

Algorithm 4 REGIME-exploration

Input: The number of total episodes K , bonus parameter β_{ex} and regularization parameter λ_{ex} .

for $k = 1, \dots, K$ **do**

 Initialize: $Q_{H+1}^k(\cdot, \cdot) \leftarrow 0, V_{H+1}^k(\cdot) \leftarrow 0$.

for $h = H, \dots, 1$ **do**

 Compute the covariance matrix: $\Lambda_h^k \leftarrow \sum_{i=1}^{k-1} \phi_h(s_h^i, a_h^i) \phi_h(s_h^i, a_h^i)^\top + \lambda_{\text{ex}} I$.

 Compute the bonus and reward:

$b_h^k(\cdot, \cdot) \leftarrow \min \{ \beta_{\text{ex}} \|\phi_h(\cdot, \cdot)\|_{(\Lambda_h^k)^{-1}}, H - h + 1 \}$ and $r_h^k = b_h^k / H$.

 Compute Q function:

$$Q_h^k(\cdot, \cdot) \leftarrow \text{Clip}_{[0, H-h+1]} \left(\text{Clip}_{[0, H-h+1]} \left((w_h^k)^\top \phi_h(\cdot, \cdot) + r_h^k(\cdot, \cdot) \right) + b_h^k(\cdot, \cdot) \right),$$

 where $w_h^k = (\Lambda_h^k)^{-1} \sum_{i=1}^{k-1} \phi_h(s_h^i, a_h^i) \cdot V_{h+1}^k(s_{h+1}^i)$.

 Compute value function and policy:

$$V_h^k(\cdot) \leftarrow \max_{a \in \mathcal{A}} Q_h^k(\cdot, a), \pi_h^k(\cdot) \leftarrow \arg \max_{a \in \mathcal{A}} Q_h^k(\cdot, a).$$

end for

 Collect a trajectory $\tau^k = (s_h^k, a_h^k, s_{h+1}^k)_{h=1}^H$ by running $\pi^k = \{\pi_h^k\}_{h=1}^H$ and add τ^k into \mathcal{D}_{ex} .

end for

Sample K states from the initial states $\{s_1^{i, \text{in}}\}_{i=1}^K$ and add them to \mathcal{D}_{in} .

Return $\mathcal{D}_{\text{ex}}, \mathcal{D}_{\text{in}}$.

A DISCUSSION

Validity of Linear Parametrization. In this work we consider linear reward parametrization, or more generally, linear trajectory embeddings. Such assumptions are commonly used in the theoretical works of PbRL (Pacchiano et al., 2021; Zhu et al., 2023) and relevant examples can be borrowed from the practical works (Pacchiano et al., 2020; Parker-Holder et al., 2020) like the behavior guided class of algorithms for policy optimization. We admit that our analysis cannot cover all kinds of reward parametrization, but we think it is a reasonable starting point to consider linear function approximation for theoretical analysis.

Applicability of REGIME. Our results are not restricted to linear MDPs and low-rank MDPs. Our main result (Theorem 1) only requires a suitable reward-free oracle for dynamics learning and in practice there exist some oracles ready for use, even when the MDPs are very complicated (Xu et al., 2022). This implies that by plugging in these general reward-free oracles, we are able to deal with general MDPs as well.

Relationship to Active Learning. The trajectory collection process in Step 1 of REGIME utilizes a similar idea to active learning. In active learning, people choose a trajectory pair to query human feedback which can maximize the information gain in each iteration. Similarly, in our algorithm, the estimated covariance matrix $\widehat{\Sigma}_n$ can be regarded as the current information in n -th iteration. Then we choose a pair of policies $(\pi^{n,0}, \pi^{n,1})$ which can maximize the information gain approximately (line 5, it is approximate because we are using an approximate dynamics \widehat{P} for evaluation).

B OMIT DETAILS IN SECTION 4

In this section we present the details of Algorithm 4 and Algorithm 5. Here $\text{Clip}_{[a,b]}(x)$ means $\min\{\max\{a, x\}, b\}$. In particular, when estimating $(\phi(\pi))_{h,j}$, we use the reward function $r_{h'}^{h,j}(s, a) = \phi_{h'}(s, a)^\top \theta_{h'}^{h,j}$ for all $h' \in [H]$ (up to an R factor) where

$$\theta_{h'}^{h,j} = \begin{cases} \frac{1}{R} \cdot e_j, & \text{if } h' = h, \\ 0, & \text{otherwise.} \end{cases}$$

Algorithm 5 REGIME-planning

Input: Dataset $\mathcal{D}_{\text{ex}} = \{(s_h^i, a_h^i, s_{h+1}^i)\}_{i=1, h=1}^{K, H}$, $\mathcal{D}_{\text{in}} = \{s_1^{i, \text{in}}\}_{i=1}^K$, bonus parameter β_{pl} and regularization parameter λ_{pl} , policy π , reward function $(r_h)_{h=1}^H$.

for $h' = H, \dots, 1$ **do**

 Compute the covariance matrix: $\Lambda_{h'} \leftarrow \sum_{i=1}^K \phi_{h'}(s_{h'}^i, a_{h'}^i) \phi_{h'}(s_{h'}^i, a_{h'}^i)^\top + \lambda_{\text{pl}} I$.

 Compute the bonus: $b_{h'}(\cdot, \cdot) \leftarrow \min \{\beta_{\text{pl}} \|\phi_{h'}(\cdot, \cdot)\|_{(\Lambda_{h'})^{-1}}, 2(H - h + 1)\}$.

end for

Initialize: $\widehat{Q}_{H+1}^{r, \pi}(\cdot, \cdot) \leftarrow 0, \widehat{V}_{H+1}^{r, \pi}(\cdot) \leftarrow 0$.

for $h' = H, \dots, 1$ **do**

 Compute Q function:

$$\widehat{Q}_{h'}^{r, \pi}(\cdot, \cdot) \leftarrow \text{Clip}_{[-(H-h+1), H-h+1]} \left(\text{Clip}_{[-(H-h+1), H-h+1]} \left((w_{h'}^{r, \pi})^\top \phi_{h'}(\cdot, \cdot) + r_{h'}(\cdot, \cdot) + b_{h'}(\cdot, \cdot) \right) \right),$$

 where $w_{h'}^{r, \pi} = (\Lambda_{h'})^{-1} \sum_{i=1}^K \phi_{h'}(s_{h'}^i, a_{h'}^i) \cdot \widehat{V}_{h'+1}^{r, \pi}(s_{h'+1}^i)$.

 Compute value function: $\widehat{V}_{h'}^{r, \pi}(\cdot) \leftarrow \mathbb{E}_{a \sim \pi_{h'}} \widehat{Q}_{h'}^{r, \pi}(\cdot, a)$.

end for

Compute $\widehat{V}^\pi(r) \leftarrow \frac{1}{K} \sum_{i=1}^K \widehat{V}_1^{r, \pi}(s_1^{i, \text{in}})$.

Return $\widehat{V}^\pi(r)$.

Here e_j is the one-hot vector whose j -th entry is 1. For simplicity, we denote $\widehat{V}^{r^{h,j}, \pi}, \widehat{Q}^{r^{h,j}, \pi}, w^{r^{h,j}, \pi}$ by $V^{h,j, \pi}, Q^{h,j, \pi}, w^{h,j, \pi}$ and let the estimation $(\widehat{\phi}(\pi))_{h,j}$ be $R\widehat{V}^\pi(r^{h,j})$.

Then we have the following formal theorem characterizing the sample complexity of Algorithm 2:

Theorem 4. *Let*

$$\lambda_{\text{ex}} = \lambda_{\text{pl}} = R^2,$$

$$\beta_{\text{ex}} = C_\beta dHR \sqrt{\log(dKHR/\delta)}, \beta_{\text{pl}} = C_\beta dHR \sqrt{\log(dKHR\mathcal{N}_\Pi(\epsilon')/\delta)}$$

$$\lambda \geq 4HR^2, K \geq \widetilde{\mathcal{O}}\left(\frac{H^8 B^2 R^4 d^4 \log(\mathcal{N}_\Pi(\epsilon')/\delta)}{\epsilon^2}\right), N \geq \widetilde{\mathcal{O}}\left(\frac{\lambda \kappa^2 B^2 R^2 H^4 d^2 \log(1/\delta)}{\epsilon^2}\right),$$

where $\epsilon' = \frac{\epsilon}{72BR^2 H \sqrt{d^H K^{H-1}}}$, $C_\beta > 0$ is a universal constant and $\kappa = 2 + \exp(2r_{\text{max}}) + \exp(-2r_{\text{max}})$. Then under Assumption 1 and 3, with probability at least $1 - \delta$, we have

$$V^{r^*, \hat{\pi}} \geq V^{r^*, * } - \epsilon.$$

The proof is deferred to Appendix E.1.

B.1 LOG-LINEAR POLICY CLASS

The sample complexity in Theorem 2 depends on the covering number of the policy class Π . Therefore we want to find a policy class for linear MDPs that is rich enough (i.e., contains near-global-optimal policies) while retains a small covering number at the same time. Indeed, the log-linear policy class (Agarwal et al., 2020b) satisfies this requirement, which is defined as follows:

$$\Pi = \left\{ \pi : \pi_h^\zeta(a|s) = \frac{\exp(\zeta_h^\top \phi_h(s, a))}{\sum_{a' \in \mathcal{A}} \exp(\zeta_h^\top \phi_h(s, a'))}, \zeta_h \in \mathbb{B}(d, W), \forall s \in \mathcal{S}, a \in \mathcal{A}, h \in [H] \right\}$$

Here $\mathbb{B}(d, W)$ is the d -dimensional ball centered at the origin with radius W . The following proposition characterizes the covering number of such log-linear policy class:

Proposition 1. *Let Π be the log-linear policy class. Then under Assumption 1, for any $\epsilon \leq 1$, we have $\log \mathcal{N}_\Pi(\epsilon) \leq Hd \log\left(\frac{12WR}{\epsilon}\right)$.*

Meanwhile, we can quantify the bias of such log-linear policy class as follows:

Proposition 2. Let $W = \frac{(B+(H+\epsilon)\sqrt{d})H \log |\mathcal{A}|}{\epsilon}$, then under Assumption 1 and 3, we have

$$V^{r^*, \pi_g} - \max_{\pi \in \Pi} V^{r^*, \pi} \leq \epsilon,$$

where π_g is the global optimal policy.

Combining Theorem 4, Proposition 1 and Proposition 2, we know that the returned policy $\hat{\pi}$ by Algorithm 2 with log-linear policy classes can indeed compete against the global optimal policy with the following sample complexities:

$$N_{\text{tra}} = K + N = \tilde{\mathcal{O}}\left(\frac{d^5 + \kappa^2 d^2}{\epsilon^2}\right), N_{\text{hum}} = \tilde{\mathcal{O}}\left(\frac{\kappa^2 d^2}{\epsilon^2}\right).$$

C PROOF OF THEOREM 1 WITH KNOWN TRANSITIONS

In this section, we consider the proof of Theorem 1 when transitions are known, i.e., $\epsilon' = 0$ and $\hat{P} = P^*$. In this case we have $\hat{\phi}(\pi) = \phi(\pi)$. We will deal with the unknown transition in Appendix D.1.

First, note that from the definition of $\hat{\pi}$, we have

$$V^{\hat{r}, \hat{\pi}} \geq V^{\hat{r}, \pi^*},$$

where π^* is the optimal policy with respect to the ground-truth reward r^* , i.e., $\pi^* = \arg \max_{\pi \in \Pi} V^{r^*, \pi}$. Therefore we can expand the suboptimality as follows:

$$\begin{aligned} V^{r^*, *} - V^{r^*, \hat{\pi}} &= (V^{r^*, *} - V^{\hat{r}, \pi^*}) + (V^{\hat{r}, \pi^*} - V^{\hat{r}, \hat{\pi}}) + (V^{\hat{r}, \hat{\pi}} - V^{r^*, \hat{\pi}}) \\ &\leq (V^{r^*, *} - V^{\hat{r}, \pi^*}) + (V^{\hat{r}, \hat{\pi}} - V^{r^*, \hat{\pi}}) \\ &= \mathbb{E}_{\tau \sim (\pi^*, P^*)}[\langle \phi(\tau), \theta^* - \hat{\theta} \rangle] - \mathbb{E}_{\tau \sim (\hat{\pi}, P^*)}[\langle \phi(\tau), \theta^* - \hat{\theta} \rangle] \\ &= \langle \phi(\pi^*) - \phi(\hat{\pi}), \theta^* - \hat{\theta} \rangle \\ &\leq \|\phi(\pi^*) - \phi(\hat{\pi})\|_{\Sigma_{N+1}^{-1}} \cdot \|\theta^* - \hat{\theta}\|_{\Sigma_{N+1}}, \end{aligned} \quad (3)$$

where $\Sigma_n := \lambda I + \sum_{i=1}^{n-1} (\phi(\pi^{i,0}) - \phi(\pi^{i,1}))(\phi(\pi^{i,0}) - \phi(\pi^{i,1}))^\top$ for all $n \in [N+1]$. Here the third step is due to the definition of value function and the last step comes from Cauchy-Schwartz inequality. Next we will bound $\|\phi(\pi^*) - \phi(\hat{\pi})\|_{\Sigma_{N+1}^{-1}}$ and $\|\theta^* - \hat{\theta}\|_{\Sigma_{N+1}}$ respectively.

First for $\|\phi(\pi^*) - \phi(\hat{\pi})\|_{\Sigma_{N+1}^{-1}}$, notice that $\Sigma_{N+1} \succeq \Sigma_n$ for all $n \in [N+1]$, which implies

$$\begin{aligned} \|\phi(\pi^*) - \phi(\hat{\pi})\|_{\Sigma_{N+1}^{-1}} &\leq \frac{1}{N} \sum_{n=1}^N \|\phi(\pi^*) - \phi(\hat{\pi})\|_{\Sigma_n^{-1}} \leq \frac{1}{N} \sum_{n=1}^N \|\phi(\pi^{n,0}) - \phi(\pi^{n,1})\|_{\Sigma_n^{-1}} \\ &\leq \frac{1}{\sqrt{N}} \sqrt{\sum_{n=1}^N \|\phi(\pi^{n,0}) - \phi(\pi^{n,1})\|_{\Sigma_n^{-1}}^2}, \end{aligned} \quad (4)$$

where the second step comes from the definition of $\pi^{n,0}$ and $\pi^{n,1}$ and the last step is due to Cauchy-Schwartz inequality. To bound the right hand side of (4), we utilize the following Elliptical Potential Lemma:

Lemma 1 (Elliptical Potential Lemma). *For any $\lambda \geq R_x^2$ and $d \geq 1$, consider a sequence of vectors $\{x^n \in \mathbb{R}^d\}_{n=1}^N$ where $\|x^n\| \leq R_x$ for all $n \in [N]$. Let $\Sigma_n = \lambda I + \sum_{i=1}^{n-1} x^i (x^i)^\top$, then we have*

$$\sum_{n=1}^N \|x^n\|_{\Sigma_n^{-1}}^2 \leq 2d \log \left(1 + \frac{N}{d}\right).$$

The proof is deferred to Appendix C.1. Since we have $\lambda \geq 4HR^2$, by Lemma 1 we know

$$\sqrt{\sum_{n=1}^N \|\phi(\pi^{n,0}) - \phi(\pi^{n,1})\|_{\Sigma_n^{-1}}^2} \leq \sqrt{2HdN \log(1 + N/(Hd))}.$$

Combining the above inequality with (4), we have

$$\|\phi(\pi^*) - \phi(\hat{\pi})\|_{\Sigma_{N+1}^{-1}} \leq \sqrt{\frac{2Hd \log(1 + N/(Hd))}{N}}. \quad (5)$$

For $\|\theta^* - \hat{\theta}\|_{\Sigma_{N+1}}$, first note that $\hat{\theta}$ is the MLE estimator. Let $\tilde{\Sigma}_n$ denote the empirical cumulative covariance matrix $\lambda I + \sum_{i=1}^{n-1} (\phi(\tau^{i,0}) - \phi(\tau^{i,1}))(\phi(\tau^{i,0}) - \phi(\tau^{i,1}))^\top$, then from the literature (Zhu et al., 2023), we know that MLE has the following guarantee:

Lemma 2 (MLE guarantee). *For any $\lambda > 0$ and $\delta \in (0, 1)$, with probability at least $1 - \delta$, we have*

$$\|\hat{\theta} - \theta^*\|_{\tilde{\Sigma}_{N+1}} \leq C_{\text{MLE}} \cdot \sqrt{\kappa^2(Hd + \log(1/\delta)) + \lambda HB^2}, \quad (6)$$

where $\kappa = 2 + \exp(2r_{\max}) + \exp(-2r_{\max})$ and $C_{\text{MLE}} > 0$ is a universal constant.

The proof is deferred to Appendix C.2. With Lemma 2, to $\|\theta^* - \hat{\theta}\|_{\Sigma_{N+1}}$ we only need to show $\tilde{\Sigma}_{N+1}$ is close to Σ_{N+1} . This can be achieved by the following concentration result from the literature:

Lemma 3 (Pacchiano et al. (2021)[Lemma 7]). *For any $\lambda > 0$ and $\delta \in (0, 1)$, with probability at least $1 - \delta$, we have*

$$\|\theta^* - \hat{\theta}\|_{\Sigma_{N+1}}^2 \leq 2\|\theta^* - \hat{\theta}\|_{\tilde{\Sigma}_{N+1}}^2 + C_{\text{CON}} H^3 d R^2 B^2 \log(N/\delta), \quad (7)$$

where $C_{\text{CON}} > 0$ is a universal constant.

Therefore, combining (7) and (6), by union bound with probability at least $1 - \delta$, we have that

$$\|\theta^* - \hat{\theta}\|_{\Sigma_{N+1}} \leq C_1 \cdot \kappa B R \sqrt{\lambda H^3 d \log(N/\delta)}, \quad (8)$$

where C_1 is a universal constant.

Thus substituting (5) and (8) into (3), we have $V^*(r^*) - V(r^*, \hat{\pi}) \leq \epsilon$ with probability at least $1 - \delta$ as long as

$$N \geq \tilde{\mathcal{O}}\left(\frac{\lambda \kappa^2 B^2 R^2 H^4 d^2 \log(1/\delta)}{\epsilon^2}\right).$$

C.1 PROOF OF LEMMA 1

Note that when $\lambda \geq R_x^2$, we have $\|x^n\|_{\Sigma_n^{-1}} \leq 1$ for all $n \in [N]$, which implies that for all $n \in [N]$, we have

$$\|x^n\|_{\Sigma_n^{-1}}^2 \leq \log\left(1 + \|x^n\|_{\Sigma_n^{-1}}^2\right).$$

On the other hand, let w^n denote $\|x^n\|_{\Sigma_n^{-1}}$, then we know for any $n \in [N - 1]$

$$\begin{aligned} \log \det \Sigma_{n+1} &= \log \det(\Sigma_n + x^n (x^n)^\top) = \log \det(\Sigma_n^{1/2} (I + \Sigma_n^{-1/2} x^n (x^n)^\top \Sigma_n^{-1/2}) \Sigma_n^{1/2}) \\ &= \log \det(\Sigma_n) + \log \det(I + (\Sigma_n^{-1/2} x^n) (\Sigma_n^{-1/2} x^n)^\top) \\ &= \log \det(\Sigma_n) + \log \det(I + (\Sigma_n^{-1/2} x^n)^\top (\Sigma_n^{-1/2} x^n)) \\ &= \log \det(\Sigma_n) + \log\left(1 + \|x^n\|_{\Sigma_n^{-1}}^2\right), \end{aligned}$$

where the fourth step is due to the property of determinants. Therefore we have

$$\begin{aligned} \sum_{n=1}^N \log\left(1 + \|x^n\|_{\Sigma_n^{-1}}^2\right) &= \log \det \Sigma_{N+1} - \log \det \Sigma_1 = \log(\det \Sigma_{N+1} / \det \Sigma_1) \\ &= \log \det\left(I + \frac{1}{\lambda} \sum_{n=1}^N x^n (x^n)^\top\right). \end{aligned}$$

Now let $\{\lambda_i\}_{i=1}^d$ denote the eigenvalues of $\sum_{n=1}^N x^n(x^n)^\top$, then we know

$$\begin{aligned} \log \det \left(I + \frac{1}{\lambda} \sum_{n=1}^N x^n(x^n)^\top \right) &= \log \left(\prod_{i=1}^d (1 + \lambda_i/\lambda) \right) \\ &\leq d \log \left(\frac{1}{d} \sum_{i=1}^d (1 + \lambda_i/\lambda) \right) \leq d \log \left(1 + \frac{NR_x^2}{d\lambda} \right) \leq d \log \left(1 + \frac{N}{d} \right), \end{aligned}$$

where the third step comes from $\sum_{i=1}^d \lambda_i = \text{Tr} \left(\sum_{n=1}^N x^n(x^n)^\top \right) = \sum_{n=1}^N \|x^n\|^2 \leq NR_x^2$ and the last step is due to the fact that $\lambda \geq R_x^2$. This concludes our proof.

C.2 PROOF OF LEMMA 2

First note that we have the following lemma from literature:

Lemma 4 (Zhu et al., 2023)[Lemma 3.1]. *For any $\lambda' > 0$, with probability at least $1 - \delta$, we have*

$$\|\hat{\theta} - \theta^*\|_{D+\lambda'I} \leq O \left(\sqrt{\frac{\kappa^2(Hd + \log(1/\delta))}{N}} + \lambda'HB^2 \right),$$

where $D = \frac{1}{N} \sum_{i=1}^N (\phi(\tau^{i,0}) - \phi(\tau^{i,1}))(\phi(\tau^{i,0}) - \phi(\tau^{i,1}))^\top$.

Therefore let $\lambda' = \frac{\lambda}{N}$ and from the above lemma we can obtain

$$\|\hat{\theta} - \theta^*\|_{\frac{\Sigma_{N+1}}{N}} \leq O \left(\sqrt{\frac{\kappa^2(Hd + \log(1/\delta))}{N}} + \frac{\lambda HB^2}{N} \right),$$

which is equivalent to

$$\|\hat{\theta} - \theta^*\|_{\Sigma_{N+1}} \leq O \left(\sqrt{\kappa^2(Hd + \log(1/\delta)) + \lambda HB^2} \right).$$

This concludes our proof.

D PROOFS IN SECTION 3

D.1 PROOF OF THEOREM 1

Note that from the proof of Theorem 1 with known transition dynamics, we have:

$$V^{r^*,*} - V^{r^*,\hat{\pi}} \leq \langle \phi(\pi^*) - \phi(\hat{\pi}), \theta^* - \hat{\theta} \rangle + (V^{\hat{r},\pi^*} - V^{\hat{r},\hat{\pi}}), \quad (9)$$

Then we have

$$\begin{aligned} V^{r^*,*} - V^{r^*,\hat{\pi}} &\leq \langle \phi(\pi^*) - \hat{\phi}(\pi^*), \theta^* - \hat{\theta} \rangle + \langle \hat{\phi}(\hat{\pi}) - \phi(\hat{\pi}), \theta^* - \hat{\theta} \rangle \\ &\quad + \langle \hat{\phi}(\pi^*) - \hat{\phi}(\hat{\pi}), \theta^* - \hat{\theta} \rangle + (V^{\hat{r},\pi^*} - V^{\hat{r},\hat{\pi}}). \end{aligned} \quad (10)$$

Now we only need to bound the three terms in the RHS of (10). For the first and second term, we need to utilize the following lemma:

Lemma 5. *Let $d_h^\pi(s, a)$ and $\hat{d}_h^\pi(s, a)$ denote the visitation measure of policy π under P^* and \hat{P} . Then with probability at least $1 - \delta/4$, we have for all $h \in [H]$ and $\pi \in \Pi$,*

$$\|d_h^\pi - \hat{d}_h^\pi\|_1 \leq h\epsilon'. \quad (11)$$

Let \mathcal{E}_1 denote the event when (11) holds. Then under event \mathcal{E}_1 , we further have the following lemma:

Lemma 6. *Under event \mathcal{E}_1 , for all policy $\pi \in \Pi$ and vector $v = [v_1, \dots, v_H]$ where $v_h \in \mathbb{R}^d$ and $\|v_h\| \leq 2B$ for all $h \in [H]$ we have,*

$$|\langle \phi(\pi) - \hat{\phi}(\pi), v \rangle| \leq BRH^2\epsilon'.$$

Substitute Lemma 6 into (10), we have

$$V^{r^*,*} - V^{r^*,\hat{\pi}} \leq \langle \hat{\phi}(\pi^*) - \hat{\phi}(\hat{\pi}), \theta^* - \hat{\theta} \rangle + 2BRH^2\epsilon' + (V^{\hat{r},\pi^*} - V^{\hat{r},\hat{\pi}}).$$

Then by Cauchy-Schwartz inequality, we have under event \mathcal{E}_1 ,

$$V^{r^*,*} - V^{r^*,\hat{\pi}} \leq \|\hat{\phi}(\pi^*) - \hat{\phi}(\hat{\pi})\|_{\hat{\Sigma}_{N+1}^{-1}} \cdot \|\theta^* - \hat{\theta}\|_{\hat{\Sigma}_{N+1}} + 2BRH^2\epsilon' + (V^{\hat{r},\pi^*} - V^{\hat{r},\hat{\pi}}). \quad (12)$$

Following the same analysis in the proof of Theorem 1 with known transition, we know

$$\|\hat{\phi}(\pi^*) - \hat{\phi}(\hat{\pi})\|_{\hat{\Sigma}_{N+1}^{-1}} \leq \sqrt{\frac{2Hd \log(1 + N/(Hd))}{N}}. \quad (13)$$

Now we only need to bound $\|\theta^* - \hat{\theta}\|_{\hat{\Sigma}_{N+1}}$. Similar to the proof of Theorem 1 with known transition, we use Σ_n and $\tilde{\Sigma}_n$ to denote $\lambda I + \sum_{i=1}^{n-1} (\phi(\pi^{i,0}) - \phi(\pi^{i,1}))(\phi(\pi^{i,0}) - \phi(\pi^{i,1}))^\top$ and $\lambda I + \sum_{i=1}^{n-1} (\phi(\tau^{i,0}) - \phi(\tau^{i,1}))(\phi(\tau^{i,0}) - \phi(\tau^{i,1}))^\top$ respectively. Then under event \mathcal{E}_1 , we have the following connection between $\hat{\Sigma}_{N+1}$ and Σ_{N+1} :

Lemma 7. *Under event \mathcal{E}_1 , we have*

$$\|\theta^* - \hat{\theta}\|_{\hat{\Sigma}_{N+1}} \leq \sqrt{2}\|\theta^* - \hat{\theta}\|_{\Sigma_{N+1}} + 2\sqrt{2}BRH^2\epsilon'.$$

Combining Lemma 7 with Lemma 2 and Lemma 3, we have under event $\mathcal{E}_1 \cap \mathcal{E}_2$,

$$\begin{aligned} \|\theta^* - \hat{\theta}\|_{\hat{\Sigma}_{N+1}} &\leq \sqrt{2}\|\theta^* - \hat{\theta}\|_{\Sigma_{N+1}} + 2\sqrt{2}BRH^2\epsilon' \\ &\leq C_2 \cdot \kappa BR \sqrt{\lambda H^3 d \log(N/\delta)} + 2\sqrt{2}BRH^2\epsilon', \end{aligned} \quad (14)$$

where $\Pr(\mathcal{E}_2) \geq 1 - \delta/2$ and $C_2 > 0$ is a universal constant.

Now we only need to bound $(V^{\hat{r},\pi^*} - V^{\hat{r},\hat{\pi}})$, which can be achieved with Lemma 6:

$$\begin{aligned} V^{\hat{r},\pi^*} - V^{\hat{r},\hat{\pi}} &= \langle \phi(\pi^*), \hat{\theta} \rangle - \langle \phi(\hat{\pi}), \hat{\theta} \rangle \\ &= \langle \phi(\pi^*) - \hat{\phi}(\pi^*), \hat{\theta} \rangle + \langle \hat{\phi}(\pi^*) - \hat{\phi}(\hat{\pi}), \hat{\theta} \rangle + \langle \hat{\phi}(\hat{\pi}) - \phi(\hat{\pi}), \hat{\theta} \rangle \leq 2BRH^2\epsilon', \end{aligned} \quad (15)$$

where the last step comes from Lemma 6 and the definition of $\hat{\pi}$.

Combining (12), (13) (14) and (15), we have $V^{r^*,*} - V^{r^*,\hat{\pi}} \leq \epsilon$ with probability at least $1 - \delta$ as long as

$$\epsilon' \leq \frac{\epsilon}{6BRH^2}, \quad N \geq \tilde{\mathcal{O}}\left(\frac{\lambda \kappa^2 B^2 R^2 H^4 d^2 \log(1/\delta)}{\epsilon^2}\right).$$

D.2 PROOF OF LEMMA 5

First notice that $d_h^\pi(s, a) = d_h^\pi(s)\pi(a|s)$ and $\hat{d}_h^\pi(s, a) = \hat{d}_h^\pi(s)\pi(a|s)$, which implies that for all $h \in [H]$

$$\begin{aligned} \|d_h^\pi - \hat{d}_h^\pi\|_1 &= \sum_{s,a} |d_h^\pi(s, a) - \hat{d}_h^\pi(s, a)| = \sum_{s,a} |d_h^\pi(s) - \hat{d}_h^\pi(s)|\pi(a|s) \\ &= \sum_s |d_h^\pi(s) - \hat{d}_h^\pi(s)| \sum_a \pi(a|s) = \sum_s |d_h^\pi(s) - \hat{d}_h^\pi(s)|. \end{aligned}$$

Therefore we only need to prove $\sum_s |d_h^\pi(s) - \hat{d}_h^\pi(s)| \leq h\epsilon'$ for all $h \in [H]$. We use induction to prove this. First for the base case, we have $\sum_s |d_1^\pi(s) - \hat{d}_1^\pi(s)| = \sum_s |P_1^*(s) - \hat{P}_1(s)| \leq \epsilon'$ according to the guarantee of the reward-free learnign oracle \mathcal{P} .

Now assume that $\sum_s |d_{h'}^\pi(s) - \hat{d}_{h'}^\pi(s)| \leq h'\epsilon'$ for all $h' \in [h]$ where $h \in [H - 1]$. Then we have

$$\sum_s |d_{h+1}^\pi(s) - \hat{d}_{h+1}^\pi(s)| = \sum_s \left| \sum_{s',a'} \hat{d}_h^\pi(s')\pi(a'|s')\hat{P}_h(s|s',a') - d_h^\pi(s')\pi(a'|s')P_h^*(s|s',a') \right|$$

$$\begin{aligned}
&\leq \left(\sum_{s,s',a'} \left| \widehat{d}_h^\pi(s') - d_h^\pi(s') \right| \pi(a'|s') \widehat{P}_h(s|s',a') \right) \\
&\quad + \left(\sum_{s,s',a'} d_h^\pi(s') \pi(a'|s') \left| \widehat{P}_h(s|s',a') - P_h^*(s|s',a') \right| \right) \\
&= \left(\sum_{s'} \left| \widehat{d}_h^\pi(s') - d_h^\pi(s') \right| \sum_{a'} \pi(a'|s') \sum_s \widehat{P}_h(s|s',a') \right) \\
&\quad + \mathbb{E}_{\pi, P^*} [\| \widehat{P}_h(\cdot|s',a') - P_h^*(\cdot|s',a') \|_1] \\
&\leq (h+1)\epsilon',
\end{aligned}$$

where the second step comes from the triangle inequality and the last step is due to the induction hypothesis and the guarantee of \mathcal{P} . Therefore, we have $\sum_s |d_{h+1}^\pi(s) - \widehat{d}_{h+1}^\pi(s)| \leq (h+1)\epsilon'$. Then by induction, we know $\sum_s |d_h^\pi(s) - \widehat{d}_h^\pi(s)| \leq h\epsilon'$ for all $h \in [H]$, which concludes our proof.

D.3 PROOF OF LEMMA 6

Note that from the definition of $\phi(\pi)$ we have

$$\langle \phi(\pi), v \rangle = \mathbb{E}_{\tau \sim (\pi, P^*)} \left[\sum_{h=1}^H \phi_h^\top(s_h, a_h) v_h \right] = \sum_{h=1}^H \sum_{s_h, a_h} d_h^\pi(s_h, a_h) \phi_h^\top(s_h, a_h) v_h.$$

Similarly, we have

$$\langle \widehat{\phi}(\pi), v \rangle = \sum_{h=1}^H \sum_{s_h, a_h} \widehat{d}_h^\pi(s_h, a_h) \phi_h^\top(s_h, a_h) v_h.$$

Therefore,

$$\begin{aligned}
|\langle \phi(\pi) - \widehat{\phi}(\pi), v \rangle| &\leq \sum_{h=1}^H \sum_{s_h, a_h} |\widehat{d}_h^\pi(s_h, a_h) - d_h^\pi(s_h, a_h)| \cdot |\phi_h^\top(s_h, a_h) v_h| \\
&\leq 2BR \sum_{h=1}^H \sum_{s_h, a_h} |\widehat{d}_h^\pi(s_h, a_h) - d_h^\pi(s_h, a_h)| \\
&\leq 2BR \sum_{h=1}^H h\epsilon' \leq BRH^2\epsilon',
\end{aligned}$$

where the first step is due to the triangle inequality and the third step comes from Lemma 5. This concludes our proof.

D.4 PROOF OF LEMMA 7

We use $\Delta\theta$ to denote $\theta^* - \hat{\theta}$ in this proof. From Lemma 6, we know that for any policy π ,

$$|\langle \phi(\pi) - \widehat{\phi}(\pi), \Delta\theta \rangle| \leq BRH^2\epsilon'.$$

By the triangle inequality, this implies that for any policy π^0, π^1 ,

$$|\langle \widehat{\phi}(\pi^0) - \widehat{\phi}(\pi^1), \Delta\theta \rangle| \leq |\langle \phi(\pi^0) - \phi(\pi^1), \Delta\theta \rangle| + 2BRH^2\epsilon'.$$

Therefore we have for any policy π^0, π^1 ,

$$|\langle \widehat{\phi}(\pi^0) - \widehat{\phi}(\pi^1), \Delta\theta \rangle|^2 \leq 2|\langle \phi(\pi^0) - \phi(\pi^1), \Delta\theta \rangle|^2 + 8(BRH^2\epsilon')^2. \quad (16)$$

Note that from the definition of $\widehat{\Sigma}_{N+1}$ and Σ_{N+1} , we have

$$\|\Delta\theta\|_{\widehat{\Sigma}_{N+1}}^2 = \Delta\theta^\top \left(\lambda I + \sum_{n=1}^N (\widehat{\phi}(\pi^{n,0}) - \widehat{\phi}(\pi^{n,1})) (\widehat{\phi}(\pi^{n,0}) - \widehat{\phi}(\pi^{n,1}))^\top \right) \Delta\theta$$

$$\begin{aligned}
&= \lambda \|\Delta\theta\|^2 + \sum_{n=1}^N |\langle \widehat{\phi}(\pi^{n,0}) - \widehat{\phi}(\pi^{n,1}), \Delta\theta \rangle|^2 \\
&\leq 2 \left(\lambda \|\Delta\theta\|^2 + \sum_{n=1}^N |\langle \phi(\pi^{n,0}) - \phi(\pi^{n,1}), \Delta\theta \rangle|^2 \right) + 8(BRH^2\epsilon')^2 \\
&= 2\|\Delta\theta\|_{\Sigma_{N+1}}^2 + 8(BRH^2\epsilon')^2,
\end{aligned}$$

where the third step comes from (16). This implies that

$$\|\Delta\theta\|_{\widehat{\Sigma}_{N+1}} \leq \sqrt{2}\|\Delta\theta\|_{\Sigma_{N+1}} + 2\sqrt{2}BRH^2\epsilon',$$

which concludes our proof.

E PROOFS IN SECTION 4 AND APPENDIX B

E.1 PROOF OF THEOREM 4

First note that Algorithm 4 provides us with the following guarantee:

Lemma 8. *We have with probability at least $1 - \delta/6$ that*

$$\mathbb{E}_{s_1 \sim P_1(\cdot)}[V_1^{b/H,*}(s_1)] \leq C_{\text{lin}} \sqrt{d^3 H^4 R^2 \cdot \log(\mathcal{N}_{\Pi}(\epsilon') dKH R/\delta)/K},$$

where b_h is defined in Algorithm 5 and $C_{\text{lin}} > 0$ is a universal constant. Here $V_1^{r,*}(s_1) := \max_{\pi \in \Pi} V_1^{r,\pi}(s_1)$.

Lemma 8 is adapted from Wang et al. (2020)[Lemma 3.2] and we highlight the difference of the proof in Appendix E.2. Then we consider a ϵ' -covering for Π , denoted by $\mathcal{C}(\Pi, \epsilon')$. Following the similar analysis in Wang et al. (2020)[Lemma 3.3], we have the following lemma:

Lemma 9. *With probability $1 - \delta/6$, for all $h' \in [H]$, policy $\pi \in \mathcal{C}(\Pi, \epsilon')$ and linear reward function r with $r_h \in [-1, 1]$, we have*

$$Q_{h'}^{r,\pi}(\cdot, \cdot) \leq \widehat{Q}_{h'}^{r,\pi}(\cdot, \cdot) \leq r_{h'}(\cdot, \cdot) + \sum_{s'} P_{h'}^*(s'|\cdot, \cdot) \widehat{V}_{h'+1}^{r,\pi}(s') + 2b_{h'}(\cdot, \cdot).$$

The proof of Lemma 9 is deferred to Appendix E.3. Denote the event in Lemma 8 and Lemma 9 by \mathcal{E}_4 and \mathcal{E}_5 respectively. Then under event $\mathcal{E}_4 \cap \mathcal{E}_5$, we have for all policy $\pi \in \mathcal{C}(\Pi, \epsilon')$ and all linear reward function r with $r_h \in [-1, 1]$,

$$\begin{aligned}
0 &\leq \mathbb{E}_{s_1 \sim P_1^*(\cdot)}[\widehat{V}_1^{r,\pi}(s_1) - V_1^{r,\pi}(s_1)] \leq 2\mathbb{E}_{s_1 \sim P_1^*(\cdot)}[V_1^{b,\pi}(s_1)] \\
&\leq 2H\mathbb{E}_{s_1 \sim P_1(\cdot)}[V_1^{b/H,*}(s_1)] \leq 2C_{\text{lin}} \sqrt{\frac{d^3 H^6 R^2 \cdot \log(dKH R \mathcal{N}_{\Pi}(\epsilon')/\delta)}{K}} \leq \epsilon_0, \quad (17)
\end{aligned}$$

where $\epsilon_0 = \frac{\epsilon}{72BR\sqrt{Hd}}$. Here the first step comes from the left part of Lemma 9 and the second step is due to the right part of Lemma 9.

Note that in the proof of Lemma 13, we calculate the covering number of the function class $\{\widehat{V}_1^{r,\pi} : r \text{ is linear and } r_h \in [-1, 1]\}$ for any fixed π in (24). Then by Azuma-Hoeffding's inequality and (24), we have with probability at least $1 - \delta/6$ that for all policy $\pi \in \mathcal{C}(\Pi, \epsilon')$ and all linear reward function r with $r_h \in [-1, 1]$ that

$$\left| \mathbb{E}_{s_1 \sim P_1^*(\cdot)}[\widehat{V}_1^{r,\pi}(s_1)] - \frac{1}{K} \sum_{i=1}^K \widehat{V}_1^{r,\pi}(s_1^{i,\text{in}}) \right| \leq C_3 H \cdot \sqrt{\frac{\log(\mathcal{N}_{\Pi}(\epsilon') H K d R/\delta)}{K}} \leq \epsilon_0, \quad (18)$$

where $C_3 > 0$ is a universal constant.

Combining (17) and (18), we have with probability at least $1 - \delta/2$ that for all policy $\pi \in \mathcal{C}(\Pi, \epsilon')$ and all linear reward function r with $r_h \in [-1, 1]$

$$|\widehat{V}^{\pi}(r) - V^{r,\pi}| \leq 2\epsilon_0. \quad (19)$$

This implies that we can estimate the value function for all $\pi \in \mathcal{C}(\Pi, \epsilon')$ and linear reward function r with $r_h \in [-1, 1]$ up to estimation error $2\epsilon_0$.

Now we consider any policy $\pi \in \Pi$. Suppose that $\pi' \in \mathcal{C}(\Pi, \epsilon')$ satisfies that

$$\max_{s \in \mathcal{S}, h \in [H]} \|\pi_h(\cdot|s) - \pi'_h(\cdot|s)\|_1 \leq \epsilon'. \quad (20)$$

Then we can bound $|\widehat{V}^\pi(r) - \widehat{V}^{\pi'}(r)|$ and $|V^{r,\pi} - V^{r,\pi'}|$ for all linear reward function r with $r_h \in [-1, 1]$ respectively.

For $|V^{r,\pi} - V^{r,\pi'}|$, note that we have the following performance difference lemma:

Lemma 10. *For any policy π, π' and reward function r , we have*

$$V^{r,\pi'} - V^{r,\pi} = \sum_{h=1}^H \mathbb{E}_{\pi', P^*} \left[\langle Q_h^{r,\pi}(s_h, \cdot), \pi'_h(\cdot|s) - \pi_h(\cdot|s) \rangle \right].$$

The proof is deferred to Appendix E.4. Therefore from Lemma 10 we have

$$\begin{aligned} |V^{r,\pi'} - V^{r,\pi}| &\leq \sum_{h'=1}^H \mathbb{E}_{\pi, P^*} \left[\left| \langle Q_{h'}^{r^{h,j},\pi'}(s_{h'}, \cdot), \pi_{h'}(\cdot|s) - \pi'_{h'}(\cdot|s) \rangle \right| \right] \\ &\leq \sum_{h'=1}^H \mathbb{E}_{\pi, P^*} \left[\|\pi_{h'}(\cdot|s) - \pi'_{h'}(\cdot|s)\|_1 \right] \leq H\epsilon'. \end{aligned} \quad (21)$$

On the other hand, we have the following lemma to bound $|\widehat{V}^\pi(r) - \widehat{V}^{\pi'}(r)|$:

Lemma 11. *Suppose (20) holds and $\widehat{V}^\pi(r), \widehat{V}^{\pi'}(r)$ are calculated as in Algorithm 5. Then for all linear reward function r with $0 \leq r(\tau) \leq r_{\max}$, we have*

$$|\widehat{V}^\pi(r) - \widehat{V}^{\pi'}(r)| \leq \epsilon_{\text{cover}} := \frac{H\epsilon'}{\sqrt{dK} - 1} \cdot (dK)^{\frac{H}{2}}. \quad (22)$$

The proof is deferred to Appendix E.5.

Combining (19),(21) and (22), we have for all policy $\pi \in \Pi$ and linear reward function r with $0 \leq r(\tau) \leq r_{\max}$,

$$|\widehat{V}^\pi(r) - V^{r,\pi}| \leq 2\epsilon_0 + H\epsilon' + \epsilon_{\text{cover}}. \quad (23)$$

In particular, since $(\phi(\pi))_{h,j} = RV^{r^{h,j},\pi}$, we have for all policy $\pi \in \Pi$ and $h \in [H], j \in [d]$,

$$|(\phi(\pi))_{h,j} - (\widehat{\phi}(\pi))_{h,j}| \leq 2R\epsilon_0 + HR\epsilon' + R\epsilon_{\text{cover}}.$$

This implies that for all policy $\pi \in \Pi$ and any v defined in Lemma 6, we have

$$|(\phi(\pi)) - (\widehat{\phi}(\pi)), v| \leq 2BH\sqrt{d}(2R\epsilon_0 + HR\epsilon' + R\epsilon_{\text{cover}}).$$

The rest of the proof is the same as Theorem 1 and thus is omitted here. The only difference is that we need to show $\widehat{\pi}$ is a near-optimal policy with respect to \widehat{r} . This can be proved as follows:

$$\begin{aligned} V^{\widehat{r},*} - V^{\widehat{r},\widehat{\pi}} &= \left(V^{\widehat{r},*} - \widehat{V}^{\widehat{\pi}}(\widehat{r}) \right) + \left(\widehat{V}^{\widehat{\pi}}(\widehat{r}) - \widehat{V}^{\pi^*(\widehat{r})}(\widehat{r}) \right) + \left(\widehat{V}^{\pi^*(\widehat{r})}(\widehat{r}) - V^{\widehat{r},\widehat{\pi}} \right) \\ &\leq 4\epsilon_0 + 2H\epsilon' + 2\epsilon_{\text{cover}}, \end{aligned}$$

where the last step comes from (23) and the definition of $\widehat{\pi}$.

E.2 PROOF OF LEMMA 8

Here we outline the difference of the proof from Wang et al. (2020)[Lemma 3.2]. First, we also have the following concentration guarantee:

Lemma 12. Fix a policy π . Then with probability at least $1 - \delta$, we have for all $h \in [H]$ and $k \in [K]$,

$$\left\| \sum_{i=1}^k \phi_h^i \left(V_{h+1}^k(s_{h+1}^i) - \sum_{s' \in \mathcal{S}} P_h^*(s'|s_h^i, a_h^i) V_{h+1}^k(s') \right) \right\|_{\Lambda_h^{-1}} \leq O(dHR\sqrt{\log(dKHR/\delta)}).$$

The proof is almost the same as Lemma 13 and thus is omitted here. Then following the same arguments in Wang et al. (2020), we have the following inequality under Lemma 12:

$$\left| \phi_h(s, a)^\top w_h^k - \sum_{s' \in \mathcal{S}} P_h^*(s'|s, a) V_{h+1}^k(s') \right| \leq \beta_{\text{ex}} \|\phi_h(s, a)\|_{(\Lambda_h^k)^{-1}}.$$

Note that $V_{h+1}^k(s) \in [0, H - h]$ for all $s \in \mathcal{S}$, which implies that

$$0 \leq \sum_{s' \in \mathcal{S}} P_h^*(s'|s, a) V_{h+1}^k(s') + r_h^k(s, a) \leq H - h + 1.$$

Note that Clip is a contraction operator, which implies that

$$\begin{aligned} & \left| \text{Clip}_{[0, H-h+1]}((w_h^k)^\top \phi_h(s, a) + r_h^k(s, a)) - \left(\sum_{s' \in \mathcal{S}} P_h^*(s'|s, a) V_{h+1}^k(s') + r_h^k(s, a) \right) \right| \\ & \leq \left| (w_h^k)^\top \phi_h(s, a) - \sum_{s' \in \mathcal{S}} P_h^*(s'|s, a) V_{h+1}^k(s') \right| \leq \beta_{\text{ex}} \|\phi_h(s, a)\|_{(\Lambda_h^k)^{-1}}. \end{aligned}$$

On the other hand,

$$\left| \text{Clip}_{[0, H-h+1]}((w_h^k)^\top \phi_h(s, a) + r_h^k(s, a)) - \left(\sum_{s' \in \mathcal{S}} P_h^*(s'|s, a) V_{h+1}^k(s') + r_h^k(s, a) \right) \right| \leq H - h + 1.$$

This implies that

$$\left| \text{Clip}_{[0, H-h+1]}((w_h^k)^\top \phi_h(s, a) + r_h^k(s, a)) - \left(\sum_{s' \in \mathcal{S}} P_h^*(s'|s, a) V_{h+1}^k(s') + r_h^k(s, a) \right) \right| \leq b_h^k(s, a).$$

The rest of the proof is the same as Wang et al. (2020) and thus is omitted.

E.3 PROOF OF LEMMA 9

In the following discussion we will use ϕ_h^i to denote $\phi_h(s_h^i, a_h^i)$. First we need the following concentration lemma which is similar to Jin et al. (2020b)[Lemma B.3]:

Lemma 13. Fix a policy π . Then with probability at least $1 - \delta$, we have for all $h \in [H]$ and linear reward functions r with $r_h \in [-1, 1]$,

$$\left\| \sum_{i=1}^K \phi_h^i \left(\widehat{V}_{h+1}^{r, \pi}(s_{h+1}^i) - \sum_{s' \in \mathcal{S}} P_h^*(s'|s_h^i, a_h^i) \widehat{V}_{h+1}^{r, \pi}(s') \right) \right\|_{\Lambda_h^{-1}} \leq O(dHR\sqrt{\log(dKHR/\delta)}).$$

The proof is deferred to Appendix E.6. Then by union bound, we know with probability $1 - \delta/6$, we have for all policy $\pi \in \mathcal{C}(\Pi, \epsilon')$, $h \in [H]$ and linear reward functions r with $r_h \in [-1, 1]$ that

$$\left\| \sum_{i=1}^K \phi_h^i \left(\widehat{V}_{h+1}^{r, \pi}(s_{h+1}^i) - \sum_{s' \in \mathcal{S}} P_h^*(s'|s_h^i, a_h^i) \widehat{V}_{h+1}^{r, \pi}(s') \right) \right\|_{\Lambda_h^{-1}} \leq O(dHR\sqrt{\log(dKHR\mathcal{N}_\Pi(\epsilon')/\delta)}).$$

Let \mathcal{E}_6 denote the event that the above inequality holds. Then under \mathcal{E}_6 , following the same analysis in (Wang et al., 2020)[Lemma 3.1], we have for all policy $\pi \in \mathcal{C}(\Pi, \epsilon')$, $(s, a) \in \mathcal{S} \times \mathcal{A}$, $h \in [H]$ and linear reward functions r with $r_h \in [-1, 1]$ that

$$\left| \phi_h(s, a)^\top w_h^{r, \pi} - \sum_{s' \in \mathcal{S}} P_h^*(s'|s, a) \widehat{V}_{h+1}^{r, \pi}(s') \right| \leq \beta_{\text{pl}} \|\phi_h(s, a)\|_{\Lambda_h^{-1}}.$$

Form the contraction property of Clip and the fact that $\sum_{s' \in \mathcal{S}} P_h^*(s'|s, a) \widehat{V}_{h+1}^{r, \pi}(s') + r_h(s, a) \in [-(H-h+1), H-h+1]$, we know

$$\left| \text{Clip}_{[-(H-h+1), H-h+1]}((w_h^{r, \pi})^\top \phi_h(s, a) + r_h(s, a)) - \sum_{s' \in \mathcal{S}} P_h^*(s'|s, a) \widehat{V}_{h+1}^{r, \pi}(s') - r_h(s, a) \right| \leq b_h(s, a)$$

Therefore, under \mathcal{E}_6 we have

$$\widehat{Q}_h^{r, \pi}(s, a) \leq r_h(s, a) + \sum_{s'} P_h^*(s'|s, a) \widehat{V}_{h+1}^{r, \pi}(s') + 2b_h(s, a).$$

Now we only need to prove under \mathcal{E}_6 , for all policy $\pi \in \mathcal{C}(\Pi, \epsilon')$, $(s, a) \in \mathcal{S} \times \mathcal{A}$, $h \in [H]$ and linear reward function r with $r_h \in [-1, 1]$, we have $Q_h^{r, \pi}(s, a) \leq \widehat{Q}_h^{r, \pi}(s, a)$. We use induction to prove this. The claim holds obviously for $h = H + 1$. Then we suppose for some $h \in [H]$, we have $Q_{h+1}^{r, \pi}(s, a) \leq \widehat{Q}_{h+1}^{r, \pi}(s, a)$ for all policy $\pi \in \mathcal{C}(\Pi, \epsilon')$, $(s, a) \in \mathcal{S} \times \mathcal{A}$ and linear reward function r with $r_h \in [-1, 1]$. Then we have:

$$V_{h+1}^{r, \pi}(s) = \mathbb{E}_{a \sim \pi_{h+1}(\cdot|s)} [Q_{h+1}^{r, \pi}(s, a)] \leq \widehat{V}_{h+1}^{r, \pi}(s) = \mathbb{E}_{a \sim \pi_{h+1}(\cdot|s)} [\widehat{Q}_{h+1}^{r, \pi}(s, a)].$$

This implies that

$$\text{Clip}_{[-(H-h+1), H-h+1]}((w_h^{r, \pi})^\top \phi_h(s, a) + r_h(s, a)) + b_h(s, a) \geq \sum_{s' \in \mathcal{S}} P_h^*(s'|s, a) V_{h+1}^{r, \pi}(s') + r_h(s, a) = Q_h^{r, \pi}(s, a).$$

On the other hand we have

$$Q_h^{r, \pi}(s, a) \leq H - h + 1.$$

Therefore we have

$$Q_h^{r, \pi}(s, a) \leq \widehat{Q}_h^{r, \pi}(s, a).$$

By induction we can prove the lemma.

E.4 PROOF OF LEMMA 10

For any two policies π' and π , it follows from the definition of $V^{r, \pi'}$ and $V^{r, \pi}$ that

$$\begin{aligned} & V^{r, \pi'} - V^{r, \pi} \\ &= \mathbb{E}_{\pi', P^*} [r_1(s_1, a_1) + V_2^{r, \pi'}(s_2)] - \mathbb{E}_{\pi', P^*} [V_1^{r, \pi}(s_1)] \\ &= \mathbb{E}_{\pi', P^*} [V_2^{r, \pi'}(s_2) - (V_1^{r, \pi}(s_1) - r_1(s_1, a_1))] \\ &= \mathbb{E}_{\pi', P^*} [V_2^{r, \pi'}(s_2) - V_2^{r, \pi}(s_2)] + \mathbb{E}_{\pi', P^*} [Q_1^{r, \pi}(s_1, a_1) - V_1^{r, \pi}(s_1)] \\ &= \mathbb{E}_{\pi', P^*} [V_2^{r, \pi'}(s_2) - V_2^{r, \pi}(s_2)] + \mathbb{E}_{\pi', P^*} [\langle Q_1^{r, \pi}(s_1, \cdot), \pi'_1(\cdot|s_1) - \pi_1(\cdot|s_1) \rangle] \\ &= \dots = \sum_{h=1}^H \mathbb{E}_{\pi', P^*} [\langle Q_h^{r, \pi}(s_h, \cdot), \pi'_h(\cdot|s) - \pi_h(\cdot|s) \rangle]. \end{aligned}$$

This concludes our proof.

E.5 PROOF OF LEMMA 11

For any $h' \in [H]$, suppose $\max_{s \in \mathcal{S}} |\widehat{V}_{h'+1}^{r, \pi}(s) - \widehat{V}_{h'+1}^{r, \pi'}(s)| \leq \epsilon_{h'+1}$, then for any $s \in \mathcal{S}$, $a \in \mathcal{A}$, we have

$$\begin{aligned} |\widehat{Q}_{h'}^{r, \pi}(s, a) - \widehat{Q}_{h'}^{r, \pi'}(s, a)| &\leq |(w_{h'}^{r, \pi} - w_{h'}^{r, \pi'})^\top \phi_{h'}(s, a)| \\ &\leq \epsilon_{h'+1} \sum_{i=1}^K |\phi_{h'}(s, a)^\top (\Lambda_{h'})^{-1} \phi_{h'}(s_{h'}^i, a_{h'}^i)| \end{aligned}$$

$$\leq \epsilon_{h'+1} \sqrt{\left[\sum_{i=1}^K \|\phi_{h'}(s, a)\|_{(\Lambda_{h'})^{-1}}^2 \right] \cdot \left[\sum_{i=1}^K \|\phi_{h'}(s_{h'}^i, a_{h'}^i)\|_{(\Lambda_{h'})^{-1}}^2 \right]} \leq \epsilon_{h'+1} \sqrt{dK}.$$

Here the final step is comes from the auxiliary Lemma 14 and the fact that $\Lambda_{h'} \geq R^2 I$ and thus $\sum_{i=1}^K \|\phi_{h'}(s, a)\|_{(\Lambda_{h'})^{-1}}^2 \leq \sum_{i=1}^K 1 \leq K$.

Therefore we have

$$\epsilon_{h'} := \max_{s \in \mathcal{S}} |\widehat{V}_{h'}^{r, \pi}(s) - \widehat{V}_{h'}^{r, \pi'}(s)| \leq H\epsilon' + \sqrt{dK}\epsilon_{h'+1}.$$

Note that $\epsilon_{H+1} = 0$, therefore we have

$$\epsilon_1 \leq \frac{H\epsilon'}{\sqrt{dK} - 1} \cdot (dK)^{\frac{H}{2}},$$

This concludes our proof.

E.6 PROOF OF LEMMA 13

The proof is almost the same as Jin et al. (2020b)[Lemma B.3] except that the function class of $V_h^{r, \pi}$ is different. Therefore we only need to bound the covering number $\mathcal{N}_{\mathcal{V}}(\epsilon)$ of $V_h^{r, \pi}$ where the distance is defined as $\text{dist}(V, V') = \sup_s |V(s) - V'(s)|$. Note that $V_h^{r, \pi}$ belongs to the following function class:

$$\mathcal{V} = \left\{ V_{w, A}(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} \left[\text{Clip}_{[-(H-h+1), H-h+1]} \left(\text{Clip}_{[-(H-h+1), H-h+1]}(w^\top \phi_{h'}(s, a)) \right. \right. \right. \\ \left. \left. \left. + \text{Clip}_{[0, 2(H-h+1)]}(\|\phi(s, a)\|_A) \right) \right], \forall s \in \mathcal{S} \right\},$$

where the parameters (w, A) satisfy $\|w\| \leq 2H\sqrt{dK/\lambda_{\text{pl}}}$, $\|A\| \leq \beta_{\text{pl}}^2 \lambda_{\text{pl}}^{-1}$.

Note that for any $V_{w_1, A_1}, V_{w_2, A_2} \in \mathcal{V}$, we have

$$\begin{aligned} \text{dist}(V_{w_1, A_1}, V_{w_2, A_2}) &\leq \sup_{s, a} \left| \left[\text{Clip}_{[-(H-h+1), H-h+1]}(w_1^\top \phi_{h'}(s, a)) + \text{Clip}_{[0, 2(H-h+1)]}(\|\phi(s, a)\|_{A_1}) \right] \right. \\ &\quad \left. - \left[\text{Clip}_{[-(H-h+1), H-h+1]}(w_2^\top \phi_{h'}(s, a)) + \text{Clip}_{[0, 2(H-h+1)]}(\|\phi(s, a)\|_{A_2}) \right] \right| \\ &\leq \sup_{s, a} \left| \text{Clip}_{[-(H-h+1), H-h+1]}(w_1^\top \phi_{h'}(s, a)) - \text{Clip}_{[-(H-h+1), H-h+1]}(w_2^\top \phi_{h'}(s, a)) \right| \\ &\quad + \sup_{s, a} \left| \text{Clip}_{[0, 2(H-h+1)]}(\|\phi(s, a)\|_{A_1}) - \text{Clip}_{[0, 2(H-h+1)]}(\|\phi(s, a)\|_{A_2}) \right| \\ &\leq R \sup_{\|\phi\| \leq 1} \left| (w_1 - w_2)^\top \phi \right| + R \sup_{\|\phi\| \leq 1} \sqrt{\left| \phi^\top (A_1 - A_2) \phi \right|} \\ &\leq R(\|w_1 - w_2\| + \sqrt{\|A_1 - A_2\|_F}), \end{aligned}$$

where the first and third step utilize the contraction property of Clip. Let \mathcal{C}_w be the $\epsilon/(2R)$ -cover of $\{w \in \mathbb{R}^d : \|w\| \leq 2r_{\max} \sqrt{dK/\lambda_{\text{pl}}}\}$ w.r.t. ℓ_2 -norm and \mathcal{C}_A be the $(\epsilon/2R)$ -cover of $\{A \in \mathbb{R}^{d \times d} : \|A\| \leq \beta_{\text{pl}}^2 \lambda_{\text{pl}}^{-1}\}$ w.r.t. the Frobenius norm, then from the literature Jin et al. (2020b)[Lemma D.5], we have

$$\mathcal{N}_{\mathcal{V}}(\epsilon) \leq \log |\mathcal{C}_w| + \log |\mathcal{C}_A| \leq d \log \left(1 + 8\sqrt{dK r_{\max}^2 R^2 / (\lambda_{\text{pl}} \epsilon^2)} \right) + d^2 \log \left[1 + 8d^{1/2} \beta_{\text{pl}}^2 R^2 / (\lambda_{\text{pl}} \epsilon^2) \right]. \quad (24)$$

The rest of the proof follows Jin et al. (2020b)[Lemma B.3] directly so we omit it here.

E.7 PROOF OF PROPOSITION 1

First consider ζ and ζ' which satisfies:

$$\|\zeta_h - \zeta'_h\| \leq \epsilon_z, \forall h \in [H].$$

Then we know for any $h \in [H], s \in \mathcal{S}, a \in \mathcal{A}$,

$$|\zeta_h^\top \phi_h(s, a) - (\zeta'_h)^\top \phi_h(s, a)| \leq \epsilon_z R. \quad (25)$$

Now fix any $h \in [H]$ and $s \in \mathcal{S}$. To simplify writing, we use $x(a)$ and $x'(a)$ to denote $\zeta_h^\top \phi_h(s, a)$ and $(\zeta'_h)^\top \phi_h(s, a)$ respectively. Without loss of generality, we assume $\sum_a \exp(x(a)) \leq \sum_a \exp(x'(a))$. Then from (25) we have

$$\sum_a \exp(x(a)) \leq \sum_a \exp(x'(a)) \leq \exp(\epsilon_z R) \sum_a \exp(x(a)).$$

Note that we have

$$\begin{aligned} \|\pi_h^\zeta(\cdot|s) - \pi_h^{\zeta'}(\cdot|s)\|_1 &= \sum_a \left| \frac{\exp(x(a))}{\sum_{a'} \exp(x(a'))} - \frac{\exp(x'(a))}{\sum_{a'} \exp(x'(a'))} \right| \\ &= \frac{\sum_a \left| \exp(x(a)) \sum_{a'} \exp(x'(a')) - \exp(x'(a)) \sum_{a'} \exp(x(a')) \right|}{\sum_{a'} \exp(x(a')) \cdot \sum_{a'} \exp(x'(a'))}. \end{aligned}$$

For any $a \in \mathcal{A}$, if $\exp(x(a)) \sum_{a'} \exp(x'(a')) - \exp(x'(a)) \sum_{a'} \exp(x(a')) \geq 0$, then

$$\begin{aligned} &\left| \exp(x(a)) \sum_{a'} \exp(x'(a')) - \exp(x'(a)) \sum_{a'} \exp(x(a')) \right| \\ &\leq \exp(\epsilon_z R) \exp(x(a)) \sum_{a'} \exp(x(a')) - \exp(-\epsilon_z R) \exp(x(a)) \sum_{a'} \exp(x'(a')) \\ &= (\exp(\epsilon_z R) - \exp(-\epsilon_z R)) \exp(x(a)) \sum_{a'} \exp(x(a')). \end{aligned}$$

Otherwise, we have

$$\begin{aligned} &\left| \exp(x(a)) \sum_{a'} \exp(x'(a')) - \exp(x'(a)) \sum_{a'} \exp(x(a')) \right| \\ &\leq \exp(\epsilon_z R) \exp(x(a)) \sum_{a'} \exp(x(a')) - \exp(x(a)) \sum_{a'} \exp(x(a')) \\ &= (\exp(\epsilon_z R) - 1) \exp(x(a)) \sum_{a'} \exp(x(a')). \end{aligned}$$

Therefore we have

$$\|\pi_h^\zeta(\cdot|s) - \pi_h^{\zeta'}(\cdot|s)\|_1 \leq \frac{(\exp(\epsilon_z R) - \exp(-\epsilon_z R)) \sum_a \exp(x(a)) \sum_{a'} \exp(x(a'))}{\sum_{a'} \exp(x(a')) \cdot \sum_{a'} \exp(x'(a'))} \leq \exp(2\epsilon_z R) - 1.$$

This implies that for any $\epsilon \leq 1$,

$$\mathcal{N}_\Pi(\epsilon) \leq \left(\mathcal{N}_{\mathbb{B}(d, W)} \left(\frac{\ln 2}{2R} \epsilon \right) \right)^H \leq \left(\frac{12WR}{\epsilon} \right)^{Hd},$$

where the first step uses $\exp(x) - 1 \leq x/\ln 2$ when $x \leq \ln 2$. This concludes our proof.

E.8 PROOF OF PROPOSITION 2

First we consider the following entropy-regularized RL problem where we try to maximize the following objective for some $\alpha > 0$:

$$\max_{\pi} V_{\alpha}(r^*, \pi) := \mathbb{E}_{\pi, P^*} \left[\sum_{h=1}^H r_h^*(s_h, a_h) - \alpha \log \pi_h(a_h | s_h) \right].$$

From the literature (Nachum et al., 2017; Cen et al., 2022), we know that we can define corresponding optimal regularized value function and Q function as follows:

$$\begin{aligned} Q_{\alpha,h}^*(s,a) &= r_h^*(s,a) + \mathbb{E}_{s_{h+1} \sim P_h^*(\cdot|s,a)} [V_{\alpha,h+1}^*], \\ V_{\alpha,h}^*(s) &= \max_{\pi_h} \mathbb{E}_{a_h \sim \pi_h(\cdot|s)} [Q_{\alpha,h}^*(s,a_h) - \alpha \log \pi_h(a_h|s)], \end{aligned}$$

where $V_{\alpha,H+1}^*(s) = 0$ for all $s \in \mathcal{S}$. Note that we have $V_{\alpha,h}^*(s) \leq H(1 + \alpha \log |\mathcal{A}|)$ for all $s \in \mathcal{S}$ and $h \in [H]$. The global optimal regularized policy is therefore

$$\pi_{\alpha,h}^*(a|s) = \frac{\exp(Q_{\alpha,h}^*(s,a)/\alpha)}{\sum_{a'} \exp(Q_{\alpha,h}^*(s,a')/\alpha)}.$$

In particular, in linear MDPs, we have

$$Q_{\alpha,h}^*(s,a) = \phi_h(s,a)^\top \left(\theta_h^* + \int_{s \in \mathcal{S}} \mu_h^*(s) V_{\alpha,h+1}^*(s) ds \right).$$

Therefore, $Q_{\alpha,h}^*(s,a) = \phi_h(s,a)^\top w_{\alpha,h}^*$ where

$$\|w_{\alpha,h}^*\| \leq B + H(1 + \alpha \log |\mathcal{A}|) \sqrt{d}.$$

This implies that π_α^* belongs to the log-linear policy class Π with $W = (B + H(1 + \alpha \log |\mathcal{A}|) \sqrt{d})/\alpha$.

On the other hand, let π_g denote the global unregularized optimal policy, then

$$\begin{aligned} V^*(r^*, \pi_g) - \max_{\pi \in \Pi} V(r^*, \pi) &\leq V^*(r^*, \pi_g) - V(r^*, \pi_\alpha^*) \\ &= (V^*(r^*, \pi_g) - V_\alpha^*(r^*, \pi_g)) + (V_\alpha^*(r^*, \pi_g) - V_\alpha^*(r^*, \pi_\alpha^*)) + (V_\alpha^*(r^*, \pi_\alpha^*) - V^*(r^*, \pi_\alpha^*)) \\ &\leq V_\alpha^*(r^*, \pi_\alpha^*) - V^*(r^*, \pi_\alpha^*) \leq \alpha H \log |\mathcal{A}|. \end{aligned}$$

Therefore we only need to let $\alpha = \frac{\epsilon}{H \log |\mathcal{A}|}$ to ensure $V(r^*, \pi_g) - \max_{\pi \in \Pi} V(r^*, \pi) \leq \epsilon$.

F PROOF OF THEOREM 3

First from performance difference lemma (Lemma 10), we have

$$\begin{aligned} V^{r^*, \hat{\pi}} - V^{r^*, *} &= \sum_{h=1}^H \mathbb{E}_{s_h \sim d_h^{\hat{\pi}}} [Q_h^*(s_h, \hat{\pi}) - Q_h^*(s_h, \pi^*)] \\ &= \sum_{h=1}^H \mathbb{E}_{s_h \sim d_h^{\hat{\pi}}} [Q_h^*(s_h, \hat{\pi}) - \hat{A}_h(s_h, \hat{\pi})] + \mathbb{E}_{s_h \sim d_h^{\hat{\pi}}} [\hat{A}_h(s_h, \hat{\pi}) - \hat{A}_h(s_h, \pi^*)] \\ &\quad + \mathbb{E}_{s_h \sim d_h^{\hat{\pi}}} [\hat{A}_h(s_h, \pi^*) - Q_h^*(s_h, \pi^*)] \\ &\geq \sum_{h=1}^H \mathbb{E}_{s_h \sim d_h^{\hat{\pi}}} [Q_h^*(s_h, \hat{\pi}) - \hat{A}_h(s_h, \hat{\pi})] + \mathbb{E}_{s_h \sim d_h^{\hat{\pi}}} [\hat{A}_h(s_h, \pi^*) - Q_h^*(s_h, \pi^*)] \\ &= \sum_{h=1}^H \mathbb{E}_{s_h \sim d_h^{\hat{\pi}}} [A_h^*(s_h, \hat{\pi}) - \hat{A}_h(s_h, \hat{\pi}) + \hat{A}_h(s_h, \pi^*) - A_h^*(s_h, \pi^*)] \\ &= \sum_{h=1}^H \mathbb{E}_{s_h \sim d_h^{\hat{\pi}}} [\langle \phi_h(s_h, \hat{\pi}), \xi_h^* - \hat{\xi}_h \rangle - \langle \phi_h(s_h, \pi^*), \xi_h^* - \hat{\xi}_h \rangle] \\ &= \sum_{h=1}^H \mathbb{E}_{s_h \sim d_h^{\hat{\pi}}} [\langle \phi_h(s_h, \hat{\pi}) - \phi_h(s_h, \pi^*), \xi_h^* - \hat{\xi}_h \rangle] \\ &\geq - \sum_{h=1}^H \|\mathbb{E}_{s_h \sim d_h^{\hat{\pi}}} [\phi_h(s_h, \hat{\pi}) - \phi_h(s_h, \pi^*)]\|_{\Sigma_{h,N+1}^{-1}} \cdot \|\xi_h^* - \hat{\xi}_h\|_{\Sigma_{h,N+1}}. \end{aligned} \tag{26}$$

Next we will bound $\|\mathbb{E}_{s_h \sim d_{\widehat{\pi}}^h}[\phi_h(s_h, \widehat{\pi}) - \phi_h(s_h, \pi^*)]\|_{\Sigma_{h,N+1}^{-1}}$ and $\|\xi^* - \widehat{\xi}\|_{\Sigma_{h,N+1}}$ respectively. First for $\|\mathbb{E}_{s_h \sim d_{\widehat{\pi}}^h}[\phi_h(s_h, \widehat{\pi}) - \phi_h(s_h, \pi^*)]\|_{\Sigma_{h,N+1}^{-1}}$, notice that $\Sigma_{h,N+1} \succeq \Sigma_{h,n}$ for all $n \in [N+1]$, which implies

$$\begin{aligned} \|\mathbb{E}_{s_h \sim d_{\widehat{\pi}}^h}[\phi_h(s_h, \widehat{\pi}) - \phi_h(s_h, \pi^*)]\|_{\Sigma_{h,N+1}^{-1}} &\leq \frac{1}{N} \sum_{n=1}^N \|\mathbb{E}_{s_h \sim d_{\widehat{\pi}}^h}[\phi_h(s_h, \widehat{\pi}) - \phi_h(s_h, \pi^*)]\|_{\Sigma_{h,n}^{-1}} \\ &\leq \frac{1}{N} \sum_{n=1}^N \|\mathbb{E}_{s_h \sim \pi^{h,n,0}}[\phi_h(s_h, \pi^{h,n,0}) - \phi_h(s_h, \pi^{h,n,1})]\|_{\Sigma_{h,n}^{-1}} \\ &\leq \frac{1}{\sqrt{N}} \sqrt{\sum_{n=1}^N \|\mathbb{E}_{s_h \sim \pi^{h,n,0}}[\phi_h(s_h, \pi^{h,n,0}) - \phi_h(s_h, \pi^{h,n,1})]\|_{\Sigma_{h,n}^{-1}}^2} \\ &\leq \sqrt{\frac{2d \log(1+N/d)}{N}}, \end{aligned} \quad (27)$$

where the third step comes from the definition of $\pi^{h,n,0}$ and $\pi^{h,n,1}$ and the last step comes from Elliptical Potential Lemma (Lemma 1) and the fact that $\lambda \geq 4R^2$.

For $\|\xi_h^* - \widehat{\xi}_h\|_{\Sigma_{h,N+1}}$, let $\widetilde{\Sigma}_{h,n}$ denote $\lambda I + \sum_{i=1}^{n-1} (\phi_h(s^{h,n}, a^{h,n,0}) - \phi_h(s^{h,n}, a^{h,n,1}))(\phi_h(s^{h,n}, a^{h,n,0}) - \phi_h(s^{h,n}, a^{h,n,1}))^\top$. Then similar to Lemma 3, we have with probability at least $1 - \delta/2$,

$$\|\xi_h^* - \widehat{\xi}_h\|_{\Sigma_{h,N+1}}^2 \leq 2\|\xi_h^* - \widehat{\xi}_h\|_{\widetilde{\Sigma}_{h,N+1}}^2 + 2C_{\text{CON}} d R^2 B^2 \log(N/\delta). \quad (28)$$

On the other hand, similar to Lemma 2, MLE guarantees us that with probability at least $1 - \delta/2$,

$$\|\widehat{\xi}_h - \xi_h^*\|_{\widetilde{\Sigma}_{h,N+1}} \leq 2C_{\text{MLE}} \cdot \sqrt{\kappa_{\text{adv}}^2 (d + \log(1/\delta)) + \lambda B^2}, \quad (29)$$

where $\kappa_{\text{adv}} = 2 + \exp(2B_{\text{adv}}) + \exp(-2B_{\text{adv}})$.

Therefore combining (28) and (29), we have with probability at least $1 - \delta$,

$$\|\xi^* - \widehat{\xi}_h\|_{\Sigma_{h,N+1}} \leq \mathcal{O}(\kappa_{\text{adv}} B R \sqrt{\lambda d \log(N/\delta)}). \quad (30)$$

Thus combining (26), (27) and (30) via union bound, we have $V^*(r^*) - V(r^*, \widehat{\pi}) \leq \epsilon$ with probability at least $1 - \delta$ as long as

$$N \geq \widetilde{\mathcal{O}}\left(\frac{\lambda \kappa_{\text{adv}}^2 B^2 R^2 H^2 d^2 \log(1/\delta)}{\epsilon^2}\right).$$

G AUXILIARY LEMMAS

Lemma 14 (Jin et al. (2020b)[Lemma D.1]). *Let $\Lambda = \lambda I + \sum_{i=1}^K \phi_i \phi_i^\top$ where $\phi_i \in \mathbb{R}^d$ and $\lambda > 0$, then we have $\sum_{i=1}^K \phi_i^\top \Lambda^{-1} \phi_i \leq d$.*