *Consider an episode $n \in [N]$ and set $\alpha^{(n)} = \Theta\left(H\tilde{A}\sqrt{n\zeta^{(n)} + d^L\lambda}\right)$. When the event \mathcal{E} holds and the policy $\pi^{(n)}$ is computed by solving NE or CCE, we have*

$$\bar{v}_i^{(n)}(s) - v_i^{\dagger, \pi^{(n)}}(s) \geq -HM\sqrt{\tilde{A}\zeta^{(n)}}, \quad \forall n \in [N], i \in [M].$$

Proof. Denote $\tilde{\mu}_{h,i}^{(n)}(\cdot|s) := \arg \max_{\mu} \left(\mathbb{D}_{\mu, \pi_{h,-i}^{(n)}} Q_{h,i}^{\dagger, \pi^{(n)}} \right)(s)$ and let $\tilde{\pi}_h^{(n)} = \tilde{\mu}_{h,i}^{(n)} \times \pi_{h,-i}^{(n)}$. Let $f_h^{(n)}(s, \mathbf{a}) = \left\| \hat{P}_h^{(n)}(\cdot|s, \mathbf{a}) - P_h^*(\cdot|s, \mathbf{a}) \right\|_1$ and $f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) = \left\| \hat{P}_{h,i}^{(n)}(\cdot|s[Z_i], \mathbf{a}_i) - P_{h,i}^*(\cdot|s[Z_i], \mathbf{a}_i) \right\|_1$. Then according to the event \mathcal{E} , we have

$$\begin{aligned} \mathbb{E}_{(s,\mathbf{a})\sim\rho_h^{(n)}} \left[\left(f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) \right)^2 \right] &\leq \zeta^{(n)}, \quad \mathbb{E}_{(s,\mathbf{a})\sim\hat{\rho}_h^{(n)}} \left[\left(f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) \right)^2 \right] \leq \zeta^{(n)}, \quad \forall n \in [N], h \in [H], i \in [M] \\ \left\| \bar{\phi}_{h,i}(s, \mathbf{a}) \right\|_{\left(\hat{\Sigma}_{h, \bar{\phi}_{h,i}}^{(n)} \right)^{-1}} &= \Theta \left(\left\| \bar{\phi}_{h,i}(s, \mathbf{a}) \right\|_{\Sigma_{n,\rho_h^{(n)}, \bar{\phi}_{h,i}}^{-1}} \right), \quad \forall n \in [N], h \in [H], \bar{\phi}_{h,i} \in \bar{\Phi}_{h,i}, i \in [M]. \end{aligned}$$

A direct conclusion of the event \mathcal{E} is we can find an absolute constant c , such that

$$\begin{aligned} \beta_h^{(n)}(s, \mathbf{a}) &= \sum_{i=1}^M \min \left\{ \alpha^{(n)} \left\| \bar{\phi}_{h,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\left(\Sigma_{h, \bar{\phi}_{h,i}}^{(n)} \right)^{-1}}, H \right\} \\ &\geq \sum_{i=1}^M \min \left\{ c\alpha^{(n)} \left\| \bar{\phi}_{h,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n,\rho_h^{(n)}, \bar{\phi}_{h,i}}^{-1}}, H \right\}, \quad \forall n \in [N], h \in [H]. \end{aligned}$$

Next, similar to the proof in Lemma B.5, we may prove

$$\mathbb{E}_{s \sim d_1} \left[\bar{V}_{1,i}^{(n)}(s) - V_{1,i}^{\dagger, \pi^{(n)}}(s) \right] \geq \sum_{h=1}^H \mathbb{E}_{(s,\mathbf{a})\sim d_{\hat{P}^{(n)},h}^\pi} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] - H \sum_{h=1}^H \mathbb{E}_{(s,\mathbf{a})\sim d_{\hat{P}^{(n)},h}^\pi} \left[\min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right].$$

For the second term, note that we have the relation $\min\{f_h^{(n)}(s, \mathbf{a}), 1\} \leq \sum_{i=1}^M \min\{f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i), 1\}$. By Lemma D.2, we have for $h = 1$,

$$\mathbb{E}_{(s,\mathbf{a})\sim d_{\hat{P}^{(n)},1}^\pi} \left[\min \left\{ f_{1,i}^{(n)}(s[Z_i], \mathbf{a}_i), 1 \right\} \right] \leq \sqrt{A\mathbb{E}_{(s,\mathbf{a})\sim\rho_1^{(n)}} \left[\left(f_{1,i}^{(n)}(s[Z_i], \mathbf{a}_i) \right)^2 \right]} \leq \sqrt{\tilde{A}\zeta^{(n)}}.$$

And $\forall h \geq 2$, we have

$$\begin{aligned} & \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\min \left\{ f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i), 1 \right\} \right] \\ & \lesssim \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\pi^{(n)}}} \left[\min \left\{ \tilde{A} \left\| \bar{\phi}_{h-1,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}_{n, \rho_{h-1}^{(n)}, \bar{\phi}_{h-1,i}^{(n)}}} \sqrt{n \mathbb{E}_{(s, \mathbf{a}) \sim \bar{\rho}_h^{(n)}} \left[\left(f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) \right)^2 \right] + d^L \lambda + n \zeta^{(n)}}, 1 \right\} \right] \\ & \lesssim \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\pi^{(n)}}} \left[\min \left\{ \tilde{A} \left\| \bar{\phi}_{h-1,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}_{n, \rho_{h-1}^{(n)}, \bar{\phi}_{h-1,i}^{(n)}}} \sqrt{n \zeta^{(n)} + d^L \lambda}, 1 \right\} \right] \end{aligned}$$

Note that we here use $\min \{ f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i), 1 \} \leq 1$, $\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) \right)^2 \right] \leq \zeta^{(n)}$ and

$\mathbb{E}_{(s, \mathbf{a}) \sim \bar{\rho}_h^{(n)}} \left[\left(f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) \right)^2 \right] \leq \zeta^{(n)}$. Then according to our choice of $\alpha^{(n)}$, we get

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\min \left\{ f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i), 1 \right\} \right] \leq \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\pi^{(n)}}} \left[\min \left\{ \frac{c\alpha^{(n)}}{H} \left\| \bar{\phi}_{h-1,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma^{-1}_{n, \rho_{h-1}^{(n)}, \bar{\phi}_{h-1,i}^{(n)}}}, 1 \right\} \right].$$

Combining all things together,

$$\begin{aligned} \bar{v}_i^{(n)} - v_i^{\dagger, \pi^{(n)}} &= \mathbb{E}_{s \sim d_1} \left[\bar{V}_{1,i}^{(n)}(s) - V_{1,i}^{\dagger, \pi^{(n)}}(s) \right] \\ & \geq \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] - H \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right] \\ & \geq \sum_{h=1}^{H-1} \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - \sum_{j=1}^M \min \left\{ c\alpha^{(n)} \left\| \bar{\phi}_{h,j}^{(n)}(s, \mathbf{a}) \right\|_{\Sigma^{-1}_{\rho_h^{(n)}, \bar{\phi}_{h,j}^{(n)}}}, H \right\} \right] - HM \sqrt{\tilde{A} \zeta^{(n)}} \\ & \geq -HM \sqrt{\tilde{A} \zeta^{(n)}}, \end{aligned}$$

which proves the inequality. \square

Lemma D.6 (Optimism for CE). *Consider an episode $n \in [N]$ and set $\alpha^{(n)} = \Theta \left(H \tilde{A} \sqrt{n \zeta^{(n)} + d^L \lambda} \right)$. When the event \mathcal{E} holds, we have*

$$\bar{v}_i^{(n)}(s) - \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi^{(n)}}(s) \geq -HM \sqrt{A \zeta^{(n)}}, \quad \forall n \in [N], i \in [M].$$

Proof. Denote $\tilde{\omega}_{h,i}^{(n)} = \arg \max_{\omega_h \in \Omega_{h,i}} \left(\mathbb{D}_{\omega_h \circ \pi_h^{(n)}} \max_{\omega \in \Omega_i} Q_{h,i}^{\omega \circ \pi^{(n)}} \right)(s)$ and let $\tilde{\pi}_h^{(n)} = \tilde{\omega}_{h,i} \circ \pi_h^{(n)}$. Let $f_h^{(n)}(s, \mathbf{a}) = \left\| \hat{P}_h^{(n)}(\cdot | s, \mathbf{a}) - P_h^*(\cdot | s, \mathbf{a}) \right\|_1$ and $f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) = \left\| \hat{P}_{h,i}^{(n)}(\cdot | s[Z_i], \mathbf{a}_i) - P_{h,i}^*(\cdot | s[Z_i], \mathbf{a}_i) \right\|_1$. Then according to the event \mathcal{E} , we have

$$\begin{aligned} \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) \right)^2 \right] &\leq \zeta^{(n)}, \quad \mathbb{E}_{(s, \mathbf{a}) \sim \bar{\rho}_h^{(n)}} \left[\left(f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) \right)^2 \right] \leq \zeta^{(n)}, \quad \forall n \in [N], h \in [H], i \in [M] \\ \left\| \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) \right\|_{\left(\hat{\Sigma}_{h, \bar{\phi}_{h,i}^{(n)}}^{(n)} \right)^{-1}} &= \Theta \left(\left\| \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) \right\|_{\Sigma^{-1}_{n, \rho_h^{(n)}, \bar{\phi}_{h,i}^{(n)}}} \right), \quad \forall n \in [N], h \in [H], \bar{\phi}_{h,i} \in \bar{\Phi}_{h,i}, i \in [M]. \end{aligned}$$

A direct conclusion of the event \mathcal{E} is we can find an absolute constant c , such that

$$\beta_h^{(n)}(s, \mathbf{a}) = \min \left\{ \alpha^{(n)} \sum_{i=1}^M \left\| \bar{\phi}_{h,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\left(\hat{\Sigma}_{h, \bar{\phi}_{h,i}^{(n)}}^{(n)} \right)^{-1}}, H \right\}$$

$$\geq c \min \left\{ \alpha^{(n)} \sum_{i=1}^M \left\| \bar{\phi}_{h,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \bar{\phi}_{h,i}^{(n)}}^{-1}}, H \right\}, \quad \forall n \in [N], h \in [H].$$

Next, similar to the proof in Lemma B.6, we may prove

$$\mathbb{E}_{s \sim d_1} \left[\bar{V}_{1,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{1,i}^{\omega \circ \pi^{(n)}}(s) \right] \geq \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\bar{\pi}^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] - H \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\bar{\pi}^{(n)}}} \left[\min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right].$$

Note that we can use exactly the same steps in the proof of Lemma D.5 to bound the second term, and we get for $h = 1$,

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, 1}^{\bar{\pi}^{(n)}}} \left[\min \left\{ f_{1,i}^{(n)}(s[Z_i], \mathbf{a}_i), 1 \right\} \right] \leq \sqrt{\tilde{A}\zeta^{(n)}}.$$

And $\forall h \geq 2$,

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\bar{\pi}^{(n)}}} \left[\min \left\{ f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i), 1 \right\} \right] \leq \frac{c\alpha^{(n)}}{H} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\bar{\pi}^{(n)}}} \left[\left\| \bar{\phi}_{h-1,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \bar{\phi}_{h-1,i}^{(n)}}^{-1}} \right].$$

Combining all things together,

$$\begin{aligned} & \bar{v}_i^{(n)} - \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi^{(n)}} \\ &= \mathbb{E}_{s \sim d_1} \left[\bar{V}_{1,i}^{(n)}(s) - \max_{\omega \in \Omega_i} V_{1,i}^{\omega \circ \pi^{(n)}}(s) \right] \\ &\geq \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\bar{\pi}^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] - H \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\bar{\pi}^{(n)}}} \left[\min \left\{ f_h^{(n)}(s, \mathbf{a}), 1 \right\} \right] \\ &\geq \sum_{h=1}^{H-1} \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\bar{\pi}^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) - \sum_{j=1}^M \min \left(c\alpha^{(n)} \left\| \bar{\phi}_{h,j}^{(n)}(s, \mathbf{a}) \right\|_{\Sigma_{\rho_h^{(n)}, \bar{\phi}_{h,j}^{(n)}}^{-1}}, H \right) \right] - HM\sqrt{\tilde{A}\zeta^{(n)}} \\ &\geq -HM\sqrt{\tilde{A}\zeta^{(n)}}, \end{aligned}$$

which proves the inequality. \square

Lemma D.7 (pessimism). *Consider an episode $n \in [N]$ and set $\alpha^{(n)} = \Theta \left(H\tilde{A}\sqrt{n\zeta^{(n)}} + d^L\lambda \right)$. When the event \mathcal{E} holds, we have*

$$\underline{v}_i^{(n)}(s) - v_i^{\pi^{(n)}}(s) \leq HM\sqrt{\tilde{A}\zeta^{(n)}}, \quad \forall n \in [N], i \in [M].$$

Proof. Let $f_h^{(n)}(s, \mathbf{a}) = \left\| \hat{P}_h^{(n)}(\cdot | s, \mathbf{a}) - P_h^*(\cdot | s, \mathbf{a}) \right\|_1$ and $f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) = \left\| \hat{P}_{h,i}^{(n)}(\cdot | s[Z_i], \mathbf{a}_i) - P_{h,i}^*(\cdot | s[Z_i], \mathbf{a}_i) \right\|_1$. Then according to the event \mathcal{E} , we have

$$\begin{aligned} & \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) \right)^2 \right] \leq \zeta^{(n)}, \quad \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_h^{(n)}(s, \mathbf{a}) \right)^2 \right] \leq \zeta^{(n)}, \quad \forall n \in [N], h \in [H], i \in [M] \\ & \left\| \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) \right\|_{\left(\hat{\Sigma}_{h, \bar{\phi}_{h,i}^{(n)}}^{(n)} \right)^{-1}} = \Theta \left(\left\| \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \bar{\phi}_{h,i}^{(n)}}^{-1}} \right), \quad \forall n \in [N], h \in [H], \bar{\phi}_{h,i} \in \bar{\Phi}_{h,i}, i \in [M]. \end{aligned}$$

A direct conclusion of the event \mathcal{E} is we can find an absolute constant c , such that

$$\beta_h^{(n)}(s, \mathbf{a}) = \sum_{i=1}^M \min \left\{ \alpha^{(n)} \left\| \bar{\phi}_{h,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\left(\hat{\Sigma}_{h, \bar{\phi}_{h,i}^{(n)}}^{(n)} \right)^{-1}}, H \right\}$$

$$\geq \sum_{i=1}^M \min \left\{ c\alpha^{(n)} \left\| \bar{\phi}_{h,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \bar{\phi}_{h,i}^{(n)}}^{-1}}, H \right\}, \quad \forall n \in [N], h \in [H].$$

Next, similar to the proof in Lemma B.7, we may prove

$$\mathbb{E}_{s \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[V_{h,i}^{(n)}(s) - V_{h,i}^{\pi^{(n)}}(s) \right] \leq \sum_{h'=h}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h'}^{\pi^{(n)}}} \left[-\hat{\beta}_{h'}^{(n)}(s, \mathbf{a}) + H \min \left\{ f_{h'}^{(n)}(s, \mathbf{a}), 1 \right\} \right], \quad \forall h \in [H]. \quad (22)$$

and we get for $h = 1$,

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, 1}^{\pi^{(n)}}} \left[\min \left\{ f_{1,i}^{(n)}(s[Z_i], \mathbf{a}_i), 1 \right\} \right] \leq \sqrt{\tilde{A}\zeta^{(n)}}.$$

And $\forall h \geq 2$,

$$\mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\min \left\{ f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i), 1 \right\} \right] \leq \frac{c\alpha^{(n)}}{H} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{\hat{P}^{(n)}, h-1}^{\pi^{(n)}}} \left[\left\| \bar{\phi}_{h-1,i}^{(n)}(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \rho_{h-1}^{(n)}, \bar{\phi}_{h-1,i}^{(n)}}^{-1}} \right].$$

Finally, we get

$$\begin{aligned} \underline{v}_i^{(n)} - v_i^{\pi^{(n)}} &= \mathbb{E}_{s \sim d_1} \left[\underline{V}_{1,i}^{(n)}(s) - V_{1,i}^{\pi^{(n)}}(s) \right] \\ &\leq \sum_{h=1}^{H-1} \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[-\hat{\beta}_h^{(n)}(s, \mathbf{a}) + \sum_{j=1}^M \min \left(c\alpha^{(n)} \left\| \bar{\phi}_{h,j}^{(n)}(s, \mathbf{a}) \right\|_{\Sigma_{\rho_h^{(n)}, \bar{\phi}_{h,j}^{(n)}}^{-1}}, H \right) \right] + HM\sqrt{\tilde{A}\zeta^{(n)}} \\ &\leq HM\sqrt{\tilde{A}\zeta^{(n)}}, \end{aligned}$$

which has finished the proof. \square

Lemma D.8. *When the event \mathcal{E} holds and $\alpha^{(n)} = \Theta \left(H\tilde{A}\sqrt{n\zeta^{(n)} + d^L\lambda} \right)$ satisfies $\alpha^{(1)} \leq \alpha^{(2)} \leq \dots \leq \alpha^{(N)}$, we have*

$$\sum_{n=1}^N \Delta^{(n)} \lesssim H^2 M d^{L^2} A^L \sqrt{N \log \left(1 + \frac{N}{d\lambda} \right)} \alpha^{(N)}.$$

Proof. Let $f_h^{(n)}(s, \mathbf{a}) = \left\| \hat{P}_h^{(n)}(\cdot | s, \mathbf{a}) - P_h^*(\cdot | s, \mathbf{a}) \right\|_1$ and $f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) = \left\| \hat{P}_{h,i}^{(n)}(\cdot | s[Z_i], \mathbf{a}_i) - P_{h,i}^*(\cdot | s[Z_i], \mathbf{a}_i) \right\|_1$. Then according to the event \mathcal{E} , we have

$$\begin{aligned} \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) \right)^2 \right] &\leq \zeta^{(n)}, \quad \mathbb{E}_{(s, \mathbf{a}) \sim \hat{\rho}_h^{(n)}} \left[\left(f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) \right)^2 \right] \leq \zeta^{(n)}, \quad \forall n \in [N], h \in [H], i \in [M] \\ \left\| \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) \right\|_{\left(\hat{\Sigma}_{h, \bar{\phi}_{h,i}^{(n)}}^{(n)} \right)^{-1}} &= \Theta \left(\left\| \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \bar{\phi}_{h,i}^{(n)}}^{-1}} \right), \quad \forall n \in [N], h \in [H], \bar{\phi}_{h,i}^{(n)} \in \bar{\Phi}_{h,i}, i \in [M]. \end{aligned}$$

By definition, we have

$$\Delta^{(n)} = \max_{i \in [M]} \left\{ \bar{v}_i^{(n)} - \underline{v}_i^{(n)} \right\} + 2HM\sqrt{\tilde{A}\zeta^{(n)}}.$$

With similar steps as those in the proof of Lemma B.8 (note that $\bar{V}_{h,i}^{(n)}(s) - \underline{V}_{h,i}^{(n)}(s)$ is upper bounded by $2H^2M$), we have

$$\mathbb{E}_{s \sim d_{\hat{P}^{(n)}, 1}^{\pi^{(n)}}} \left[\bar{V}_{1,i}^{(n)}(s) - \underline{V}_{1,i}^{(n)}(s) \right] \leq \underbrace{2 \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right]}_{(a)} + \underbrace{2H^2M \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{\hat{P}^{(n)}, h}^{\pi^{(n)}}} \left[f_h^{(n)}(s, \mathbf{a}) \right]}_{(b)}. \quad (23)$$

First, we calculate the first term (a) in Inequality equation 23. Following Lemma D.4, we have

$$\begin{aligned}
& \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] \\
& \lesssim \sum_{h=1}^H \sum_{i=1}^M \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[\min \left(\alpha^{(n)} \left\| \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \bar{\phi}_{h,i}^{(n)}}^{-1}}, H \right) \right] \\
& \lesssim \sum_{h=1}^{H-1} \sum_{i=1}^M \tilde{A}^L \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[\left\| \tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_h^{(n)}, \tilde{\phi}_{h,i}^*}^{-1}} \right] \\
& \quad \cdot \sqrt{n (\alpha^{(n)})^2 \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left\| \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \bar{\phi}_{h,i}^{(n)}}^{-1}}^2 \right] + H^2 d L^2 \lambda} \\
& \quad + \sqrt{\tilde{A}^L (\alpha^{(n)})^2 \mathbb{E}_{(s, \mathbf{a}) \sim \rho_1^{(n)}} \left[\left\| \bar{\phi}_{1,i}^{(n)}(s, \mathbf{a}) \right\|_{\Sigma_{n, \rho_1^{(n)}, \bar{\phi}_{1,i}^{(n)}}^{-1}}^2 \right]}.
\end{aligned}$$

Note that we use the fact that $B = H$ when applying Lemma D.3. In addition, we have

$$\begin{aligned}
& n \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left\| \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) \right\|_{\Sigma_{n, \rho_h^{(n)}, \bar{\phi}_{h,i}^{(n)}}^{-1}}^2 \right] \\
& = n \text{Tr} \left(\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a})^\top \right] \left(n \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\bar{\phi}_{h,i}^{(n)}(s, \mathbf{a}) \bar{\phi}_{h,i}^{(n)}(s, \mathbf{a})^\top \right] + \lambda I_{d|Z_i|} \right)^{-1} \right) \\
& \leq d^L.
\end{aligned}$$

Then,

$$\begin{aligned}
& \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[\hat{\beta}_h^{(n)}(s, \mathbf{a}) \right] \\
& \leq \sum_{h=1}^{H-1} \sum_{i=1}^M \tilde{A}^L \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[\left\| \tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_h^{(n)}, \tilde{\phi}_{h,i}^*}^{-1}} \right] \sqrt{d^L (\alpha^{(n)})^2 + H^2 d L^2 \lambda} \\
& \quad + \sqrt{d^L \tilde{A}^L (\alpha^{(n)})^2 / n}.
\end{aligned}$$

Second, we calculate the term (b) in inequality equation 23. Following Lemma D.3 and noting $f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i)$ is upper-bounded by 2 (i.e., $B = 2$ in Lemma D.3), we have

$$\begin{aligned}
& \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[f_h^{(n)}(s, \mathbf{a}) \right] \\
& \leq \sum_{i=1}^M \sum_{h=1}^H \mathbb{E}_{(s, \mathbf{a}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) \right] \\
& \leq \sum_{i=1}^M \sum_{h=1}^{H-1} \tilde{A} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi^{(n)}}} \left[\left\| \tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_h^{(n)}, \tilde{\phi}_{h,i}^*}^{-1}} \right] \sqrt{n \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) \right)^2 \right] + d^L \lambda} \\
& \quad + \sqrt{\tilde{A} \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_1^{(n)}(s[Z_j], \mathbf{a}_j) \right)^2 \right]}
\end{aligned}$$

$$\begin{aligned} &\leq \sum_{i=1}^M \sum_{h=1}^{H-1} \tilde{A} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_h^{(n)}, \tilde{\phi}_{h,i}^*}^{-1}} \right] \sqrt{n\zeta^{(n)} + d^L \lambda} + \sqrt{\tilde{A}\zeta^{(n)}} \\ &\lesssim \frac{\alpha^{(n)}}{H} \sum_{i=1}^M \sum_{h=1}^{H-1} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_h^{(n)}, \tilde{\phi}_{h,i}^*}^{-1}} \right] + \sqrt{\tilde{A}\zeta^{(n)}}, \end{aligned}$$

where in the second inequality, we use $\mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} \left[\left(f_{h,i}^{(n)}(s[Z_i], \mathbf{a}_i) \right)^2 \right] \leq \zeta^{(n)}$, and in the last line,

recall $\tilde{A} \sqrt{n\zeta^{(n)} + d^L \lambda} \lesssim \alpha^{(n)}/H$. Then, by combining the above calculation of the term (a) and term (b) in inequality equation 23, we have:

$$\begin{aligned} &\bar{v}_i^{(n)} - \underline{v}_i^{(n)} \\ &= \mathbb{E}_{s \sim d_{P^*, 1}^{\pi(n)}} \left[\bar{V}_{1,i}^{(n)}(s) - \underline{V}_{1,i}^{(n)}(s) \right] \\ &\lesssim \sum_{i=1}^M \sum_{h=1}^{H-1} \left(\tilde{A}^L \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_h^{(n)}, \tilde{\phi}_{h,i}^*}^{-1}} \right] \sqrt{d^L (\alpha^{(n)})^2 + H^2 d^L \lambda} + \sqrt{\frac{d^L \tilde{A}^L (\alpha^{(n)})^2}{n}} \right) \\ &\quad + H^2 M \sum_{i=1}^M \sum_{h=1}^{H-1} \left(\frac{\alpha^{(n)}}{H} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_h^{(n)}, \tilde{\phi}_{h,i}^*}^{-1}} \right] + \sqrt{\tilde{A}\zeta^{(n)}} \right). \end{aligned}$$

Taking maximum over i on both sides and use the definition of $\Delta^{(n)}$, we get

$$\begin{aligned} \Delta^{(n)} &= \max_{i \in [M]} \left\{ \bar{v}_i^{(n)} - \underline{v}_i^{(n)} \right\} + 2HM \sqrt{\tilde{A}\zeta^{(n)}} \\ &\lesssim \sum_{i=1}^M \sum_{h=1}^{H-1} \left(\tilde{A}^L \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_h^{(n)}, \tilde{\phi}_{h,i}^*}^{-1}} \right] \sqrt{d^L (\alpha^{(n)})^2 + H^2 d^L \lambda} + \sqrt{\frac{d^L \tilde{A}^L (\alpha^{(n)})^2}{n}} \right) \\ &\quad + H^2 M \sum_{i=1}^M \sum_{h=1}^{H-1} \left(\frac{\alpha^{(n)}}{H} \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_h^{(n)}, \tilde{\phi}_{h,i}^*}^{-1}} \right] + \sqrt{\tilde{A}\zeta^{(n)}} \right). \end{aligned}$$

Hereafter, we take the dominating term out. Note that

$$\begin{aligned} &\sum_{n=1}^N \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_h^{(n)}, \tilde{\phi}_{h,i}^*}^{-1}} \right] \\ &\leq \sqrt{N \sum_{n=1}^N \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}})^\top \Sigma_{n, \gamma_h^{(n)}, \tilde{\phi}_{h,i}^*}^{-1} \tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}}) \right]} \quad (\text{CS inequality}) \\ &\lesssim \sqrt{N \left(\log \det \left(\lambda I_{d^{\cup_{j \in Z_i} Z_j}} + \sum_{n=1}^N \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}}) \tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}})^\top \right] \right) - \log \det(\lambda I_{d^{\cup_{j \in Z_i} Z_j}}) \right)} \\ &\quad (\text{Lemma E.2}) \\ &\leq \sqrt{d^L N \log \left(1 + \frac{N}{d\lambda} \right)}. \end{aligned}$$

(Potential function bound, Lemma E.3 noting $\|\phi_{h,i}^*(s[Z_i], \mathbf{a}_i)\|_2 \leq 1$ for any (s, \mathbf{a}) .)

Similarly, we have

$$\sum_{n=1}^N \mathbb{E}_{(\tilde{s}, \tilde{\mathbf{a}}) \sim d_{P^*, h}^{\pi(n)}} \left[\left\| \tilde{\phi}_{h,i}^*(\tilde{s}, \tilde{\mathbf{a}}) \right\|_{\Sigma_{n, \gamma_h^{(n)}, \tilde{\phi}_{h,i}^*}^{-1}} \right] \leq \sqrt{d^L N \log \left(1 + \frac{N}{d\lambda} \right)}.$$

Finally,

$$\sum_{n=1}^N \Delta^{(n)} \lesssim HM \left(\sqrt{d^L N \log \left(1 + \frac{N}{d\lambda} \right)} \tilde{A}^L \sqrt{d^L (\alpha^{(N)})^2 + H^2 d^L \lambda} + \sum_{n=1}^N \sqrt{\frac{d^L \tilde{A}^L (\alpha^{(n)})^2}{n}} \right)$$

$$\begin{aligned}
& + H^3 M^2 \left(\frac{1}{H} \sqrt{d^L N \log \left(1 + \frac{N}{d\lambda} \right)} \alpha^{(N)} + \sum_{n=1}^N \sqrt{\tilde{A} \zeta^{(n)}} \right) \\
& \lesssim H^2 M^2 d^{L^2} \tilde{A}^L \sqrt{N \log \left(1 + \frac{N}{d\lambda} \right)} \alpha^{(N)}.
\end{aligned}$$

(Some algebra. We take the dominating term out. Note that $\alpha^{(n)}$ is increasing in n)

This concludes the proof. \square

D.3 PROOF OF THE MAIN THEOREMS

Lemma D.9. *For the model-based algorithm, when we pick $\lambda = \Theta \left(L d^L \log \frac{N H M |\Phi|}{\delta} \right)$, $\alpha^{(n)} = \Theta \left(H \tilde{A} \sqrt{n \zeta^{(n)} + d^L \lambda} \right)$ and $\zeta^{(n)} = \Theta \left(\frac{1}{n} \log \frac{|\mathcal{M}| H N M}{\delta} \right)$, with probability $1 - \delta$, we have*

$$\sum_{n=1}^N \Delta^{(n)} \lesssim H^3 M^2 d^{(L+1)^2} A^{\frac{L+1}{2}} N^{\frac{1}{2}} \log \frac{|\mathcal{M}| H N M}{\delta}.$$

Proof. The result of Lemma D.1 implies with our choice of λ and $\zeta^{(n)}$, the event \mathcal{E} holds with probability at least $1 - \delta$. In this case, we have

$$\alpha^{(n)} = \Theta \left(H \tilde{A} \sqrt{\log \frac{|\mathcal{M}| H N M}{\delta} + L d^{2L} \log \frac{N H M |\Phi|}{\delta}} \right), \quad (24)$$

which is a constant unrelated with n . Therefore, using the result of Lemma D.8, we get

$$\sum_{n=1}^N \Delta^{(n)} \lesssim H^2 d^{L^2} \tilde{A}^L M^2 \sqrt{N \log \left(1 + \frac{N}{d\lambda} \right)} \alpha^{(N)} \lesssim H^3 M^2 d^{(L+1)^2} A^{L+1} L^{\frac{1}{2}} N^{\frac{1}{2}} \log \frac{|\mathcal{M}| H N M}{\delta},$$

which has finished the proof. \square

Proof of Theorem 4.1

Proof. For any fixed episode n and agent i , by Lemma D.5, Lemma D.6 and Lemma D.7, we have

$$v_i^{\dagger, \pi_{-i}^{(n)}} - v_i^{\pi^{(n)}} \left(\text{or } \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi^{(n)}} - v_i^{\pi^{(n)}} \right) \leq \bar{v}_i^{(n)} - v_i^{(n)} + 2 H M \sqrt{\tilde{A} \zeta^{(n)}} \leq \Delta^{(n)}.$$

Taking maximum over i on both sides, we have

$$\max_{i \in [M]} \left\{ v_i^{\dagger, \pi_{-i}^{(n)}} - v_i^{\pi^{(n)}} \right\} \left(\text{or } \max_{i \in [M]} \left\{ \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi^{(n)}} - v_i^{\pi^{(n)}} \right\} \right) \leq \Delta^{(n)}. \quad (25)$$

From Lemma B.8, with probability $1 - \delta$, we can ensure

$$\sum_{n=1}^N \Delta^{(n)} \lesssim H^3 M^2 d^{(L+1)^2} A^{L+1} L^{\frac{1}{2}} N^{\frac{1}{2}} \log \frac{|\mathcal{M}| H N M}{\delta}.$$

Therefore, according to Lemma E.4, when we pick N to be

$$O \left(\frac{L^5 M^4 H^6 d^{2(L+1)^2} \tilde{A}^{2(L+1)}}{\varepsilon^2} \log^2 \left(\frac{H d A L M |\mathcal{M}|}{\delta \varepsilon} \right) \right),$$

we have

$$\frac{1}{N} \sum_{n=1}^N \Delta^{(n)} \leq \varepsilon.$$

On the other hand, from equation 25, we have

$$\begin{aligned} & \max_{i \in [M]} \left\{ v_i^{\dagger, \hat{\pi}^{-i}} - v_i^{\hat{\pi}} \right\} \left(\text{or } \max_{i \in [M]} \left\{ \max_{\omega \in \Omega_i} v_i^{\omega \circ \hat{\pi}} - v_i^{\hat{\pi}} \right\} \right) \\ &= \max_{i \in [M]} \left\{ v_i^{\dagger, \pi^{(n^*)}^{-i}} - v_i^{\pi^{(n^*)}} \right\} \left(\text{or } \max_{i \in [M]} \left\{ \max_{\omega \in \Omega_i} v_i^{\omega \circ \pi^{(n^*)}} - v_i^{\pi^{(n^*)}} \right\} \right) \\ &\leq \Delta^{(n^*)} = \min_{n \in [N]} \Delta^{(n)} \leq \frac{1}{N} \sum_{n=1}^N \Delta^{(n)} \leq \varepsilon, \end{aligned}$$

which has finished the proof, noting our assumption that $L = O(1)$. \square

E AUXILIARY LEMMAS

Lemma E.1 (Concentration of the bonus term (Zanette et al. (2021), Lemma 39)). *Set $\lambda^{(n)} \geq \Theta(d \log(nH|\Phi|/\delta))$ for any n . Define*

$$\Sigma_{n, \rho_h^{(n)}, \phi} = n \mathbb{E}_{(s, \mathbf{a}) \sim \rho_h^{(n)}} [\phi(s, \mathbf{a}) \phi^\top(s, \mathbf{a})] + \lambda^{(n)} I_d, \quad \hat{\Sigma}_{h, \phi}^{(n)} = \sum_{i=1}^n \phi(s_h^{(i)}, \mathbf{a}_h^{(i)}) \phi^\top(s_h^{(i)}, \mathbf{a}_h^{(i)}) + \lambda^{(n)} I_d.$$

With probability $1 - \delta$, we have

$$\forall n \in \mathbb{N}^+, \forall h \in [H], \forall \phi \in \Phi, \quad c_1 \|\phi(s, \mathbf{a})\|_{\Sigma_{n, \rho_h^{(n)}, \phi}^{-1}} \leq \|\phi(s, \mathbf{a})\|_{(\hat{\Sigma}_{h, \phi}^{(n)})^{-1}} \leq c_2 \|\phi(s, \mathbf{a})\|_{\rho_h^{(n)}, \phi}^{-1}.$$

Lemma E.2 (Agarwal et al. (2020a), Lemma G.2). *Consider the following process. For $n = 1, \dots, N$, $M_n = M_{n-1} + G_n$ with $M_0 = \lambda_0 I$ and G_n being a positive semidefinite matrix with eigenvalues upper bounded by 1. We have*

$$2 \log \det(M_N) - 2 \log \det(\lambda_0 I) \geq \sum_{n=1}^N \text{Tr}(G_n M_{n-1}^{-1}).$$

Lemma E.3 (Potential function lemma). *Suppose $\text{Tr}(G_n) \leq B^2$.*

$$2 \log \det(M_N) - 2 \log \det(\lambda_0 I) \leq d \log \left(1 + \frac{NB^2}{d\lambda_0} \right)$$

Proof. Let $\sigma_1, \dots, \sigma_d$ be the set of singular values of M_N recalling M_N is a positive semidefinite matrix. Then, by the AM-GM inequality,

$$\log \det(M_N) / \det(\lambda_0 I) = \log \prod_{i=1}^d (\sigma_i / \lambda_0) \leq \log d \left(\frac{1}{d} \sum_{i=1}^d (\sigma_i / \lambda_0) \right)$$

Since we have $\sum_i \sigma_i = \text{Tr}(M_N) \leq d\lambda_0 + NB^2$, the statement is concluded. \square

Lemma E.4. *For parameters A, B, ε such that $\frac{A^2 B}{\varepsilon^2}$ is larger than some absolute constant, when we pick $N = \frac{A^2}{\varepsilon^2} \log^2 \frac{A^4 B^2}{\varepsilon^4} = O\left(\frac{A^2}{\varepsilon^2} \log^2 \frac{AB}{\varepsilon}\right)$, we have*

$$\frac{A}{\sqrt{N}} \log(BN) \leq \varepsilon.$$

Proof. We have

$$\frac{A}{\sqrt{N}} \log(BN) = \varepsilon \frac{\log\left(\frac{A^2 B}{\varepsilon^2} \log^2 \frac{A^4 B^2}{\varepsilon^4}\right)}{\log \frac{A^4 B^2}{\varepsilon^4}}$$

Note that

$$\frac{A^2 B}{\varepsilon^2} \log^2 \frac{A^4 B^2}{\varepsilon^4} \leq \frac{A^4 B^2}{\varepsilon^4} \Leftrightarrow \log^2 \frac{A^4 B^2}{\varepsilon^4} \leq \frac{A^2 B}{\varepsilon^2}$$

where the right hand side is always true whenever $\frac{A^2 B}{\varepsilon^2}$ is larger than some given constant. Therefore, we get

$$\frac{A}{\sqrt{N}} \log(BN) \leq \varepsilon.$$

□

F EXPERIMENT DETAILS

F.1 DETAILED ENVIRONMENT SETUP

In this section we introduce the details of the environment construction of the Block Markov games. For completeness we repeat certain details already introduced in the main text. We design our Block Markov game by first randomly generating a tabular Markov game with horizon H , 3 states, 2 players each with 3 actions, and random reward matrix $R_h \in (0, 1)^{3 \times 3^2 \times H}$ and random transition matrix $T_h(s_h, a_h) \in \Delta(\mathcal{S}_{h+1})$. For the reward generalization, for each $r(s, a, s')$ entry in the reward matrix, we assign it with a random number sampled from a uniform distribution from -1 to 1. For the probability matrix generation, for each conditional distribution $T(\cdot|s, a)$, we randomly sample 3 numbers from a uniform distribution from -1 to 1 and form the probability simplex by normalization. For the generation of rich observation (emission distribution), we follow the experiment design of (Misra et al., 2020): the dimension of the observation is $2^{\lceil \log(H+|\mathcal{S}|+1) \rceil}$. For an observation o that emitted from state s and time step h , we concatenate the one-hot vector of s and h , adding i.i.d. Gaussian noise $\mathcal{N}(0, 0.1)$ on each entry, pend zero at the end if necessary, and finally multiply with a Hadamard matrix. In our setting, we have variants with different horizons H .

F.2 IMPLEMENTATION DETAILS

For the implementation of GERL_MG2, we break down the introduction into two parts: the implementation of Alg. 3 and the implementation of game solving algorithm with current features (line. 10 and line. 11 of Algorithm. 1). For the implementation of Algorithm. 3, we follow the same function approximation as (Zhang et al. (2022)) and adapt their open-sourced code at <https://github.com/yudasong/briee>. We include an overview of the function class for completeness: we adopt a two layer neural network with tanh non-linearity as the function class as the discriminator class. For the decoder, we let $\psi(o) = \text{softmax}(A^\top o)$, where $A \in \mathbb{R}^{|\mathcal{O}| \times 3}$, and we let $\phi(o, \mathbf{a}) = \psi(o) \otimes \mathbf{a}$. Here \mathbf{a} denotes the one-hot encoding in the joint action space.

Different from Zhang et al. (2022), we solve the optimization problem by directly solving the min-max-min problem instead of using an iterative method. We show the implementation in Algorithm. 5. We first perform minibatch stochastic gradient descent aggressively on the discriminator selection step (line. 5, on $\hat{\phi}$ and f) and the feature selection step (line. 6, on ϕ), where in each step we first compute the linear weight w and \hat{w} closed-formly and then perform gradient descent/ascend on the features and discriminators. Note that here the number of iteration T is very small.

For solving the Markov games, in addition to following Algorithm. 1, to solve line.10 (i.e., solving equation 2 or equation 3 or equation 4), we implement the NE/CCE solvers based on the public repository: <https://github.com/quantumiracle/MARS>. Note that the essential difference lies in that (Xie et al., 2020) assumes that the algorithm has the access to the ground-truth feature but our algorithm needs to utilize the different features we learn for each iteration. We also adopt the Deep RL baseline from the same public repository.

F.3 ZERO-SUM EXPERIMENT TRAINING CURVES

In this section we provide the training curves of GERL_MG2 and Deep RL baseline in the zero-sum setting in Figure. 1.

F.4 GENERAL-SUM EXPERIMENT DETAILS

In this section we complete the remaining details for the general-sum experiment. We include the training curve in Fig. 2.

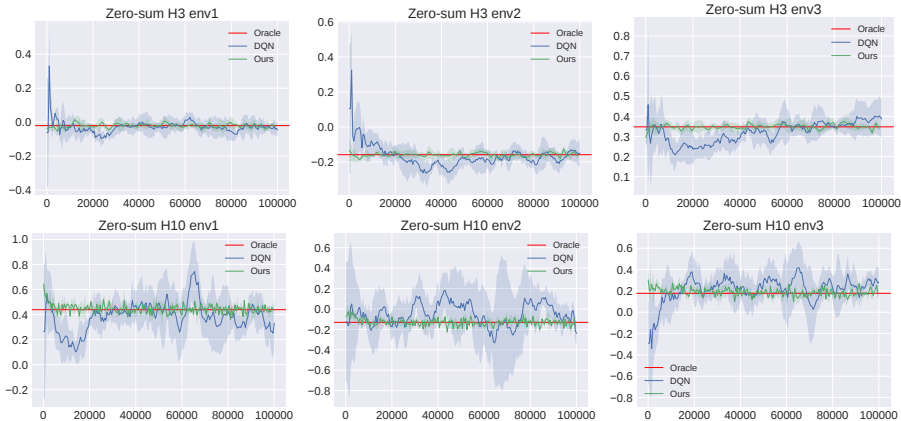


Figure 1: Training curve in the zero-sum setting. We evaluate each method over 5 random seeds and report the mean and standard deviation of the moving average of evaluation returns, wherein for each evaluation we perform 1000 runs. We use “Oracle” to denote the ground truth NE values of the Markov game. The x-axis denotes the number of episodes and the y-axis denotes the value of returns.

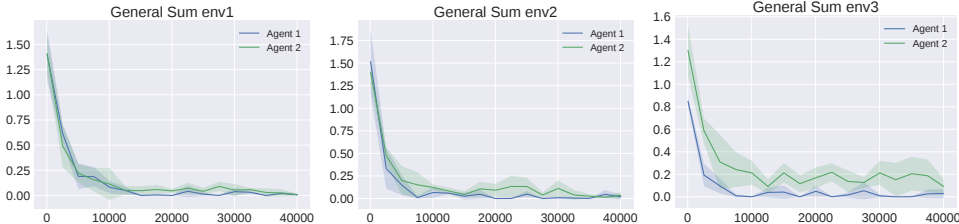


Figure 2: Training curve of GERL_MG2 in the general sum setting. In this setting, the y-axis denotes exploitability instead of raw returns.

F.5 HYPERPARAMETERS

In this section, we include the hyperparameter for GERL_MG2 in Table. 2, and the hyperparameter for DQN in Table. 3 and Table. 4.

Table 2: Hyperparameters for GERL_MG2.

	Value Considered	Final Value
Decoder ϕ learning rate	{1e-2}	1e-2
Discriminator f learning rate	{1e-}	1e-2
Discriminator f hidden layer size	{128,256,512}	256
RepLearn Iteration T	{10,20,30,50}	10
Decoder ϕ number of gradient steps	{64,128,256}	256
Discriminator f number of gradient steps	{64,128,256}	256
Decoder ϕ batch size	{128,256,512}	512
Discriminator f batch size	{128,256,512}	512
RepLearn regularization coefficient λ	{0.01}	0.01
Decoder ϕ softmax temperature	{1,0.5,0.1}	1
LSVI bonus coefficient β	{0.1,0.5,1}	0.1
LSVI regularization coefficient λ	{1}	1
Warm up samples	{0,200}	0

Algorithm 5 Model-free Representation Learning in Practice

-
- 1: **Input:** Dataset \mathcal{D} , step h , regularization λ , iterations T .
 - 2: Denote least squares loss: $\mathcal{L}_{\lambda, \mathcal{D}}(\phi, \theta, f) := \mathbb{E}_{\mathcal{D}} \left[(\phi(s, \mathbf{a})^\top \theta - f(s))^2 \right] + \lambda \|\theta\|_2^2$.
 - 3: Initialize $\phi_0 \in \Phi_h$ arbitrarily;
 - 4: **for** $t = 0, 1, \dots, T$ **do**
 - 5: Discriminator selection: $f_t = \arg \max_{f \in \mathcal{F}_h} [\min_{\theta} \mathcal{L}_{\lambda, \mathcal{D}}(\phi_t, \theta, f) - \min_{\tilde{\phi} \in \Phi, \tilde{\theta}} \mathcal{L}_{\lambda, \mathcal{D}}(\tilde{\phi}, \tilde{\theta}, f)]$
 - 6: Feature selection: $\phi_{t+1} = \arg \min_{\phi \in \Phi_h} \sum_{i=1}^t \min_{\theta_i} \mathcal{L}_{\lambda, \mathcal{D}}(\phi, \theta_i, f_i)$, $\hat{\phi} \leftarrow \phi_{t+1}$
 - 7: **end for**
 - 8: **Return** $\hat{\phi}$, \hat{P} where \hat{P} is calculated from equation 1.
-

Table 3: Hyperparameters for DQN in short horizon environment.

	Value considered	Final Value
Target update interval	{1000}	1000
ϵ_0	{1}	1
ϵ_N	{0.01}	0.01
ϵ decay frequency	{8000}	8000
Batch size	{8000}	8000
Optimizer	{Adam}	Adam
Learning Rate	{0.0001}	0.0001
Hidden layer	{[32,32,32]}	[32,32,32]
Self-play δ	{1.5}	1.5

Table 4: Hyperparameters for DQN in long horizon environment.

	Value considered	Final Value
Target update interval	{1000}	1000
ϵ_0	{1}	1
ϵ_N	{0.01}	0.01
ϵ decay frequency	{8000}	8000
Batch size	{8000}	8000
Optimizer	{Adam}	Adam
Learning Rate	{0.0001}	0.0001
Hidden layer	{[32,32,32]}	[32,32,32]
Self-play δ	{1.5,2}	2