

Appendices

A. Omitted Proof in Section 4

A.1. PROOF OF LEMMA 1

Proof. By the optimality of $\hat{\mathbf{B}}_t$ and $\hat{\mathbf{W}}_t = [\hat{\mathbf{w}}_{t,1}, \dots, \hat{\mathbf{w}}_{t,M}]$, we know that $\sum_{i=1}^M \left\| \mathbf{y}_{t-1,i} - \mathbf{X}_{t-1,i}^\top \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} \right\|_2^2 \leq \sum_{i=1}^M \left\| \mathbf{y}_{t-1,i} - \mathbf{X}_{t-1,i}^\top \mathbf{B} \mathbf{w}_i \right\|_2^2$. Since $\mathbf{y}_{t-1,i} = \mathbf{X}_{t-1,i}^\top \mathbf{B} \mathbf{w}_i + \boldsymbol{\eta}_{t-1,i}$, we have

$$\sum_{i=1}^M \left\| \mathbf{X}_{t-1,i}^\top \left(\hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right) \right\|_2^2 \leq 2 \sum_{i=1}^M \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \left(\hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right). \quad (28)$$

We firstly analyse the non-trivial setting where $d \geq 2k$. Note that both $\Theta = \mathbf{B} \mathbf{W}$ and $\hat{\Theta}_t = \hat{\mathbf{B}}_t \hat{\mathbf{W}}_t$ are low-rank matrix with rank upper bounded by k , which indicates that $\text{rank}(\hat{\Theta}_t - \Theta) \leq 2k$. In that case, we can write $\hat{\Theta}_t - \Theta = \mathbf{U}_t \mathbf{R}_t = [\mathbf{U}_t \mathbf{r}_{t,1}, \mathbf{U}_t \mathbf{r}_{t,2}, \dots, \mathbf{U}_t \mathbf{r}_{t,M}]$, where $\mathbf{U}_t \in \mathbb{R}^{d \times 2k}$ is an orthonormal matrix with $\|\mathbf{U}_t\|_F = \sqrt{2k}$, and $\mathbf{R}_t \in \mathbb{R}^{2k \times M}$ satisfies $\|\mathbf{r}_{t,i}\|_2 \leq \sqrt{k}$. In other words, we can write $\hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i = \mathbf{U}_t \mathbf{r}_{t,i}$ for certain \mathbf{U}_t and $\mathbf{r}_{t,i}$.

Define $\tilde{\mathbf{V}}_{t-1,i}(\lambda) \stackrel{\text{def}}{=} (\mathbf{U}_t^\top \mathbf{X}_{t-1,i}) (\mathbf{U}_t^\top \mathbf{X}_{t-1,i})^\top + \lambda \mathbf{I}$. We have:

$$\sum_{i=1}^M \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2 \quad (29)$$

$$= \sum_{i=1}^M \left\| \mathbf{X}_{t-1,i}^\top \left(\hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right) \right\|_2^2 + \sum_{i=1}^M \lambda \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_2^2 \quad (30)$$

$$\leq 2 \sum_{i=1}^M \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \left(\hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right) + 4M\lambda \quad (31)$$

$$= 2 \sum_{i=1}^M \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \mathbf{U}_t \mathbf{r}_{t,i} + 4M\lambda \quad (32)$$

$$\leq 2 \sum_{i=1}^M \left\| \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \mathbf{U}_t \right\|_{\mathbf{V}_{t-1,i}^{-1}(\lambda)} \|\mathbf{r}_{t,i}\|_{\mathbf{V}_{t-1,i}(\lambda)} + 4M\lambda \quad (33)$$

$$\leq 2 \sqrt{\sum_{i=1}^M \left\| \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \mathbf{U}_t \right\|_{\mathbf{V}_{t-1,i}^{-1}(\lambda)}^2} \sqrt{\sum_{i=1}^M \|\mathbf{r}_{t,i}\|_{\mathbf{V}_{t-1,i}(\lambda)}^2} + 4M\lambda \quad (34)$$

$$= 2 \sqrt{\sum_{i=1}^M \left\| \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \mathbf{U}_t \right\|_{\mathbf{V}_{t-1,i}^{-1}(\lambda)}^2} \sqrt{\sum_{i=1}^M \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2} + 4M\lambda \quad (35)$$

Eqn 31 is due to Eqn 28, $\left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} \right\| \leq 1$ and $\|\mathbf{B} \mathbf{w}_i\| \leq 1$. Eqn 34 is due to Cauchy-Schwarz inequality. Eqn 35 is from

$$\sum_{i=1}^M \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2 = \sum_{i=1}^M \left\| \mathbf{U}_t \mathbf{r}_{t,i} \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2 = \sum_{i=1}^M \|\mathbf{r}_{t,i}\|_{\mathbf{U}_t^\top \tilde{\mathbf{V}}_{t-1,i}(\lambda) \mathbf{U}_t}^2 = \sum_{i=1}^M \|\mathbf{r}_{t,i}\|_{\mathbf{V}_{t-1,i}(\lambda)}^2.$$

The main problem is how to bound $\left\| \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \mathbf{U}_t \right\|_{\mathbf{V}_{t-1,i}^{-1}(\lambda)} = \left\| \sum_{n=1}^{t-1} \eta_{n,i} \mathbf{U}_n^\top \mathbf{x}_{n,i} \right\|_{\mathbf{V}_{t-1,i}^{-1}(\lambda)}$. Note that for a fixed $\mathbf{U}_t = \bar{\mathbf{U}}$, we can regard $\bar{\mathbf{U}}^\top \mathbf{x}_{n,i} \in \mathbb{R}^k$ as the corresponding ‘‘action’’ chosen in step t . With this observation, if \mathbf{U}_t is fixed, we can bound this term following the arguments of the self-normalized bound for vector-valued martingales (Abbasi-Yadkori et al., 2011).

Lemma 2. For a fixed $\bar{\mathbf{U}}$, define $\bar{\mathbf{V}}_{t,i}(\lambda) \stackrel{\text{def}}{=} (\bar{\mathbf{U}}^\top \mathbf{X}_{t,i}) (\bar{\mathbf{U}}^\top \mathbf{X}_{t,i})^\top + \lambda \mathbf{I}$, then any $\delta > 0$, with probability at least $1 - \delta$, for all $t \geq 0$,

$$\sum_{i=1}^M \|\bar{\mathbf{U}}^\top \mathbf{X}_{t,i} \boldsymbol{\eta}_{t,i}\|_{\bar{\mathbf{V}}_{t,i}^{-1}}^2 \quad (36)$$

$$\leq 2 \log \left(\frac{\prod_{i=1}^M (\det(\bar{\mathbf{V}}_{t,i})^{1/2} \det(\lambda \mathbf{I})^{-1/2})}{\delta} \right). \quad (37)$$

We defer the proof of Lemma 2 to Appendix A.3. We set $\lambda = 1$. By Lemma 2, we know that for a fixed $\bar{\mathbf{U}}$, with probability at least $1 - \delta_1$,

$$\sum_{i=1}^M \left\| \sum_{n=1}^{t-1} \eta_{n,i} \bar{\mathbf{U}}^\top x_{n,i} \right\|_{\bar{\mathbf{V}}_{t,i}^{-1}(\lambda)}^2 \leq 2 \log \left(\frac{\prod_{i=1}^M \det(\bar{\mathbf{V}}_{t,i}(\lambda))^{1/2} \det(\lambda \mathbf{I})^{-1/2}}{\delta_1} \right) \leq 2Mk + 2 \log(1/\delta_1). \quad (38)$$

The above analysis shows that we can bound $\|\boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \mathbf{U}_t\|_{\bar{\mathbf{V}}_{t-1,i}^{-1}(\lambda)}$ if \mathbf{U}_t is fixed as $\bar{\mathbf{U}}$. Following this idea, we prove the lemma by the construction of ϵ -net over all possible \mathbf{U}_t . To apply the trick of ϵ -net, we need to slightly modify the derivation of Eqn 29. For a fixed matrix $\bar{\mathbf{U}} \in \mathbb{R}^{d \times 2k}$, we have

$$\sum_{i=1}^M \|\hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i\|_{\bar{\mathbf{V}}_{t-1,i}^{-1}(\lambda)}^2 \quad (39)$$

$$\leq 2 \sum_{i=1}^M \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \mathbf{U}_t \mathbf{r}_{t,i} + 4M\lambda \quad (40)$$

$$= 2 \sum_{i=1}^M \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \bar{\mathbf{U}} \mathbf{r}_{t,i} + 2 \sum_{i=1}^M \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top (\mathbf{U}_t - \bar{\mathbf{U}}) \mathbf{r}_{t,i} + 4M\lambda \quad (41)$$

$$\leq 2 \sum_{i=1}^M \|\boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \bar{\mathbf{U}}\|_{\bar{\mathbf{V}}_{t-1,i}^{-1}(\lambda)} \|\mathbf{r}_{t,i}\|_{\bar{\mathbf{V}}_{t-1,i}(\lambda)} + 2 \sum_{i=1}^M \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top (\mathbf{U}_t - \bar{\mathbf{U}}) \mathbf{r}_{t,i} + 4M\lambda \quad (42)$$

$$= 2 \sum_{i=1}^M \|\boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \bar{\mathbf{U}}\|_{\bar{\mathbf{V}}_{t-1,i}^{-1}(\lambda)} \|\mathbf{r}_{t,i}\|_{\bar{\mathbf{V}}_{t-1,i}(\lambda)} + 2 \sum_{i=1}^M \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top (\mathbf{U}_t - \bar{\mathbf{U}}) \mathbf{r}_{t,i} \quad (43)$$

$$+ 2 \sum_{i=1}^M \|\boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \bar{\mathbf{U}}\|_{\bar{\mathbf{V}}_{t-1,i}^{-1}(\lambda)} \left(\|\mathbf{r}_{t,i}\|_{\bar{\mathbf{V}}_{t-1,i}(\lambda)} - \|\mathbf{r}_{t,i}\|_{\bar{\mathbf{V}}_{t-1,i}(\lambda)} \right) + 4M\lambda \quad (44)$$

$$\leq 2 \sqrt{\sum_{i=1}^M \|\boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \bar{\mathbf{U}}\|_{\bar{\mathbf{V}}_{t-1,i}^{-1}(\lambda)}^2} \sqrt{\sum_{i=1}^M \|\hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i\|_{\bar{\mathbf{V}}_{t-1,i}^{-1}(\lambda)}^2} + 2 \sum_{i=1}^M \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top (\mathbf{U}_t - \bar{\mathbf{U}}) \mathbf{r}_{t,i} \quad (45)$$

$$+ 2 \sum_{i=1}^M \|\boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \bar{\mathbf{U}}\|_{\bar{\mathbf{V}}_{t-1,i}^{-1}(\lambda)} \left(\|\mathbf{r}_{t,i}\|_{\bar{\mathbf{V}}_{t-1,i}(\lambda)} - \|\mathbf{r}_{t,i}\|_{\bar{\mathbf{V}}_{t-1,i}(\lambda)} \right) + 4M\lambda \quad (46)$$

Eqn 40, 42 and 45 follow the same idea of Eqn 32, 33 and 35.

We construct an ϵ -net \mathcal{E} in Frobenius norm over the matrix set $\{\mathbf{U} \in \mathbb{R}^{d \times 2k} : \|\mathbf{U}\|_F \leq k\}$. It is not hard to see that $|\mathcal{E}| \leq \left(\frac{6\sqrt{2k}}{\epsilon}\right)^{2kd}$. By the union bound over all possible $\bar{\mathbf{U}} \in \mathcal{E}$, we know that with probability $1 - |\mathcal{E}|\delta_1$, Eqn 38 holds for any $\bar{\mathbf{U}} \in \mathcal{E}$. For each \mathbf{U}_t , we choose an $\bar{\mathbf{U}} \in \mathcal{E}$ with $\|\mathbf{U}_t - \bar{\mathbf{U}}\|_F \leq \epsilon$, and we have

$$2 \sqrt{\sum_{i=1}^M \|\boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \bar{\mathbf{U}}\|_{\bar{\mathbf{V}}_{t-1,i}^{-1}(\lambda)}^2} \leq 2\sqrt{2Mk + 2 \log(1/\delta_1)} \quad (47)$$

Since $\|\mathbf{U}_t - \bar{\mathbf{U}}\|_F \leq \epsilon$, we have

$$2 \sum_{i=1}^M \left\| \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \bar{\mathbf{U}} \right\|_{\tilde{\mathbf{V}}_{t-1,i}^{-1}(\lambda)} \left(\|\mathbf{r}_{t,i}\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)} - \|\mathbf{r}_{t,i}\|_{\mathbf{V}_{t-1,i}(\lambda)} \right) \leq 2\sqrt{Mk\epsilon(2Mk + 2\log(1/\delta_1))}. \quad (48)$$

For the term $2 \sum_{i=1}^M \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top (\mathbf{U}_t - \bar{\mathbf{U}}) \mathbf{r}_{t,i}$, the following inequality holds for any step $t \in [T]$ with probability $1 - MT\delta_2$,

$$2 \sum_{i=1}^M \boldsymbol{\eta}_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top (\mathbf{U}_t - \bar{\mathbf{U}}) \mathbf{r}_{t,i} \leq 2 \sum_{i=1}^M \|\boldsymbol{\eta}_{t-1,i}\|_2 \left\| \mathbf{X}_{t-1,i}^\top (\mathbf{U}_t - \bar{\mathbf{U}}) \mathbf{r}_{t,i} \right\|_2 \quad (49)$$

$$\leq 2 \sum_{i=1}^M \|\boldsymbol{\eta}_{t-1,i}\|_2 \sqrt{kT\epsilon} \quad (50)$$

$$\leq 2M\sqrt{2\log(2/\delta_2)kT^2\epsilon} \quad (51)$$

The last inequality follows from the fact that $|\eta_{n,i}| \leq \sqrt{2\log(2/\delta_2)}$ with probability $1 - \delta_2$ for fixed n, i , and apply a union bound over $n \in [t-1], i \in [M]$. Plugging Eqn. 47, 48 and 49 back to Eqn. 45, the following inequality holds for any $t \in [T]$ with probability at least $1 - |\mathcal{E}|\delta_1 - MT\delta_2$:

$$\sum_{i=1}^M \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2 \quad (52)$$

$$\leq 2\sqrt{Mk + 2\log(1/\delta_1)} \sqrt{\sum_{i=1}^M \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2} \quad (53)$$

$$+ 2M\sqrt{2\log(2/\delta_2)kT^2\epsilon} + 2\sqrt{Mk\epsilon(2Mk + 2\log(1/\delta_1))} + 4M\lambda \quad (54)$$

By solving the above inequality, we know that

$$\sum_{i=1}^M \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2 \leq 32(Mk + \log(1/\delta_1)) + 4M\sqrt{2\log(2/\delta_2)kT^2\epsilon} \quad (55)$$

$$+ 4\sqrt{Mk\epsilon(2Mk + 2\log(1/\delta_1))} + 8M\lambda \quad (56)$$

Setting $\lambda = 1$, $\epsilon = \frac{1}{kM^2T^2}$, $\delta_1 = \frac{\delta}{2\left(\frac{6\sqrt{2k}}{\epsilon}\right)^{2kd}} \leq \frac{\delta}{2|\mathcal{E}|}$, and $\delta_2 = \frac{\delta}{2MT}$, the following inequality holds with probability $1 - \delta$:

$$\sum_{i=1}^M \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2 \leq L \stackrel{\text{def}}{=} 48(Mk + 5kd\log(kMT)) + 32\log(4MT) + 76\log(1/\delta) \quad (57)$$

At last we talk about the trivial setting where $k < d < 2k$. In this case, we can write $\hat{\boldsymbol{\Theta}}_t - \boldsymbol{\Theta} = \mathbf{R}_t$ where $\mathbf{R}_t \in \mathbb{R}^{d \times M}$. The proof then follows the same framework as the case when $d \geq 2k$, except that we don't need to consider \mathbf{U}_t and construct ϵ -net over all possible \mathbf{U}_t . It is not hard to show that $\sum_{i=1}^M \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2 \leq 24(Md + 2\log(Tk/\delta))$ in this case, which is also less than L since $d < 2k$. \square

A.2. PROOF OF THEOREM 1

With Lemma 1, we are ready to prove Theorem 1.

Proof. Let $\tilde{\mathbf{V}}_{t,i}(\lambda) = \mathbf{X}_{t,i}\mathbf{X}_{t,i}^\top + \lambda\mathbf{I}_d$ for some $\lambda > 0$.

$$\text{Reg}(T) = \sum_{t=1}^T \sum_{i=1}^M \langle \boldsymbol{\theta}_i, \mathbf{x}_{t,i}^* - \mathbf{x}_{t,i} \rangle \quad (58)$$

$$\leq \sum_{t=1}^T \sum_{i=1}^M \langle \tilde{\boldsymbol{\theta}}_{t,i} - \boldsymbol{\theta}_i, \mathbf{x}_{t,i} \rangle \quad (59)$$

$$= \sum_{t=1}^T \sum_{i=1}^M \langle \tilde{\boldsymbol{\theta}}_{t,i} - \hat{\boldsymbol{\theta}}_{t,i} + \hat{\boldsymbol{\theta}}_{t,i} - \boldsymbol{\theta}_i, \mathbf{x}_{t,i} \rangle \quad (60)$$

$$\leq \sum_{t=1}^T \sum_{i=1}^M \left(\|\tilde{\boldsymbol{\theta}}_{t,i} - \hat{\boldsymbol{\theta}}_{t,i}\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)} + \|\hat{\boldsymbol{\theta}}_{t,i} - \boldsymbol{\theta}_i\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)} \right) \|\mathbf{x}_{t,i}\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)^{-1}} \quad (61)$$

$$\leq \left(\sqrt{\sum_{t=1}^T \sum_{i=1}^M \|\tilde{\boldsymbol{\theta}}_{t,i} - \hat{\boldsymbol{\theta}}_{t,i}\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2} + \sqrt{\sum_{i=1}^M \|\hat{\boldsymbol{\theta}}_{t,i} - \boldsymbol{\theta}_i\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2} \right) \cdot \sqrt{\sum_{t=1}^T \sum_{i=1}^M \|\mathbf{x}_{t,i}\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)^{-1}}^2} \quad (62)$$

$$\leq 2\sqrt{T(L+4\lambda M)} \cdot \sqrt{\sum_{i=1}^M \sum_{t=1}^T \|\mathbf{x}_{t,i}\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)^{-1}}^2} \quad (63)$$

where the first inequality is due to $\sum_{i=1}^M \langle \boldsymbol{\theta}_i, \mathbf{x}_{t,i}^* \rangle \leq \langle \tilde{\boldsymbol{\theta}}_{t,i}, \mathbf{x}_{t,i} \rangle$ from the optimistic choice of $\tilde{\boldsymbol{\theta}}_{t,i}$ and $\mathbf{x}_{t,i}$. By Lemma 11 of Abbasi-Yadkori et al. (2011), as long as $\lambda \geq 1$ we have

$$\sum_{t=1}^T \|\mathbf{x}_{t,i}\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda')^{-1}}^2 \leq 2 \log \frac{\det(\tilde{\mathbf{V}}_{T,i}(\lambda'))}{\det(\lambda' \mathbf{I}_d)} \leq 2d \log \left(1 + \frac{T}{\lambda d} \right) \quad (64)$$

Therefore, we can finally bound the regret by choosing $\lambda = 1$

$$\text{Reg}(T) \leq 2\sqrt{T(L+4M)} \cdot \sqrt{\sum_{i=1}^M \sum_{t=1}^T \|\mathbf{x}_{t,i}\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda')^{-1}}^2} \quad (65)$$

$$\leq 2\sqrt{T(L+4M)} \cdot \sqrt{Md \log \left(1 + \frac{T}{d} \right)} \quad (66)$$

$$= \tilde{O} \left(M\sqrt{dkT} + d\sqrt{kMT} \right). \quad (67)$$

□

A.3. PROOF OF LEMMA 2

The proof of Lemma 2 follows the similar idea of Theorem 1 in Abbasi-Yadkori et al. (2011). We consider the σ -algebra $F_t = \sigma(\{\mathbf{x}_{1,i}\}_{i=1}^M, \{\mathbf{x}_{2,i}\}_{i=1}^M, \dots, \{\mathbf{x}_{t+1,i}\}_{i=1}^M, \{\eta_{1,i}\}_{i=1}^M, \{\eta_{2,i}\}_{i=1}^M, \dots, \{\eta_{t,i}\}_{i=1}^M)$, then $\{\mathbf{x}_{t,i}\}_{i=1}^M$ is F_{t-1} -measurable, and $\{\eta_{t,i}\}_{i=1}^M$ is F_t -measurable.

Define $\bar{\mathbf{x}}_{t,i} = \mathbf{U}^\top \mathbf{x}_{t,i}$ and $\mathbf{S}_{t,i} = \sum_{n=1}^t \bar{\mathbf{U}}^\top \mathbf{x}_{t,i} \eta_{t,i}$. Let

$$M_t(\mathbf{Q}) = \exp\left(\sum_{n=1}^t \sum_{i=1}^M \left[\eta_{t,i} \langle \mathbf{q}_i, \bar{\mathbf{x}}_{t,i} \rangle - \frac{1}{2} \langle \mathbf{q}_i, \bar{\mathbf{x}}_{t,i} \rangle^2\right]\right), \quad \mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_M] \in \mathbb{R}^{2k \times M} \quad (68)$$

Lemma 3. *Let τ be a stopping time w.r.t the filtration $\{F_t\}_{t=0}^\infty$. Then $M_t(\mathbf{Q})$ is almost surely well-defined and $\mathbb{E}[M_t(\mathbf{Q})] \leq 1$.*

Proof. Let $D_t(\mathbf{Q}) = \exp\left(\sum_{i=1}^M \left[\eta_{t,i} \langle \mathbf{q}_i, \bar{\mathbf{x}}_{t,i} \rangle - \frac{1}{2} \langle \mathbf{q}_i, \bar{\mathbf{x}}_{t,i} \rangle^2\right]\right)$. By the sub-Gaussianity of $\eta_{t,i}$, we have

$$\mathbb{E}\left[\exp\left(\left[\eta_{t,i} \langle \mathbf{q}_i, \bar{\mathbf{x}}_{t,i} \rangle - \frac{1}{2} \langle \mathbf{q}_i, \bar{\mathbf{x}}_{t,i} \rangle^2\right]\right) \mid F_{t-1}\right] \leq 1. \quad (69)$$

Then we have $\mathbb{E}[D_t(\mathbf{Q}) \mid F_{t-1}] \leq 1$. Further,

$$\mathbb{E}[M_t(\mathbf{Q}) \mid F_{t-1}] = \mathbb{E}[M_1(\mathbf{Q}) \cdots D_{t-1}(\mathbf{Q}) D_t(\mathbf{Q}) \mid F_{t-1}] \quad (70)$$

$$= D_1(\mathbf{Q}) \cdots D_{t-1}(\mathbf{Q}) \mathbb{E}[D_t(\mathbf{Q}) \mid F_{t-1}] \leq M_{t-1}(\mathbf{Q}) \quad (71)$$

This shows that $\{M_t(\mathbf{Q})\}_{t=0}^\infty$ is a supermartingale and $\mathbb{E}[M_t(\mathbf{Q})] \leq 1$.

Following the same argument of Lemma 8 in Abbasi-Yadkori et al. (2011), we show that $M_\tau(\mathbf{Q})$ is almost surely well-defined. By the convergence theorem for nonnegative supermartingales, $M_\infty(\mathbf{Q}) = \lim_{t \rightarrow \infty} M_t(\mathbf{Q})$ is almost surely well-defined. Therefore, $M_\tau(\mathbf{Q})$ is indeed well-defined independently of whether $\tau < \infty$ or not. Let $W_t(\mathbf{Q}) = M_{\min\{\tau, t\}}(\mathbf{Q})$ be a stopped version of $(M_t(\mathbf{Q}))_t$. By Fatou's Lemma, $\mathbb{E}[M_\tau(\mathbf{Q})] = \mathbb{E}[\liminf_{t \rightarrow \infty} W_t(\mathbf{Q})] \leq \liminf_{t \rightarrow \infty} \mathbb{E}[W_t(\mathbf{Q})] \leq 1$. This shows that $\mathbb{E}[M_\tau(\mathbf{Q})] \leq 1$. \square

The next lemma uses the ‘‘method of mixtures’’ technique to bound $\sum_{i=1}^M \|\mathbf{S}_{t,i}\|_{\bar{\mathbf{V}}_{t,i}^{-1}(\lambda)}^2$.

Lemma 4. *Let τ be a stopping time w.r.t the filtration $\{F_t\}_{t=0}^\infty$. Then, for $\delta > 0$, with probability $1 - \delta$,*

$$\sum_{i=1}^M \|\mathbf{S}_{\tau,i}\|_{\bar{\mathbf{V}}_{\tau,i}^{-1}(\lambda)}^2 \leq 2 \log\left(\frac{\prod_{i=1}^M (\det(\bar{\mathbf{V}}_{\tau,i})^{1/2} \det(\lambda \mathbf{I})^{-1/2})}{\delta}\right). \quad (72)$$

Proof. For each $i \in [M]$, let $\mathbf{\Lambda}_i$ be a \mathbb{R}^{2k} Gaussian random variable which is independent of all the other random variables and whose covariance is $\lambda^{-1} \mathbf{I}$. Define $M_t = \mathbb{E}[M_t(\mathbf{\Lambda}_1, \dots, \mathbf{\Lambda}_M) \mid F_\infty]$. We still have $\mathbb{E}[M_\tau] = \mathbb{E}[\mathbb{E}[M_\tau(\mathbf{\Lambda}_1, \dots, \mathbf{\Lambda}_M) \mid \{\mathbf{\Lambda}_i\}_{i=1}^M]] \leq 1$.

Now we calculate M_t . Define $M_{t,i}(\mathbf{q}_i) \stackrel{\text{def}}{=} \exp\left(\sum_{n=1}^t \left[\eta_{t,i} \langle \mathbf{q}_i, \bar{\mathbf{x}}_{t,i} \rangle - \frac{1}{2} \langle \mathbf{q}_i, \bar{\mathbf{x}}_{t,i} \rangle^2\right]\right)$, then we have $M_t = \mathbb{E}\left[\prod_{i=1}^M M_{t,i}(\mathbf{\Lambda}_i) \mid F_\infty\right] = \prod_{i=1}^M \mathbb{E}[M_{t,i}(\mathbf{\Lambda}_i) \mid F_\infty]$, where the second equality is due to the fact that $\{M_{t,i}(\mathbf{\Lambda}_i)\}_{i=1}^M$ are relatively independent given F_∞ . We only need to calculate $\mathbb{E}[M_{t,i}(\mathbf{\Lambda}_i) \mid F_\infty]$ for each $i \in [M]$.

Following the proof of Lemma 9 in Abbasi-Yadkori et al. (2011), we know that

$$\mathbb{E}[M_{t,i}(\mathbf{\Lambda}_i) \mid F_\infty] = \left(\frac{\det(\lambda \mathbf{I})}{\det(\bar{\mathbf{V}}_{t,i})}\right)^{1/2} \exp\left(\frac{1}{2} \|\mathbf{S}_{t,i}\|_{\bar{\mathbf{V}}_{t,i}^{-1}(\lambda)}^2\right). \quad (73)$$

Then we have

$$M_t = \prod_{i=1}^M \left(\left(\frac{\det(\lambda \mathbf{I})}{\det(\bar{\mathbf{V}}_{t,i})}\right)^{1/2}\right) \exp\left(\frac{1}{2} \sum_{i=1}^M \|\mathbf{S}_{t,i}\|_{\bar{\mathbf{V}}_{t,i}^{-1}(\lambda)}^2\right). \quad (74)$$

Since $\mathbb{E}[M_\tau] \leq 1$, we have

$$\begin{aligned}
 & \Pr \left[\sum_{i=1}^M \|\mathbf{S}_{\tau,i}\|_{\bar{\mathbf{V}}_{\tau,i}^{-1}(\lambda)}^2 > 2 \log \left(\frac{\prod_{i=1}^M (\det(\bar{\mathbf{V}}_{\tau,i})^{1/2} \det(\lambda \mathbf{I})^{-1/2})}{\delta} \right) \right] \\
 &= \Pr \left[\frac{\exp \left(\frac{1}{2} \sum_{i=1}^M \|\mathbf{S}_{\tau,i}\|_{\bar{\mathbf{V}}_{\tau,i}^{-1}(\lambda)}^2 \right)}{\delta^{-1} \left(\prod_{i=1}^M (\det(\bar{\mathbf{V}}_{\tau,i})^{1/2} \det(\lambda \mathbf{I})^{-1/2}) \right)} > 1 \right] \\
 &\leq \mathbb{E} \left[\frac{\exp \left(\sum_{i=1}^M \|\mathbf{S}_{\tau,i}\|_{\bar{\mathbf{V}}_{\tau,i}^{-1}(\lambda)}^2 \right)}{\delta^{-1} \left(\prod_{i=1}^M (\det(\bar{\mathbf{V}}_{\tau,i})^{1/2} \det(\lambda \mathbf{I})^{-1/2}) \right)} \right] \\
 &= \mathbb{E}[M_\tau] \delta \leq \delta.
 \end{aligned}$$

□

Proof. (Proof of Lemma 2) The only remaining issue is the stopping time construction. Define the bad event

$$B_t(\delta) \stackrel{\text{def}}{=} \left\{ \omega \in \Omega : \sum_{i=1}^M \|\mathbf{S}_{t,i}\|_{\bar{\mathbf{V}}_{t,i}^{-1}(\lambda)}^2 > 2 \log \left(\frac{\prod_{i=1}^M (\det(\bar{\mathbf{V}}_{t,i})^{1/2} \det(\lambda \mathbf{I})^{-1/2})}{\delta} \right) \right\} \quad (75)$$

Consider the stopping time $\tau(\omega) = \min\{t \geq 0 : \omega \in B_t(\delta)\}$, we have $\bigcup_{t \geq 0} B_t(\delta) = \{\omega : \tau(\omega) < \infty\}$.

By lemma 4, we have

$$\Pr \left[\bigcup_{t \geq 0} B_t(\delta) \right] = \Pr[\tau < \infty] \quad (76)$$

$$= \Pr \left[\sum_{i=1}^M \|\mathbf{S}_{\tau,i}\|_{\bar{\mathbf{V}}_{\tau,i}^{-1}(\lambda)}^2 > 2 \log \left(\frac{\prod_{i=1}^M (\det(\bar{\mathbf{V}}_{\tau,i})^{1/2} \det(\lambda \mathbf{I})^{-1/2})}{\delta} \right), \tau \leq \infty \right] \quad (77)$$

$$\leq \Pr \left[\sum_{i=1}^M \|\mathbf{S}_{\tau,i}\|_{\bar{\mathbf{V}}_{\tau,i}^{-1}(\lambda)}^2 > 2 \log \left(\frac{\prod_{i=1}^M (\det(\bar{\mathbf{V}}_{\tau,i})^{1/2} \det(\lambda \mathbf{I})^{-1/2})}{\delta} \right) \right] \quad (78)$$

$$\leq \delta. \quad (79)$$

□

A.4. PROOF OF THEOREM 2

Proof. The proof follows the same idea of that for Theorem 1. The only difference is that, in our setting, we have $y_{t,i} = \mathbf{x}_{t,i}^\top \mathbf{B} \mathbf{w}_i + \eta_{t,i} + \Delta_{t,i}$, where $\theta_i = \mathbf{B} \mathbf{w}_i$ is the best approximator for task $i \in [M]$ such that $\left| \mathbb{E}[y_i | \mathbf{x}_i] - \langle \mathbf{x}_i, \hat{\mathbf{B}} \hat{\mathbf{w}}_i \rangle \right| \leq \zeta$, and $\|\Delta_{t,i}\| \leq \zeta$. Define $\Delta_{t,i} = [\Delta_{1,i}, \Delta_{2,i}, \dots, \Delta_{t,i}]$. Similarly, by the optimality of $\hat{\mathbf{B}}_t$ and $\hat{\mathbf{W}}_t = [\hat{\mathbf{w}}_{t,1}, \dots, \hat{\mathbf{w}}_{t,M}]$, we know that $\sum_{i=1}^M \left\| \mathbf{y}_{t-1,i} - \mathbf{X}_{t-1,i}^\top \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} \right\|_2^2 \leq \sum_{i=1}^M \left\| \mathbf{y}_{t-1,i} - \mathbf{X}_{t-1,i}^\top \mathbf{B} \mathbf{w}_i \right\|_2^2$. Since $\mathbf{y}_{t-1,i} = \mathbf{X}_{t-1,i}^\top \mathbf{B} \mathbf{w}_i + \eta_{t-1,i} + \Delta_{t,i}$, thus we have

$$\sum_{i=1}^M \left\| \mathbf{X}_{t-1,i}^\top \left(\hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right) \right\|^2 \quad (80)$$

$$\leq 2 \sum_{i=1}^M \eta_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \left(\hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right) + 2 \sum_{i=1}^M \Delta_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \left(\hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right) \quad (81)$$

$$\leq 2 \sum_{i=1}^M \eta_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \left(\hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right) + 2 \sum_{i=1}^M \left\| \mathbf{X}_{t-1,i} \Delta_{t-1,i} \right\|_{\tilde{\mathbf{V}}_{t-1,i}^{-1}(\lambda)} \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)} \quad (82)$$

$$\leq 2 \sum_{i=1}^M \eta_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \left(\hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right) + 2 \sum_{i=1}^M \sqrt{T} \zeta \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)} \quad (83)$$

$$\leq 2 \sum_{i=1}^M \eta_{t-1,i}^\top \mathbf{X}_{t-1,i}^\top \left(\hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right) + 2\sqrt{MT} \zeta \sqrt{\sum_{i=1}^M \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2} \quad (84)$$

The third inequality follows from Projection Bound (Lemma 8) in Zanette et al. (2020a). The first term of Eqn 84 shares the same form of Eqn 28. Following the same proof idea of Lemma 1, we know that with probability $1 - \delta$,

$$\sum_{i=1}^M \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2 \quad (85)$$

$$\leq \left(2\sqrt{Mk + 8kd \log(kMT/\delta)} + 2\sqrt{MT} \zeta \right) \sqrt{\sum_{i=1}^M \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2} + 4M + 4\sqrt{\log(4MT/\delta)} \quad (86)$$

Solving for $\sum_{i=1}^M \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2$, we know that the true parameter $\mathbf{B} \mathbf{W}$ is always contained in the confidence set, i.e.

$$\sum_{i=1}^M \left\| \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,i} - \mathbf{B} \mathbf{w}_i \right\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2 \leq L', \quad (87)$$

where $L' = 2L + 32MT\zeta^2$.

Thus we have

$$\text{Reg}(T) = \sum_{t=1}^T \sum_{i=1}^M (y_{t,i}^* - y_{t,i}) \quad (88)$$

$$\leq 2MT\zeta + \sum_{t=1}^T \sum_{i=1}^M \langle \boldsymbol{\theta}_i, \mathbf{x}_{t,i}^* - \mathbf{x}_{t,i} \rangle \quad (89)$$

$$\leq 2MT\zeta + \sum_{t=1}^T \sum_{i=1}^M \langle \tilde{\boldsymbol{\theta}}_{t,i} - \boldsymbol{\theta}_i, \mathbf{x}_{t,i} \rangle \quad (90)$$

$$= 2MT\zeta + \sum_{t=1}^T \sum_{i=1}^M \langle \tilde{\boldsymbol{\theta}}_{t,i} - \hat{\boldsymbol{\theta}}_{t,i} + \hat{\boldsymbol{\theta}}_{t,i} - \boldsymbol{\theta}_i, \mathbf{x}_{t,i} \rangle \quad (91)$$

$$\leq 2MT\zeta + \sum_{t=1}^T \sum_{i=1}^M \left(\|\tilde{\boldsymbol{\theta}}_{t,i} - \hat{\boldsymbol{\theta}}_{t,i}\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)} + \|\hat{\boldsymbol{\theta}}_{t,i} - \boldsymbol{\theta}_i\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)} \right) \|\mathbf{x}_{t,i}\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)^{-1}} \quad (92)$$

$$\leq 2MT\zeta + \left(\sqrt{\sum_{t=1}^T \sum_{i=1}^M \|\tilde{\boldsymbol{\theta}}_{t,i} - \hat{\boldsymbol{\theta}}_{t,i}\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2} + \sqrt{\sum_{i=1}^M \|\hat{\boldsymbol{\theta}}_{t,i} - \boldsymbol{\theta}_i\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)}^2} \right) \cdot \sqrt{\sum_{t=1}^T \sum_{i=1}^M \|\mathbf{x}_{t,i}\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)^{-1}}^2} \quad (93)$$

$$\leq 2MT\zeta + 2\sqrt{T(L' + 4\lambda M)} \cdot \sqrt{\sum_{i=1}^M \sum_{t=1}^T \|\mathbf{x}_{t,i}\|_{\tilde{\mathbf{V}}_{t-1,i}(\lambda)^{-1}}^2} \quad (94)$$

$$\leq 2MT\zeta + 2\sqrt{T(L' + 4\lambda M)} \sqrt{Md \log\left(1 + \frac{T}{d}\right)} \quad (95)$$

$$= \tilde{O}(M\sqrt{dkT} + d\sqrt{kMT} + MT\sqrt{d\zeta}), \quad (96)$$

where the second inequality is due to $\sum_{i=1}^M \langle \boldsymbol{\theta}_i, \mathbf{x}_{t,i}^* \rangle \leq \langle \tilde{\boldsymbol{\theta}}_{t,i}, \mathbf{x}_{t,i} \rangle$ from the optimistic choice of $\tilde{\boldsymbol{\theta}}_{t,i}$ and $\mathbf{x}_{t,i}$. The third inequality is due to Eqn 87. The last inequality is from Eqn 64. \square

A.5. PROOF OF THEOREM 3

Since our setting is strictly harder than the setting of multi-task linear bandit with infinite arms in Yang et al. (2020), we can prove the following lemma directly from their Theorem 4 by reduction.

Lemma 5. *Under the setting of Theorem 3, the regret of any Algorithm \mathcal{A} is lower bounded by $\Omega\left(Mk\sqrt{T} + d\sqrt{kMT}\right)$.*

In order to prove Theorem 3, we only need to show that the following lemma is true.

Lemma 6. *Under the setting of Theorem 3, the regret of any Algorithm \mathcal{A} is lower bounded by $\Omega\left(MT\sqrt{d}\zeta\right)$.*

Proof. (Proof of Lemma 6)

To prove Lemma 6, we leverage the lower bound for misspecified linear bandits in the single-task setting. We restate the following lemma from the previous literature with a slight modification of notations.

Lemma 7. (Proposition 6 in Zanette et al. (2020a)). *There exists a feature map $\phi : \mathcal{A} \rightarrow \mathbb{R}^d$ that defines a misspecified linear bandits class \mathcal{M} such that every bandit instance in that class has reward response:*

$$\mu_a = \phi_a^\top \theta + z_a$$

for any action a (Here $z_a \in [0, \zeta]$ is the deviation from linearity and $\mu_a \in [0, 1]$) and such that the expected regret of any algorithm on at least a member of the class up to round T is $\Omega(\sqrt{d}\zeta T)$.

Suppose M can be exactly divided by k , we construct the following instances to prove lemma 6. We divide M tasks into k groups. Each group shares the same parameter θ_i . To be more specific, we let $\mathbf{w}_1 = \mathbf{w}_2 = \dots = \mathbf{w}_{M/k} = \mathbf{e}_1$, $\mathbf{w}_{M/k+1} = \mathbf{w}_{M/k+2} = \dots = \mathbf{w}_{2M/k} = \mathbf{e}_2$, \dots , $\mathbf{w}_{(k-1)M/k+1} = \mathbf{w}_{(k-1)M/k+2} = \dots = \mathbf{w}_M = \mathbf{e}_k$. Under this construction, the parameters θ_i for these tasks are exactly the same in each group, but relatively independent among different groups. That is to say, the expected regret lower bound is at least the summation of the regret lower bounds in all k groups.

Now we consider the regret lower bound for group $j \in [k]$. Since the parameters are shared in the same group, the regret of running an algorithm for M/k tasks with T steps each is at least the regret of running an algorithm for single-task linear bandit with $M/k \cdot T$ steps. By Lemma 7, the regret for single-task linear bandit with MT/k steps is at least $\Omega(\sqrt{d}\zeta MT/k)$. Summing over all k groups, we can prove that the regret lower bound is $\Omega(\sqrt{d}\zeta MT)$. \square

Combining Lemma 5 and Lemma 6, we complete the proof of Theorem 3.

B. Proof of Theorem 4

B.1. DEFINITIONS AND FIRST STEP ANALYSIS

Before presenting the proof of theorem 4, we will make a first step analysis on the low-rank least-square estimator in equation 15.

For any $\{Q_{h+1}^i\}_{i=1}^M \in \mathcal{Q}_{h+1}$, there exists $\{\hat{\theta}_h^i(Q_{h+1}^i)\}_{i=1}^M \in \Theta_h$ that

$$\Delta_h^i(Q_{h+1}^i)(s, a) = \mathcal{T}_h^i(Q_{h+1}^i)(s, a) - \phi(s, a)^\top \hat{\theta}_h^i(Q_{h+1}^i) \quad (97)$$

where the approximation error $\|\Delta_h^i(Q_{h+1}^i)\|_\infty \leq \mathcal{I}$ is small for each $i \in [M]$. We also use $\hat{B}_h \hat{w}_h^i(Q_{h+1}^i)$ in place of $\hat{\theta}_h^i(Q_{h+1}^i)$ in the following sections since we can write $\hat{\theta}_h^i$ as $\hat{B}_h \hat{w}_h^i$.

In the multi-task low-rank least-square regression (equation 15), we are actually trying to recover $\hat{\theta}_h^i$. However, due to the noise and representation error (i.e. the inherent Bellman error), we can only obtain an approximate solution $\hat{\theta}_h^i = \hat{B}_h \hat{w}_h^i$ (see the global optimization problem in Definition 1).

$$\left(\hat{\theta}_h^1, \dots, \hat{\theta}_h^M\right) = \hat{B}_h \left[\hat{w}_h^1 \quad \hat{w}_h^2 \quad \dots \quad \hat{w}_h^M\right] \quad (98)$$

$$= \underset{\|\mathbf{B}_h \mathbf{w}_h^i\|_2 \leq D}{\operatorname{argmin}} \sum_{i=1}^M \sum_{j=1}^{t-1} \left(\phi(s_{hj}^i, a_{hj}^i)^\top \mathbf{B}_h \mathbf{w}_h^i - R(s_{hj}^i, a_{hj}^i) - \max_a Q_{h+1}^i(s_{h+1,j}^i) \right)^2 \quad (99)$$

$$= \underset{\|\mathbf{B}_h \mathbf{w}_h^i\|_2 \leq D}{\operatorname{argmin}} \sum_{i=1}^M \sum_{j=1}^{t-1} \left(\phi(s_{hj}^i, a_{hj}^i)^\top \mathbf{B}_h \mathbf{w}_h^i - \mathcal{T}_h^i(Q_{h+1}^i)(s_{hj}^i, a_{hj}^i) - z_{hj}^i(Q_{h+1}^i)(s_{hj}^i, a_{hj}^i) \right)^2 \quad (100)$$

where $z_{hj}^i(Q_{h+1}^i)(s_{hj}^i, a_{hj}^i) \stackrel{\text{def}}{=} R(s_{hj}^i, a_{hj}^i) + \max_a Q_{h+1}^i(s_{h+1,j}^i, a) - \mathcal{T}_h^i(Q_{h+1}^i)(s_{hj}^i, a_{hj}^i)$.

Define $\Phi_{ht}^i \in \mathbb{R}^{(t-1) \times d}$ to be the collection of linear features up to episode $t-1$ in task i , i.e. the j -th row of Φ_{ht}^i is $\phi(s_{hj}^i, a_{hj}^i)^\top$. Let $\mathbf{Y}_{ht}^i \in \mathbb{R}^{t-1}$ be a vector whose j -th dimension is $\mathcal{T}_h^i(Q_{h+1}^i)(s_{hj}^i, a_{hj}^i) + z_{hj}^i(Q_{h+1}^i)(s_{hj}^i, a_{hj}^i)$. Then the objective in (100) can be written as

$$\underset{\|\mathbf{B}_h \mathbf{w}_h^i\|_2 \leq D}{\operatorname{argmin}} \sum_{i=1}^M \|\Phi_{ht}^i \mathbf{B}_h \mathbf{w}_h^i - \mathbf{Y}_{ht}^i\|_2^2 \quad (101)$$

Therefore, we have

$$\sum_{i=1}^M \left\| \Phi_{ht}^i \hat{B}_h \hat{w}_h^i(Q_{h+1}^i) - \mathbf{Y}_{ht}^i \right\|_2^2 \leq \sum_{i=1}^M \left\| \Phi_{ht}^i \hat{B}_h \hat{w}_h^i(Q_{h+1}^i) - \mathbf{Y}_{ht}^i \right\|_2^2 \quad (102)$$

which implies

$$\sum_{i=1}^M \left\| \Phi_{ht}^i \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \Phi_{ht}^i \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_2^2 \quad (103)$$

$$\leq 2 \sum_{i=1}^M (\Delta_{ht}^i)^\top \Phi_{ht}^i \left(\hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right) \quad (104)$$

$$+ 2 \sum_{i=1}^M (\mathbf{z}_{ht}^i)^\top \Phi_{ht}^i \left(\hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right) \quad (105)$$

where $\Delta_{ht}^i \stackrel{\text{def}}{=} \left[\Delta_{h1}^i(Q_{h+1}^i)(s_{h1}^i, a_{h1}^i) \quad \Delta_{h2}^i(Q_{h+1}^i)(s_{h2}^i, a_{h2}^i) \quad \cdots \quad \Delta_{h,t-1}^i(Q_{h+1}^i)(s_{h,t-1}^i, a_{h,t-1}^i) \right] \in \mathbb{R}^{t-1}$, and $\mathbf{z}_{ht}^i \stackrel{\text{def}}{=} \left[z_{h1}^i(Q_{h+1}^i)(s_{h1}^i, a_{h1}^i) \quad \cdots \quad z_{h,t-1}^i(Q_{h+1}^i)(s_{h,t-1}^i, a_{h,t-1}^i) \right] \in \mathbb{R}^{t-1}$.

In the next sections we will show how to bound 104 and 105.

B.2. FAILURE EVENT

Define the failure event at step h in episode t as

Definition 2 (Failure Event).

$$E_{ht} \stackrel{\text{def}}{=} I \left[\exists \{Q_{h+1}^i\}_{i=1}^M \in \mathcal{Q}_{h+1} \quad \sum_{i=1}^M (z_{ht}^i)^\top \Phi_{ht}^i \left(\hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right) > \right] \quad (106)$$

$$F_h^1 \sqrt{\sum_{i=1}^M \left\| \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_{\tilde{\mathbf{V}}_{ht}^i(\lambda)}^2} + F_h^2 \quad (107)$$

where F_h^1 and F_h^2 will be specified later.

We have the following lemma to bound the probability of E_{ht} .

Lemma 8. *For the input parameter $\delta > 0$, there exists F_h^1 and F_h^2 such that*

$$\mathbb{P} \left(\bigcup_{t=1}^T \bigcup_{h=1}^H E_{ht} \right) \leq \frac{\delta}{2} \quad (108)$$

Proof. According to Lemma A.5 of Du et al. (2020), there exists an ϵ -net \mathcal{E}_{h+1}^o over $\mathcal{O}^{d \times k}$ (with regards to the Frobenius norm) such that $|\mathcal{E}_{h+1}^o| \leq (6\sqrt{k}/\epsilon')^{kd}$. Moreover, there exists an ϵ -net \mathcal{E}_{h+1}^b over \mathcal{B}^k that $|\mathcal{E}_{h+1}^b| \leq (1 + 2/\epsilon')^k$. We can show a corresponding ϵ -net $\mathcal{E}_{h+1}^{\text{mul}} \stackrel{\text{def}}{=} \mathcal{E}_{h+1}^o \times (\mathcal{E}_{h+1}^b)^M$ over Θ_{h+1} .

For any $(Q_{h+1}^1(\mathbf{B}_{h+1}\mathbf{w}_{h+1}^1), \dots, Q_{h+1}^M(\mathbf{B}_{h+1}\mathbf{w}_{h+1}^M)) \in \mathcal{Q}_{h+1}$, there exists $\bar{\mathbf{B}}_{h+1} \in \mathcal{E}_{h+1}^o$ and $(\bar{\mathbf{w}}_{h+1}^1, \dots, \bar{\mathbf{w}}_{h+1}^M) \in (\mathcal{E}_{h+1}^b)^M$ such that

$$\|\mathbf{B}_{h+1} - \bar{\mathbf{B}}_{h+1}\|_F \leq \epsilon' \quad \|\mathbf{w}_{h+1}^i - \bar{\mathbf{w}}_{h+1}^i\|_2 \leq \epsilon', \forall i \in [M]$$

Therefore,

$$\|\mathbf{B}_{h+1}\mathbf{w}_{h+1}^i - \bar{\mathbf{B}}_{h+1}\bar{\mathbf{w}}_{h+1}^i\|_2 \leq 2\epsilon', \forall i \in [M]$$

Define \bar{Q}_{h+1}^i to be $Q_{h+1}^i(\bar{\mathbf{B}}_{h+1}\bar{\mathbf{w}}_{h+1}^i)$, and let $\bar{z}_{ht}^i \stackrel{\text{def}}{=} [z_{h1}^i(\bar{Q}_{h+1}^i)(s_{h1}^i, a_{h1}^i) \quad \dots \quad z_{h,t-1}^i(\bar{Q}_{h+1}^i)(s_{h,t-1}^i, a_{h,t-1}^i)] \in \mathbb{R}^{t-1}$, then

$$\sum_{i=1}^M (z_{ht}^i)^\top \Phi_{ht}^i \left(\hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right) \quad (109)$$

$$= \sum_{i=1}^M (\bar{z}_{ht}^i)^\top \Phi_{ht}^i \left(\hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right) \quad (110)$$

$$+ \sum_{i=1}^M (z_{ht}^i - \bar{z}_{ht}^i)^\top \Phi_{ht}^i \left(\hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right) \quad (111)$$

For fixed $\{\bar{\mathbf{B}}_{h+1}\bar{\mathbf{w}}_{h+1}^i\}_{i=1}^M \in \mathcal{E}_{h+1}^{\text{mul}}$, $z_{h,j}^i(\bar{Q}_{h+1}^i)(s_{h,j}^i, a_{h,j}^i)$ is zero-mean 1-subgaussian conditioned on $\mathcal{F}_{h,j}$ according to Assumption 3. Thus, we can use exactly the same argument as in Lemma 1 to show that

$$\sum_{i=1}^M (\bar{\mathbf{z}}_{ht}^i)^\top \Phi_{ht}^i \left(\hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right) \quad (112)$$

$$\leq \sqrt{Mk + 5kd \log(kMT) + 2 \log(1/\delta')} \sqrt{\sum_{i=1}^M \left\| \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_{\hat{\mathbf{V}}_{ht}^i(\lambda)}^2} \quad (113)$$

$$+ \sqrt{2 \log(2MT/\delta')} + \sqrt{k + 3kd \log(kMT) + \log(1/\delta')} \quad (114)$$

by setting $\epsilon = \frac{1}{kM^2T^2}$, $\delta_1 = \frac{\delta'}{2 \left(\frac{6\sqrt{2k}}{\epsilon} \right)^{2kd}}$, and $\delta_2 = \frac{\delta'}{2MT}$ in equation 54. Thus, we have that with probability $1 - \delta'$ the inequality above holds for any $h \in [H], t \in [T]$. Take $\delta = \frac{\delta'}{2 \lceil \mathcal{E}_{h+1}^{\text{mul}} \rceil}$, by union bound we know the above inequality holds with probability $1 - \delta$ for any $\{\bar{\mathbf{B}}_{h+1} \bar{\mathbf{w}}_{h+1}^i\}_{i=1}^M \in \mathcal{E}_{h+1}^{\text{mul}}$ and any $h \in [H], t \in [T]$.

Since it holds that $|Q_{h+1}^i(\mathbf{B}_{h+1} \mathbf{w}_{h+1}^i)(s, a) - Q_{h+1}^i(\bar{\mathbf{B}}_{h+1} \bar{\mathbf{w}}_{h+1}^i)(s, a)| \leq 2\epsilon'$ for any $(s, a) \in \mathcal{S} \times \mathcal{A}, i \in [M]$, we have

$$|z_{hj}^i(\bar{Q}_{h+1}^i)(s_{hj}^i, a_{hj}^i) - z_{hj}^i(Q_{h+1}^i)(s_{hj}^i, a_{hj}^i)| \leq 8\epsilon' \quad (115)$$

Then we have

$$\sum_{i=1}^M (\mathbf{z}_{ht}^i - \bar{\mathbf{z}}_{ht}^i)^\top \Phi_{ht}^i \left(\hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right) \quad (116)$$

$$\leq \sum_{i=1}^M \left\| (\Phi_{ht}^i)^\top (\mathbf{z}_{ht}^i - \bar{\mathbf{z}}_{ht}^i) \right\|_{\hat{\mathbf{V}}_{ht}^i(\lambda)^{-1}} \left\| \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_{\hat{\mathbf{V}}_{ht}^i(\lambda)} \quad (117)$$

$$\leq 8\epsilon' \sqrt{T} \sum_{i=1}^M \left\| \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_{\hat{\mathbf{V}}_{ht}^i(\lambda)} \quad (118)$$

$$\leq 8\epsilon' \sqrt{MT} \sqrt{\sum_{i=1}^M \left\| \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_{\hat{\mathbf{V}}_{ht}^i(\lambda)}^2} \quad (119)$$

for arbitrary $\{Q_{h+1}^i\}$ and any $h \in [H], t \in [T]$. The second inequality follows from the Projection Bound (Lemma 8) in Zanette et al. (2020a).

Take $\epsilon' = 1/8\sqrt{MT}$, we finally finish the proof by setting

$$F_h^1 \stackrel{\text{def}}{=} \sqrt{9kd \log(kMT) + 5Mk \log(MT) + 2 \log(2/\delta)} \quad (120)$$

$$F_h^2 \stackrel{\text{def}}{=} \sqrt{4kd \log(kMT) + 5Mk \log(MT) + 2 \log(2/\delta)} \quad (121)$$

$$+ \sqrt{k + 5kd \log(kMT) + 2Mk \log(MT) + \log(2/\delta)} \quad (122)$$

□

In the next sections we assume the failure event $\bigcup_{t=1}^T \bigcup_{h=1}^H E_{ht}$ won't happen.

B.3. BELLMAN ERROR

Outside the failure event, we can bound the estimation error of the least-square regression 15.

Lemma 9. For any episode $t \in [T]$ and step $h \in [H]$, any $\{Q_{h+1}^i\}_{i=1}^M \in \mathcal{Q}_{h+1}$, we have

$$\sum_{i=1}^M \left\| \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_{\tilde{\mathbf{V}}_{ht}^i(\lambda)}^2 \leq \alpha_{ht} \stackrel{\text{def}}{=} \left(2\sqrt{MT\mathcal{I}} + 2F_h^1 + \sqrt{2F_h^2 + 4MD^2\lambda} \right)^2 \quad (123)$$

Proof. Recall that

$$\sum_{i=1}^M \left\| \Phi_{ht}^i \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \Phi_{ht}^i \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_2^2 \quad (124)$$

$$\leq 2 \sum_{i=1}^M (\Delta_{ht}^i)^\top \Phi_{ht}^i \left(\hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right) \quad (125)$$

$$+ 2 \sum_{i=1}^M (\mathbf{z}_{ht}^i)^\top \Phi_{ht}^i \left(\hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right) \quad (126)$$

For the first term, we have

$$\sum_{i=1}^M (\Delta_{ht}^i)^\top \Phi_{ht}^i \left(\hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right) \quad (127)$$

$$\leq \sum_{i=1}^M \left\| (\Phi_{ht}^i)^\top \Delta_{ht}^i \right\|_{\tilde{\mathbf{V}}_{ht}^i(\lambda)^{-1}} \left\| \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_{\tilde{\mathbf{V}}_{ht}^i(\lambda)} \quad (128)$$

$$\leq \sqrt{T\mathcal{I}} \sum_{i=1}^M \left\| \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_{\tilde{\mathbf{V}}_{ht}^i(\lambda)} \quad (129)$$

$$\leq \sqrt{MT\mathcal{I}} \sqrt{\sum_{i=1}^M \left\| \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_{\tilde{\mathbf{V}}_{ht}^i(\lambda)}^2} \quad (130)$$

The second inequality follows from the Projection Bound (Lemma 8) in Zanette et al. (2020a), and the last inequality is due to Cauchy-Schwarz.

Outside the failure event, we have

$$\sum_{i=1}^M \left\| \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_{\tilde{\mathbf{V}}_{ht}^i(\lambda)}^2 \quad (131)$$

$$\leq \sum_{i=1}^M \left\| \Phi_{ht}^i \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \Phi_{ht}^i \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_2^2 + 4MD^2\lambda \quad (132)$$

$$\leq \left(2\sqrt{MT\mathcal{I}} + 2F_h^1 \right) \sqrt{\sum_{i=1}^M \left\| \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_{\tilde{\mathbf{V}}_{ht}^i(\lambda)}^2} + 2F_h^2 + 4MD^2\lambda \quad (133)$$

which implies

$$\sum_{i=1}^M \left\| \hat{\mathbf{B}}_h \hat{\mathbf{w}}_h^i(Q_{h+1}^i) - \dot{\mathbf{B}}_h \dot{\mathbf{w}}_h^i(Q_{h+1}^i) \right\|_{\tilde{\mathbf{V}}_{ht}^i(\lambda)}^2 \quad (134)$$

$$\leq \left(2\sqrt{MT\mathcal{I}} + 2F_h^1 \right)^2 + 2F_h^2 + 4MD^2\lambda + \left(2\sqrt{MT\mathcal{I}} + 2F_h^1 \right) \sqrt{2F_h^2 + 4MD^2\lambda} \quad (135)$$

$$\leq \left(2\sqrt{MT\mathcal{I}} + 2F_h^1 + \sqrt{2F_h^2 + 4MD^2\lambda} \right)^2 \quad (136)$$

□

Lemma 10 (Bound on Bellman Error). *Outside the failure event, for any feasible solution $\{Q_h^i(\bar{\theta}_h^i)\}_h^i$ (\bar{Q}_h^i for short, with a little abuse of notations) of the global optimization procedure in definition 1, for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, any $h \in [H]$, $t \in [T]$*

$$\sum_{i=1}^M |\bar{Q}_h^i(s, a) - \mathcal{T}_h^i \bar{Q}_{h+1}^i(s, a)| \leq M\mathcal{I} + 2\sqrt{\alpha_{ht} \cdot \sum_{i=1}^M \|\phi(s, a)\|_{\tilde{\mathbf{V}}_{ht}^i(\lambda)-1}^2} \quad (137)$$

Proof.

$$\sum_{i=1}^M |\bar{Q}_h^i(s, a) - \mathcal{T}_h^i \bar{Q}_{h+1}^i(s, a)| = \sum_{i=1}^M \left| \phi(s, a)^\top \bar{\theta}_h^i - \phi(s, a)^\top \dot{\theta}_h^i(\bar{Q}_{h+1}^i) - \Delta_h^i(\bar{Q}_{h+1}^i)(s, a) \right| \quad (138)$$

$$\leq M\mathcal{I} + \sum_{i=1}^M \left| \phi(s, a)^\top \bar{\theta}_h^i - \phi(s, a)^\top \dot{\theta}_h^i(\bar{Q}_{h+1}^i) \right| \quad (139)$$

$$\leq M\mathcal{I} + \sum_{i=1}^M \left(\left| \phi(s, a)^\top \dot{\theta}_h^i(\bar{Q}_{h+1}^i) - \phi(s, a)^\top \hat{\theta}_h^i \right| + \left| \phi(s, a)^\top \hat{\theta}_h^i - \phi(s, a)^\top \bar{\theta}_h^i \right| \right) \quad (140)$$

$$\leq M\mathcal{I} + \sum_{i=1}^M \|\phi(s, a)\|_{\tilde{\mathbf{V}}_{ht}^i(\lambda)-1} \left(\left\| \dot{\theta}_h^i(\bar{Q}_{h+1}^i) - \hat{\theta}_h^i \right\|_{\tilde{\mathbf{V}}_{ht}^i(\lambda)} + \left\| \hat{\theta}_h^i - \bar{\theta}_h^i \right\|_{\tilde{\mathbf{V}}_{ht}^i(\lambda)} \right) \quad (141)$$

$$\leq M\mathcal{I} + 2\sqrt{\alpha_{ht} \cdot \sum_{i=1}^M \|\phi(s, a)\|_{\tilde{\mathbf{V}}_{ht}^i(\lambda)-1}^2} \quad (142)$$

The first equality is due to the definition of $\Delta_h^i(\bar{Q}_{h+1}^i)(s, a)$. The last inequality is due to lemma 9. □

B.4. OPTIMISM

We can find the "best" approximator of optimal value functions in our function class recursively defined as

$$(\boldsymbol{\theta}_h^{1*}, \boldsymbol{\theta}_h^{2*}, \dots, \boldsymbol{\theta}_h^{M*}) \stackrel{\text{def}}{=} \underset{(\boldsymbol{\theta}_h^1, \boldsymbol{\theta}_h^2, \dots, \boldsymbol{\theta}_h^M) \in \Theta_h}{\operatorname{argmin}} \sup_{s, a, i} |(\boldsymbol{\phi}(s, a)^\top \boldsymbol{\theta}_h^i - \mathcal{T}_h^i Q_{h+1}^i(\boldsymbol{\theta}_{h+1}^{i*})) (s, a)| \quad (143)$$

with $\boldsymbol{\theta}_{H+1}^{i*} = \mathbf{0}, \forall i \in [M]$

For the accuracy of this best approximator, we have

Lemma 11. For any $h \in [H]$,

$$\sup_{(s, a) \in \mathcal{S} \times \mathcal{A}, i \in [M]} |Q_h^{i*}(s, a) - \boldsymbol{\phi}(s, a)^\top \boldsymbol{\theta}_h^i| \leq (H - h + 1)\mathcal{I}$$

where Q_h^{i*} is the optimal value function for task i . This lemma is derived directly from Lemma 6 in Zanette et al. (2020a).

For our solution of the problem in Definition 1 in episode t , we have the following lemma:

Lemma 12. $\{(\boldsymbol{\theta}_h^{1*}, \boldsymbol{\theta}_h^{2*}, \dots, \boldsymbol{\theta}_h^{M*})\}_{h=1}^H$ is a feasible solution of the problem in Definition 1. Moreover, denote the solution of the problem in Definition 1 in episode t by $\bar{\boldsymbol{\theta}}_{ht}^i$ for $h \in [H], i \in [M]$, it holds that

$$\sum_{i=1}^M V_1^i(\bar{\boldsymbol{\theta}}_{1t}^i)(s_{1t}^i) \geq \sum_{i=1}^M V_1^{i*}(s_{1t}^i) - MH\mathcal{I} \quad (144)$$

Proof. First we show that $\{(\boldsymbol{\theta}_h^{1*}, \boldsymbol{\theta}_h^{2*}, \dots, \boldsymbol{\theta}_h^{M*})\}_{h=1}^H$ is a feasible solution. We can construct $\{\bar{\boldsymbol{\xi}}_h^i\}_{i=1}^M$ so that $\bar{\boldsymbol{\theta}}_h^i = \boldsymbol{\theta}_h^{i*}$ and no other constraints are violated. We use an inductive construction, and the base case when $\bar{\boldsymbol{\theta}}_{H+1}^i = \boldsymbol{\theta}_{H+1}^{i*} = \mathbf{0}$ is trivial.

Now suppose we have $\{\bar{\boldsymbol{\xi}}_y^i\}_{i=1}^M$ for $y = h + 1, \dots, H$ such that $\bar{\boldsymbol{\theta}}_y^i = \boldsymbol{\theta}_y^{i*}$ for $y = h + 1, \dots, H$ and $i \in [M]$, we show we can find $\{\bar{\boldsymbol{\xi}}_h^i\}_{i=1}^M$ so $\bar{\boldsymbol{\theta}}_h^i = \boldsymbol{\theta}_h^{i*}$ for $i \in [M]$, and no constraints are violated. From the definition of $\boldsymbol{\theta}_h^{i*}$ we can set (with a little abuse of notations)

$$\hat{\boldsymbol{\theta}}_h^i(\boldsymbol{\theta}_{h+1}^{i*}) = \boldsymbol{\theta}_h^{i*} \quad (145)$$

According to lemma 9 we have

$$\sum_{i=1}^M \left\| \hat{\boldsymbol{\theta}}_h^i(\boldsymbol{\theta}_{h+1}^{i*}) - \hat{\boldsymbol{\theta}}_h^i(\boldsymbol{\theta}_{h+1}^{i*}) \right\|_{\bar{\mathbf{V}}_{ht}^i(\lambda)}^2 \leq \alpha_{ht} \quad (146)$$

Therefore, set $\bar{\boldsymbol{\xi}}_h^i = \hat{\boldsymbol{\theta}}_h^i(\boldsymbol{\theta}_{h+1}^{i*}) - \hat{\boldsymbol{\theta}}_h^i(\boldsymbol{\theta}_{h+1}^{i*})$, then

$$\bar{\boldsymbol{\theta}}_h^i = \hat{\boldsymbol{\theta}}_h^i(\bar{\boldsymbol{\theta}}_{h+1}^i) + \bar{\boldsymbol{\xi}}_h^i \quad (147)$$

$$= \hat{\boldsymbol{\theta}}_h^i(\boldsymbol{\theta}_{h+1}^{i*}) + \hat{\boldsymbol{\theta}}_h^i(\boldsymbol{\theta}_{h+1}^{i*}) - \hat{\boldsymbol{\theta}}_h^i(\boldsymbol{\theta}_{h+1}^{i*}) \quad (148)$$

$$= \boldsymbol{\theta}_h^{i*} \quad (149)$$

Finally, we can verify $(\bar{\boldsymbol{\theta}}_h^1, \dots, \bar{\boldsymbol{\theta}}_h^M) \in \Theta_h$ from $(\boldsymbol{\theta}_h^{1*}, \dots, \boldsymbol{\theta}_h^{M*}) \in \Theta_h$.

Since $\bar{\theta}_{1t}^i$ is the optimal solution, we can finish the proof by showing

$$\sum_{i=1}^M V_1^i(\bar{\theta}_{1t}^i)(s_{1t}^i) = \sum_{i=1}^M \max_a \phi(s_{1t}^i, a)^\top \bar{\theta}_{1t}^i \quad (150)$$

$$\geq \sum_{i=1}^M \max_a \phi(s_{1t}^i, a)^\top \theta_1^{i*} \quad (\text{since } \theta_1^{i*} \text{ is the feasible solution}) \quad (151)$$

$$\geq \sum_{i=1}^M \phi(s_{1t}^i, \pi_1^{i*}(s_{1t}^i))^\top \theta_1^{i*} \quad (152)$$

$$\geq \sum_{i=1}^M Q_h^{i*}(s_{1t}^i, \pi_1^{i*}(s_{1t}^i)) - MHT \quad (\text{by Lemma 11}) \quad (153)$$

$$\geq \sum_{i=1}^M V_h^{i*}(s_{1t}^i) - MHT \quad (154)$$

□

B.5. REGRET BOUND

We are ready to present the proof of our regret bound.

From Lemma 8 we know that the failure event $\bigcup_{t=1}^T \bigcup_{h=1}^H E_{ht}$ happens with probability at most $\delta/2$, so we assume it does not happen. Then we can decompose the regret as

$$\text{Reg}(T) = \sum_{t=1}^T \sum_{i=1}^M \left(V_1^{i*} - V_1^{\pi_t^i} \right) (s_{1t}^i) \quad (155)$$

$$= \sum_{t=1}^T \sum_{i=1}^M \left(V_1^{i*} - V_1^i(\bar{\theta}_{1t}^i) \right) (s_{1t}^i) + \sum_{t=1}^T \sum_{i=1}^M \left(V_1^i(\bar{\theta}_{1t}^i) - V_1^{\pi_t^i} \right) (s_{1t}^i) \quad (156)$$

$$\leq \sum_{t=1}^T \sum_{i=1}^M \left(V_1^i(\bar{\theta}_{1t}^i) - V_1^{\pi_t^i} \right) (s_{1t}^i) + MHTI \quad (\text{by Lemma 12}) \quad (157)$$

Let $a_{ht}^i = \pi_t^i(s_{ht}^i)$, and denote $Q_h^i(\bar{\theta}_{ht}^i)(V_h^i(\bar{\theta}_{ht}^i))$ by $\bar{Q}_{ht}^i(\bar{V}_{ht}^i)$ for short, we have

$$\sum_{i=1}^M \left(\bar{V}_{ht}^i - V_h^{\pi_t^i} \right) (s_{ht}^i) = \sum_{i=1}^M \left(\bar{Q}_{ht}^i - Q_h^{\pi_t^i} \right) (s_{ht}^i, a_{ht}^i) \quad (158)$$

$$= \sum_{i=1}^M \left(\bar{Q}_{ht}^i - \mathcal{T}_h^i \bar{Q}_{h+1,t}^i \right) (s_{ht}^i, a_{ht}^i) + \sum_{i=1}^M \left(\mathcal{T}_h^i \bar{Q}_{h+1,t}^i - Q_h^{\pi_t^i} \right) (s_{ht}^i, a_{ht}^i) \quad (159)$$

$$\leq M\mathcal{I} + 2\sqrt{\alpha_{ht} \cdot \sum_{i=1}^M \|\phi(s_{ht}^i, a_{ht}^i)\|_{\bar{V}_{ht}^i(\lambda)^{-1}}^2} + \sum_{i=1}^M \mathbb{E}_{s' \sim p_h^i(s_{ht}^i, a_{ht}^i)} \left[\left(\bar{V}_{h+1,t}^i - V_{h+1}^{\pi_t^i} \right) (s') \right] \quad (160)$$

$$\leq \sum_{i=1}^M \left(\bar{V}_{h+1,t}^i - V_{h+1}^{\pi_t^i} \right) (s_{h+1,t}^i) + M\mathcal{I} + 2\sqrt{\alpha_{ht} \cdot \sum_{i=1}^M \|\phi(s_{ht}^i, a_{ht}^i)\|_{\bar{V}_{ht}^i(\lambda)^{-1}}^2} + \sum_{i=1}^M \zeta_{ht}^i \quad (161)$$

where ζ_{ht}^i is a martingale difference with regards to the filtration $\mathcal{F}_{h,t}$ defined as

$$\zeta_{ht}^i \stackrel{\text{def}}{=} \left(\bar{V}_{h+1,t}^i - V_{h+1}^{\pi_t^i} \right) (s_{h+1,t}^i) - \mathbb{E}_{s' \sim p_h^i(s_{ht}^i, a_{ht}^i)} \left[\left(\bar{V}_{h+1,t}^i - V_{h+1}^{\pi_t^i} \right) (s') \right] \quad (162)$$

According to assumption 3 we know $|\zeta_{ht}^i| \leq 4$, so we can apply Azuma-Hoeffding's inequality that with probability $1 - \delta/2$ for any $t \in [T]$ and $i \in [M]$

$$\sum_{j=1}^t \zeta_{ht}^i \leq 4\sqrt{2t \ln \left(\frac{2T}{\delta} \right)} \quad (163)$$

By applying inequality 161 recursively, we can bound the regret as

$$\text{Reg}(T) \leq \sum_{t=1}^T \sum_{i=1}^M \left(\bar{V}_{1t}^i - V_1^{\pi^i} \right) (s_{1t}^i) + MHT\mathcal{I} \quad (164)$$

$$\leq 2MHT\mathcal{I} + \sum_{t=1}^T \sum_{h=1}^H 2 \sqrt{\alpha_{ht} \cdot \sum_{i=1}^M \|\phi(s_{ht}^i, a_{ht}^i)\|_{\bar{\mathbf{V}}_{ht}^i(\lambda)^{-1}}^2} + \sum_{i=1}^M \sum_{h=1}^H \sum_{t=1}^T \zeta_{ht}^i \quad (165)$$

The last inequality is due to $\bar{V}_{H+1}^i(s) = \max_a \phi(s, a)^\top \bar{\boldsymbol{\theta}}_{H+1,t}^i = 0$, $V_{H+1}^{\pi^i}(s) = 0$.

The Lemma 11 of Abbasi-Yadkori et al. (2011) gives that for any $i \in [M]$ and $h \in [H]$

$$\sum_{t=1}^T \|\phi(s_{ht}^i, a_{ht}^i)\|_{\bar{\mathbf{V}}_{ht}^i(\lambda)^{-1}}^2 = \tilde{O}(d) \quad (166)$$

Moreover, by the definition of α_{ht} (see Lemma 9) we know that for any $h \in [H]$ and $t \in [T]$

$$\alpha_{ht} = \tilde{O}(Mk + kd + MTT^2) \quad (167)$$

Take all of above we can show the final regret bound.

$$\text{Reg}(T) \leq 2MHT\mathcal{I} + \sum_{t=1}^T \sum_{h=1}^H 2 \sqrt{\alpha_{ht} \cdot \sum_{i=1}^M \|\phi(s_{ht}^i, a_{ht}^i)\|_{\bar{\mathbf{V}}_{ht}^i(\lambda)^{-1}}^2} + \sum_{i=1}^M \sum_{h=1}^H \sum_{t=1}^T \zeta_{ht}^i \quad (168)$$

$$= \tilde{O} \left(MHT\mathcal{I} + \tilde{O} \left(\sqrt{Mk + kd + MTT^2} \right) \sum_{h=1}^H \sum_{t=1}^T \sqrt{\sum_{i=1}^M \|\phi(s_{ht}^i, a_{ht}^i)\|_{\bar{\mathbf{V}}_{ht}^i(\lambda)^{-1}}^2} + MH\sqrt{T} \right) \quad (169)$$

$$= \tilde{O} \left(MHT\mathcal{I} + \tilde{O} \left(\sqrt{Mk + kd + MTT^2} \right) \sum_{h=1}^H \sqrt{T} \cdot \sqrt{\sum_{t=1}^T \sum_{i=1}^M \|\phi(s_{ht}^i, a_{ht}^i)\|_{\bar{\mathbf{V}}_{ht}^i(\lambda)^{-1}}^2} + MH\sqrt{T} \right) \quad (170)$$

$$= \tilde{O} \left(MHT\mathcal{I} + \tilde{O} \left(\tilde{O} \left(\sqrt{Mk + kd + MTT^2} \right) \cdot H\sqrt{MTd} \right) + MH\sqrt{T} \right) \quad (171)$$

$$= \tilde{O} \left(HM\sqrt{dkT} + Hd\sqrt{MkT} + HMT\sqrt{d\mathcal{I}} \right) \quad (172)$$

C. Proof of Theorem 5

To prove the lower bound for multi-task RL, our idea is to connect the lower bound for the multi-task learning problem to the lower bound in the single-task LSVI setting (Zanette et al., 2020a). In the paper of Zanette et al. (2020a), they assumed the feature dimension d can be varied among different steps, which is denoted as d_h for step h . They proved the lower bound for linear RL in this setting is $\Omega\left(\sum_{h=1}^H d_h \sqrt{T} + \sum_{h=1}^H \sqrt{d_h} \mathcal{I}T\right)$. However, this lower bound is derived by the hard instance with $d_1 = \sum_{h=2}^H d_h$. If we set $d_1 = d_2 = \dots = d_H = d$ like our setting, we can only obtain the lower bound of $\Omega\left(d\sqrt{T} + \sqrt{d} \mathcal{I}T\right)$ following their proof idea. In fact, the dependence on H in this lower bound can be further improved. In order to obtain a tighter lower bound, we consider the lower bound for single-task misspecified linear MDP. This setting can be proved to be strictly simpler than the LSVI setting following the idea of Proposition 3 in Zanette et al. (2020a). The lower bound for misspecified linear MDP can thus be applied to LSVI setting.

C.1. LOWER BOUNDS FOR SINGLE-TASK RL

This subsection focus on the lower bound for misspecified linear MDP setting, in which the transition kernel and the reward function are assume to be approximately linear.

Assumption 5. (Assumption B in Jin et al. (2020)) For any $\zeta \leq 1$, we say that $\text{MDP}(\mathcal{S}, \mathcal{A}, p, r, H)$ is a ζ -approximate linear MDP with a feature map $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$, if for any $h \in [H]$, there exist d unknown measures $\theta_h = (\theta_h^{(1)}, \dots, \theta_h^{(d)})$ over \mathcal{S} and an unknown vector $\nu_h \in \mathbb{R}^d$ such that for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, we have

$$\|p_h(\cdot|s, a) - \langle \phi(s, a), \theta_h(\cdot) \rangle\|_{\text{TV}} \leq \zeta \quad (173)$$

$$|r_h(s, a) - \langle \phi(s, a), \nu_h \rangle| \leq \zeta \quad (174)$$

For regularity, we assume that Assumption 3 still holds, and we also assume that there exists a constant D such that $\|\theta_h(s)\| \leq D$ for all $s \in \mathcal{S}, h \in [H]$, $\|\nu_h\| \leq D$ for all $h \in [H]$. $D \geq 4$ suffices in our hard instance construction.

For misspecified linear MDP, we can prove the following lower bound.

Proposition 1. Suppose $T \geq \frac{d^2 H}{4}$, $d \geq 10$, $H \geq 10$ and $\zeta \leq \frac{1}{4H}$, there exist a ζ -approximate linear MDP class such that the expected regret of any algorithm on at least a member of the MDP class is at least $\Omega\left(d\sqrt{HT} + HT\sqrt{d}\right)$.

To prove the lower bound, our basic idea is to connect the problem to $\frac{H}{2}$ linear bandit problems. Similar hard instance construction has been used in Zhou et al. (2020a;b). In our construction, the state space \mathcal{S} consists of $H + 2$ states, which is denoted as x_1, x_2, \dots, x_{H+2} . The agent starts the episode in state x_1 . In x_h , it can either transits to x_{h+1} or x_{H+2} with certain transition probability. If the agent enters x_{H+2} , it will stay in this state in the remaining steps, i.e. x_{H+2} is an absorbing state. For each state, there are 2^{d-4} actions and $\mathcal{A} = \{-1, 1\}^{d-4}$. Suppose the agent takes action $\mathbf{a} \in \{-1, 1\}^{d-4}$ in state s_h , the transition probability to state s_{h+1} and s_{H+2} is $1 - \zeta_h(\mathbf{a}) - \delta - \mu_h^\top \mathbf{a}$ and $\delta + \zeta_h(\mathbf{a}) + \mu_h^\top \mathbf{a}$ respectively. Here $|\zeta_h(\mathbf{a})| \leq \zeta$ denotes the approximation error of linear representation, $\delta = 1/H$ and $\mu_h \in \{-\Delta, \Delta\}^{d-4}$ with $\Delta = \sqrt{\delta/T}/(4\sqrt{2})$ so that the probability is well-defined. The reward can only be obtained in x_{H+2} , with $r_h(x_{H+2}, a) = 1/H$ for any h, a . We assume the reward to be deterministic.

We can check that this construction satisfies Assumption 5 with ϕ and θ defined in the following way:

$$\phi(s, \mathbf{a}) = \begin{cases} (0, \alpha, \alpha\delta, 0, \beta\mathbf{a}^\top)^\top & s = x_1, x_2, \dots, x_H \\ (0, 0, 0, \alpha, \mathbf{0}^\top)^\top & s = x_{H+1} \\ (\alpha, 0, 0, \alpha, \mathbf{0}^\top)^\top & s = x_{H+2} \end{cases}$$

$$\theta_h(s') = \begin{cases} \left(0, \frac{1}{\alpha}, -\frac{1}{\alpha}, 0, -\frac{\mu_h^\top}{\beta}\right)^\top & s' = x_{h+1} \\ \left(0, 0, \frac{1}{\alpha}, \frac{1}{\alpha}, \frac{\mu_h^\top}{\beta}\right)^\top & s = x_{H+2} \\ \mathbf{0} & \text{otherwise} \end{cases}$$

ν_h is defined to be $(\frac{1}{H\alpha}, \mathbf{0}^\top)^\top$, and $\alpha = \sqrt{1/(2 + \Delta(d-4))}$, $\beta = \sqrt{\Delta/(2 + \Delta(d-4))}$. Note that $\|\phi(s, a)\| \leq 1$, $\|\theta_h(s')\| \leq D$ and $\|\nu_h\| \leq D$ hold for any s, a, s', h when $T \geq d^2 H/4$.

Since the rewarding state is only x_{H+2} , the optimal strategy in state x_h ($h \leq H$) is to take an action that maximizes the probability of entering x_{H+2} , i.e., to maximize $\mu_h^\top \mathbf{a} + \zeta(\mathbf{a})$. That is to say, we can regard the problem of finding the optimal action in state s_h and step h as finding the optimal arm for a $d - 4$ -dimensional approximately (misspecified) linear bandits problem. Thanks to the choice of δ such that $(1 - \delta)^{H/2}$ is a constant, there is sufficiently high probability of entering state x_h for any $h \leq H/2$. Therefore, we can show that this problem is harder than solving $H/2$ misspecified linear bandit problems. This following lemma characterizes this intuition. The lemma follows the same idea of Lemma C.7 in Zhou et al. (2020a), though our setting is more difficult since we consider misspecified case.

Lemma 13. Suppose $H \geq 10$, $d \geq 10$ and $(d-4)\Delta \leq \frac{1}{2H}$. We define $r_h^b(\mathbf{a}) = \boldsymbol{\mu}^\top \mathbf{a} + \zeta_h(\mathbf{a})$, which can be regarded as the corresponding reward for the equivalent linear bandit problem in step h . Fix $\boldsymbol{\mu} \in (\{-\Delta, \Delta\}^{d-4})^H$. Fix a possibly history dependent policy π . Letting V^* and V^π be the optimal value function and the value function of policy π respectively, we have

$$V_1^*(s_1) - V_1^\pi(s_1) \geq 0.02 \sum_{h=1}^{H/2} \left(\max_{\mathbf{a} \in \mathcal{A}} r_h^b(\mathbf{a}) - \sum_{\mathbf{a} \in \mathcal{A}} \pi_h(\mathbf{a}|s_h) r_h^b(\mathbf{a}) \right) \quad (175)$$

Proof. Note that the only rewarding state is x_{H+2} with $r_h(x_{H+2}, \mathbf{a}) = \frac{1}{H}$. Therefore, the value function of a certain policy π can be calculated as:

$$V_1^\pi(x_1) = \sum_{h=1}^{H-1} \frac{H-h}{H} \mathbb{P}(N_h|\pi) \quad (176)$$

where N_h denotes the event of visiting state x_h in step h and then transits to x_{H+2} , i.e. $N_h = \{s_h = x_h, s_{h+1} = x_{H+2}\}$. Suppose $\omega_h^\pi = \sum_{\mathbf{a} \in \mathcal{A}} \pi_h(\mathbf{a}|s_h) r_h^b(\mathbf{a})$ and $\omega_h^* = \max_{\mathbf{a} \in \mathcal{A}} r_h^b(\mathbf{a})$. By the law of total probability and the Markov property, we have

$$\mathbb{P}(N_h|\pi) = (\delta + \omega_h^\pi) \prod_{j=1}^{h-1} (1 - \delta - \omega_j^\pi) \quad (177)$$

Thus we have

$$V_1^\pi(x_1) = \sum_{h=1}^{H-1} \frac{H-h}{H} (\delta + \omega_h^\pi) \prod_{j=1}^{h-1} (1 - \delta - \omega_j^\pi) \quad (178)$$

Similarly, for the value function of the optimal policy, we have

$$V_1^*(x_1) = \sum_{h=1}^{H-1} \frac{H-h}{H} (\delta + \omega_h^*) \prod_{j=1}^{h-1} (1 - \delta - \omega_j^*) \quad (179)$$

Define $S_i = \sum_{h=i}^{H-1} \frac{H-h}{H} (\delta + \omega_h^\pi) \prod_{j=i}^{h-1} (1 - \delta - \omega_j^\pi)$ and $T_i = \sum_{h=i}^{H-1} \frac{H-h}{H} (\delta + \omega_h^*) \prod_{j=i}^{h-1} (1 - \delta - \omega_j^*)$. Then we have $V_1^*(x_1) - V_1^\pi(x_1) = T_1 - S_1$. Notice that

$$S_i = \frac{H-i}{H} (\omega_i^\pi + \delta) + S_{i+1} (1 - \omega_i^\pi - \delta) \quad (180)$$

$$T_i = \frac{H-i}{H} (\omega_i^* + \delta) + T_{i+1} (1 - \omega_i^* - \delta) \quad (181)$$

Thus we have

$$T_i - S_i = \left(\frac{H-i}{H} - T_{i+1} \right) (\omega_i^* - \omega_i^\pi) + (T_{i+1} - S_{i+1}) (1 - \omega_i^\pi - \delta) \quad (182)$$

By induction, we get

$$T_1 - S_1 = \sum_{h=1}^{H-1} (\omega_h^* - \omega_h^\pi) \left(\frac{H-h}{H} - T_{h+1} \right) \prod_{j=1}^{h-1} (1 - \omega_j^\pi - \delta) \quad (183)$$

Since the reward is non-negative and only occurs in x_{H+2} , we know that $V_1^*(x_1) \geq V_2^*(x_2) \geq \dots \geq V_1^*(x_H)$. Thus we have $T_h \leq T_1 = V_1^*(x_1) \leq \sum_{h=1}^H \mathbb{P}(N_h|\pi^*)$. If N_h doesn't happen for any $h \in [H]$, then the agent must enter x_{H+1} . The

probability of this event has the following form:

$$\mathbb{P}\left(\neg\left(\cup_{h \in [H]} N_h | \pi^*\right)\right) = 1 - \prod_{h=1}^H \mathbb{P}(N_h | \pi^*) \quad (184)$$

$$= \prod_{h \in [H]} (1 - \delta - \omega_h^*) \quad (185)$$

$$\geq \prod_{h \in [H]} \left(1 - \frac{1}{H} + \frac{1}{2H}\right) \quad (186)$$

$$= \left(1 - \frac{1}{2H}\right)^H \quad (187)$$

$$\geq 0.6 \quad (188)$$

The first inequality is due to $\delta = \frac{2}{H}$ and $|\omega_h^*| \leq \frac{1}{H}$. The above discussion indicates that $T_h \leq 0.4$, thus $\frac{H-h}{H} - T_{h+1} \geq 0.1$ for $h \leq H/2$. Similarly, $\prod_{j=1}^{h-1} (1 - \omega_j^\pi - \delta) \geq (1 - \frac{3}{2H})^{H-1} \geq 0.2$. Combining with Eqn 183, we have

$$T_1 - S_1 \geq 0.02 \sum_{h=1}^{\frac{H}{2}} (\omega_h^* - \omega_h^\pi) = 0.02 \sum_{h=1}^{H/2} \left(\max_{\mathbf{a} \in \mathcal{A}} r_h^b(\mathbf{a}) - \sum_{\mathbf{a} \in \mathcal{A}} \pi_h(\mathbf{a} | s_h) r_h^b(\mathbf{a}) \right) \quad (189)$$

Combining with the definition of T_1 and S_1 , we can prove the lemma. \square

After proving Lemma 13, we are ready to prove Proposition 1.

Proof. (proof of Proposition 1) By Lemma 13, we know that we can decompose the sub-optimality gap of a policy π in the following way:

$$V_1^*(s_1) - V_1^\pi(s_1) \geq 0.02 \sum_{h=1}^{H/2} \left(\max_{\mathbf{a} \in \mathcal{A}} r_h^b(\mathbf{a}) - \sum_{\mathbf{a} \in \mathcal{A}} \pi_h(\mathbf{a} | s_h) r_h^b(\mathbf{a}) \right) \quad (190)$$

where $r_h^b(\mathbf{a}) = \boldsymbol{\mu}^\top \mathbf{a} + \zeta_h(\mathbf{a})$, which can be regarded as a reward function for misspecified linear bandit. To prove Theorem 1, the only remaining problem is to derive the lower bound for misspecified linear bandits. We directly apply the following two lower bounds for linear bandits.

Lemma 14. (Lemma C.8 in Zhou et al. (2020a)) Fix a positive real $0 < \delta \leq 1/3$, and positive integers T, d and assume that $T \geq d^2/(2\delta)$ and consider the linear bandit problem \mathcal{L}_μ parametrized with a parameter vector $\boldsymbol{\mu} \in \{-\Delta, \Delta\}^d$ and action set $\mathcal{A} = \{-1, 1\}^d$ so that the reward distribution for taking action $\mathbf{a} \in \mathcal{A}$ is a Bernoulli distribution $B(\delta + (\boldsymbol{\mu}^*)^\top \mathbf{a})$. Then for any bandit algorithm \mathcal{B} , there exists a $\boldsymbol{\mu}^* \in \{-\Delta, \Delta\}^d$ such that the expected pseudo-regret of \mathcal{B} over T steps on bandit $\mathcal{L}_{\boldsymbol{\mu}^*}$ is lower bounded by $\frac{d\sqrt{T\delta}}{8\sqrt{2}}$.

Lemma 15. (Proposition 6 in Zanette et al. (2020a)) There exists a feature map $\phi : \mathcal{A} \rightarrow \mathbb{R}^d$ that defines a misspecified linear bandits class \mathcal{M} such that every bandit instance in that class has reward response:

$$\mu_a = \phi_a^\top \theta + z_a$$

for any action a (Here $z_a \in [0, \zeta]$ is the deviation from linearity and $\mu_a \in [0, 1]$) and such that the expected regret of any algorithm on at least a member of the class up to round T is $\Omega(\sqrt{d\zeta T})$.

Lemma 14 is used to prove the lower bound for linear mixture MDPs in Zhou et al. (2020a), which states that the lower bound for linear bandits with approximation error $\zeta = 0$, while Lemma 15 mainly consider the influence of ζ to the lower bound. Combining these two lemmas, the regret lower bound for misspecified linear bandit is $\Omega(\max(d\sqrt{T\delta}, \sqrt{d\zeta T})) = \Omega(d\sqrt{T\delta} + \sqrt{d\zeta T})$. Since here our problem can reduce from $H/2$ misspecified linear bandit, we know that the regret lower bound is $\Omega(Hd\sqrt{T\delta} + H\sqrt{d\zeta T}) = \Omega(d\sqrt{HT} + H\sqrt{d\zeta T})$ \square

Now we obtain the regret lower bound for misspecified linear MDP. We can prove the corresponding lower bound for the LSVI setting [Zanette et al. \(2020a\)](#) since LSVI setting is strictly harder than linear MDP setting. The following lemma states this relation between two settings.

Lemma 16. *If an MDP $(\mathcal{S}, \mathcal{A}, p, r, H)$ is a misspecified linear MDP with approximation error ζ , then this MDP satisfies the low inherent Bellman error assumption with $\mathcal{I} = 2\zeta$.*

Proof. If an MDP is an ζ -approximate linear MDP, then we have

$$\|p_h(\cdot|s, a) - \langle \phi(s, a), \boldsymbol{\theta}_h(\cdot) \rangle\|_{\text{TV}} \leq \zeta \quad (191)$$

$$|r_h(s, a) - \langle \phi(s, a), \boldsymbol{\nu}_h \rangle| \leq \zeta \quad (192)$$

For any $\theta_{h+1} \in \mathbb{R}^d$, we have $\mathcal{T}_h(Q_{h+1}(\theta_{h+1}))(s, a) = r_h(s, a) + \mathbb{E}_{s' \sim p_h(\cdot|s, a)} V_{h+1}(\theta_{h+1})(s')$. Since $V_{h+1}(\theta_{h+1})(s') \leq 1$, plugging the approximately linear form of $r_h(s, a)$ and $p_h(\cdot|s, a)$, we have

$$|\mathcal{T}_h(Q_{h+1}(\theta_{h+1}))(s, a) - \left\langle \phi(s, a), \sum_{s'} \boldsymbol{\theta}_h(s') V_{h+1}(\theta_{h+1})(s') + \boldsymbol{\nu}_h \right\rangle| \leq 2\zeta \quad (193)$$

□

By lemma 16, we can directly apply the hard instance construction and the lower bound for misspecified linear MDP to LSVI setting.

Proposition 2. *There exist function feature maps ϕ_1, \dots, ϕ_H that define an MDP class \mathcal{M} such that every MDP in that class satisfies low inherent Bellman error at most \mathcal{I} and such that the expected reward on at least a member of the class (for $|\mathcal{A}| \geq 3, d, k, H \geq 10, T = \Omega(d^2 H), \mathcal{I} \leq \frac{1}{4H}$) is $\Omega(d\sqrt{HT} + \sqrt{dHT})$.*

C.2. LOWER BOUND FOR MULTI-TASK RL

In order to prove Theorem 5, we need to prove and then combine the following two lemmas.

Lemma 17. *Under the setting of Theorem 5, the expected regret of any algorithm \mathcal{A} is lower bounded by $\Omega(Mk\sqrt{HT})$.*

Lemma 18. *Under the setting of Theorem 5, the expected regret of any algorithm \mathcal{A} is lower bounded by $\Omega(d\sqrt{kMHT} + HMT\sqrt{d\mathcal{I}})$.*

These two lemmas are proved by reduction from Proposition 2, which is a lower bound we proved for the single-task LSVI setting.

Proof. (Proof of Lemma 17) The lemma is proved by contradiction. Suppose there is an algorithm \mathcal{A} that achieves $\sup_{M \in \mathcal{M}} \mathbb{E}[\text{Reg}(T)] \leq CMk\sqrt{HT}$ for a constant C . Then there must exist a task $i \in [M]$, such that the expected regret for this single task is at most $Ck\sqrt{HT}$. However, by Proposition 2, the expected regret for MDPs with dimension k in horizon h is at least $\Omega(k\sqrt{HT} + \sqrt{kHTT})$. This leads to a contradiction. \square

Proof. (Proof of Lemma 18) The hard instance construction follows the same idea of the proof for our Lemma 6, as well as the hard instance to prove Lemma 19 in Yang et al. (2020). Without loss of generality, we assume that M can be exactly divided by k .

We divide M tasks into k groups. Each group shares the same parameter $\{\theta_h^i\}_{h=1}^H$. To be more specific, we let $w_h^1 = w_h^2 = \dots = w_h^{M/k} = e_h^1$, $w_h^{M/k+1} = w_h^{M/k+2} = \dots = w_h^{2M/k} = e_h^2$, \dots , $w_h^{(k-1)M/k+1} = w_h^{(k-1)M/k+2} = \dots = w_h^M = e_h^k$. Under this construction, the parameters θ_h^i for these tasks are exactly the same in each group, but relatively independent among different groups. That is to say, the expected regret lower bound is at least the summation of the regret lower bounds in all k groups.

Now we consider the regret lower bound for group $j \in [k]$. Since the parameters are shared in the same group, the regret of running an algorithm for M/k tasks with T episodes each is at least the regret of running an algorithm for single-task linear bandit with $M/k \cdot T$ episodes. By Proposition 2, the regret for single-task linear bandit with MT/k episodes is at least $\Omega(d\sqrt{MHT/k} + \sqrt{d\mathcal{I}HMT/k})$. Summing over all k groups, we can prove that the regret lower bound is $\Omega(d\sqrt{kHMT} + \sqrt{d\mathcal{I}HMT})$. \square