

A. Proof of Theorem 1

To start, we list some technical lemmas. Then we prove that our estimation of the Q -function is optimistic with high probability, *i.e.*, $\tilde{Q}_h^k(s, a) \geq Q_h^*(s, a)$. We begin by proving Lemma 1 and cite a lemma in (Zhang et al., 2020a).

Lemma 2 (Bennet's inequality.). *Let Z, Z_1, \dots, Z_n be i.i.d. random variables with values in $[0, 1]$ and let $\delta > 0$. Define $\mathbb{V}Z = \mathbb{E}[(Z - \mathbb{E}Z)^2]$. Then we have*

$$\mathbb{P}\left[\left|\mathbb{E}[Z] - \frac{1}{n} \sum_{i=1}^n Z_i\right| > \sqrt{\frac{2\mathbb{V}Z \ln(2/\delta)}{n}} + \frac{\ln(2/\delta)}{n}\right] \leq \delta.$$

Lemma 3 (Lemma 10 in (Zhang et al., 2020a)). *Let $(M_n)_{n \geq 0}$ be a martingale such that $M_0 = 0$ and $|M_n - M_{n-1}| \leq c$ for some $c > 0$ and any $n \geq 1$. Let $\text{Var}_n = \sum_{k=1}^n \mathbb{E}[(M_k - M_{k-1})^2 | \mathcal{F}_{k-1}]$ for $n \geq 0$, where $\mathcal{F}_k = \sigma(M_1, \dots, M_k)$. Then for any positive integer n , and any $\epsilon, \delta > 0$, we have that*

$$\mathbb{P}[|M_n| \geq 2\sqrt{2}\sqrt{\text{Var}_n \ln(1/\delta)} + 2\sqrt{\epsilon \ln(1/\delta)} + 2c \ln(1/\delta)] \leq 2(\log_2(\frac{nc^2}{\epsilon}) + 1)\delta.$$

A.1. Proof of Lemma 1

Proof. Recall that $Q_h(s, a) - \hat{Q}_h(s, a) = (\tilde{r}(s, a) - \hat{r}(s, a)) + (\tilde{P}_{s,a} - \hat{P}_{s,a})V_{h+1} + (\tilde{b}_h(s, a) - \hat{b}_h(s, a))$, where

$$\begin{aligned} \hat{b}_h &= c_1 \sqrt{\frac{\mathbb{V}(\hat{P}, V_{h+1})\iota}{\max\{\hat{n}(s, a), 1\}}} + c_2 \sqrt{\frac{\hat{r}(s, a)\iota}{\max\{\hat{n}(s, a), 1\}}} + c_3 \frac{\iota}{\max\{\hat{n}(s, a), 1\}}, \\ \tilde{b}_h &= c_1 \min\left\{\sqrt{\frac{\mathbb{V}(\tilde{P}, V_{h+1})\iota}{|\tilde{n}(s, a) - C|}}, 1\right\} + c_2 \min\left\{\sqrt{\frac{\tilde{r}\iota}{|\tilde{n}(s, a) - C|}}, 1\right\} + c_3 \min\left\{\frac{\iota}{|\tilde{n}(s, a) - C|}, 1\right\} \\ &\quad + 2 \min\left\{\frac{2C}{|\tilde{n}(s, a) - C|}, 1\right\} + (c_1 + c_2) \min\left\{\frac{\sqrt{C}\iota}{|\tilde{n}(s, a) - C|}, 1\right\}. \end{aligned}$$

We prove the next lemma that characterize the the difference between biased and unbiased estimators,

Lemma 4. *For any vector $V \in R^S$, $V(s) \in [0, 1]$ for any $s \in S$, it holds that*

$$\|\tilde{P}_{s,a} - \hat{P}_{s,a}\|_1 \leq 2 \min\left\{\frac{C}{|\tilde{n}(s, a) - C|}, 1\right\}, \quad (4)$$

$$|\mathbb{V}(\tilde{P}_{s,a}, V) - \mathbb{V}(\hat{P}_{s,a}, V)| \leq 6 \min\left\{\frac{C}{|\tilde{n}(s, a) - C|}, 1\right\}, \quad (5)$$

$$|\tilde{r} - \hat{r}| \leq \min\left\{\frac{C}{|\tilde{n} - C|}, 1\right\}. \quad (6)$$

For (4), we calculate that

$$\begin{aligned} \|\tilde{P}_{s,a} - \hat{P}_{s,a}\|_1 &= \sum_{s'} |\tilde{P}_{s,a,s'} - \hat{P}_{s,a,s'}| \\ &= \sum_{s'} \left| \frac{\tilde{n}(s, a, s')}{\tilde{n}(s, a)} - \frac{\hat{n}(s, a, s')}{\hat{n}(s, a)} \right| \\ &\leq \sum_{s'} \left| \frac{\tilde{n}(s, a, s')}{\tilde{n}(s, a)} - \frac{\hat{n}(s, a, s')}{\tilde{n}(s, a) - C} \right| \\ &\leq \sum_{s'} \frac{\tilde{n}(s, a) |\tilde{n}(s, a, s') - \hat{n}(s, a, s')|}{\tilde{n}(s, a) |\tilde{n}(s, a) - C|} + \frac{C \tilde{n}(s, a, s')}{\tilde{n}(s, a) |\tilde{n}(s, a) - C|} \\ &\leq \frac{2C}{|\tilde{n}(s, a) - C|}. \end{aligned}$$

Here the second inequality holds because $\sum_{s'} |\hat{n}(s, a, s') - \tilde{n}(s, a, s')| \leq C$ for any s, a . Note that $\|\tilde{P}_{s,a} - \hat{P}_{s,a}\|_1 \leq 2$, we have (4) holds. For (5),

$$\begin{aligned} & |\mathbb{V}(\tilde{P}_{s,a}, V) - \mathbb{V}(\hat{P}_{s,a}, V)| \\ & \leq |\tilde{P}_{s,a}(V)^2 - \hat{P}_{s,a}(V)^2| + |(\tilde{P}_{s,a}V)^2 - (\hat{P}_{s,a}V)^2| \\ & \leq \|\tilde{P}_{s,a} - \hat{P}_{s,a}\|_1 + \|\tilde{P}_{s,a} - \hat{P}_{s,a}\|_1(\tilde{P}V + \hat{P}V) \\ & \leq 3\|\tilde{P}_{s,a} - \hat{P}_{s,a}\|_1 \end{aligned}$$

Hence by (4), we finished the proof of (5). (6) comes straight from the definition of \tilde{r} and \hat{r} . Finally, by plugging (4), (5) and (6) into $Q_h - \hat{Q}_h$ and using $(\tilde{P}_{s,a} - \hat{P}_{s,a})V_{h+1} \leq \|\tilde{P}_{s,a} - \hat{P}_{s,a}\|_1$, we finished the proof of Lemma 1. \square

Recall that the Lemma 1 in (Zhang et al., 2020a) states that $\hat{Q}_h^k(s, a) \geq Q_h^*(s, a)$ holds with probability $1 - 2SA(\log_2(KH) + 1)\delta$, as a result, $Q_h^k(s, a) \geq Q_h^*(s, a)$, which implies that our estimation of Q -function is optimistic. We denote by \mathcal{E}_1 the event that $Q_h^k(s, a) \geq Q_h^*(s, a)$ for any (k, h, s, a) .

A.2. Bounding Bellman Error

We begin with the following lemma that bounds the Bellman error induced by the Q -function.

Lemma 5. *With probability $1 - 3S^2AH(\log_2(KH) + 1)\delta$, for any $1 \leq k \leq K$, $1 \leq h \leq H$ and (s, a) , it holds that*

$$\begin{aligned} & Q_h^k(s, a) - r(s, a) - P_{s,a}V_{h+1}^k \\ & \leq \min\left\{2\tilde{b}_h^k(s, a) + \frac{2C}{\tilde{n}(s, a) + C} + \sqrt{\frac{2\mathbb{V}(P_{s,a}, V_{h+1}^*)^\ell}{\hat{n}^k(s, a)}} + \sqrt{\frac{2S\mathbb{V}(P_{s,a}, V_{h+1}^k - V_{h+1}^*)^\ell}{\hat{n}^k(s, a)}} + \frac{2S\ell}{3\hat{n}^k(s, a)}, 1\right\}. \end{aligned} \quad (7)$$

Proof. Under event \mathcal{E}_1 , we have that with probability $1 - SAH(\log_2(KH) + 1)\delta$, for all $(s, a, k, h) \in \mathcal{S} \times \mathcal{A} \times [K] \times [H]$,

$$\begin{aligned} & Q_h^k(s, a) - r(s, a) - P_{s,a}V_{h+1}^k \\ & \leq \tilde{r}_h^k(s, a) - r(s, a) + \tilde{b}_h^k(s, a) + (\tilde{P}_{s,a} - \hat{P}_{s,a})V_{h+1}^k + (\hat{P}_{s,a} - P_{s,a})(V_{h+1}^k - V_{h+1}^*) + (\hat{P}_{s,a} - P_{s,a})V_{h+1}^* \\ & \leq 2\tilde{b}_h^k(s, a) + (\tilde{P}_{s,a} - \hat{P}_{s,a})V_{h+1}^k + (\hat{P}_{s,a} - P_{s,a})(V_{h+1}^k - V_{h+1}^*) + (\hat{P}_{s,a} - P_{s,a})V_{h+1}^*. \end{aligned} \quad (8)$$

By Bennet's inequality (Lemma 2) we have that for each s' ,

$$\mathbb{P}\left[|\hat{P}_{s,a,s'} - P_{s,a,s'}| > \sqrt{\frac{2P_{s,a,s'}^\ell}{\hat{n}^k(s, a)}} + \frac{\ell}{3\hat{n}^k(s, a)}\right] \leq \delta.$$

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

$$\begin{aligned} & (\hat{P}_{s,a} - P_{s,a})(V_{h+1}^k - V_{h+1}^*) = \sum_{s'} (\hat{P}_{s,a,s'} - P_{s,a,s'})(V_{h+1}^k(s') - V_{h+1}^*(s') - P_{s,a}(V_{h+1}^k - V_{h+1}^*)) \\ & \leq \sum_{s'} \sqrt{\frac{2P_{s,a,s'}^\ell}{\hat{n}^k(s, a)}} |V_{h+1}^k(s') - V_{h+1}^*(s') - P_{s,a}(V_{h+1}^k - V_{h+1}^*)| + \frac{S\ell}{3\hat{n}^k(s, a)} \\ & \leq \sqrt{\frac{\hat{n}^k(s, a)}{2S\mathbb{V}(P_{s,a}, V_{h+1}^k - V_{h+1}^*)}} + \frac{S\ell}{3\hat{n}^k(s, a)}, \end{aligned} \quad (9)$$

where the first equality holds because $\sum_{s'} \hat{P}_{s,a,s'} = \sum_{s'} P_{s,a,s'} = 1$ and the last inequality holds by Cauchy-Schwartz inequality. By Bennet's inequality again (Lemma 2), we obtain that

$$\mathbb{P}\left[|(\hat{P}_{s,a}^k - P_{s,a})V_{h+1}^*| > \sqrt{\frac{2\mathbb{V}(P_{s,a}, V_{h+1}^*)^\ell}{\hat{n}^k(s, a)}} + \frac{\ell}{3\hat{n}^k(s, a)}\right] \leq \delta. \quad (10)$$

On the other hand, we have that

$$\begin{aligned}
 (\tilde{P} - \hat{P})V_{h+1}^k &= \sum_{s'} \left(\frac{\tilde{n}(s, a, s')}{\tilde{n}(s, a)} - \frac{\hat{n}(s, a, s')}{\hat{n}(s, a)} \right) V_{h+1}^k(s') \\
 &\leq \sum_{s'} \left| \frac{\tilde{n}(s, a, s')}{\tilde{n}(s, a)} - \frac{\hat{n}(s, a, s')}{\tilde{n}(s, a) + C} \right| \\
 &\leq \sum_{s'} \frac{\tilde{n}(s, a) |\tilde{n}(s, a, s') - \hat{n}(s, a, s')|}{\tilde{n}(s, a)(\tilde{n}(s, a) + C)} + \frac{C\tilde{n}(s, a, s')}{\tilde{n}(s, a)(\tilde{n}(s, a) + C)} \\
 &\leq \frac{2C}{\tilde{n}(s, a) + C}, \tag{11}
 \end{aligned}$$

where the first inequality holds because $|\hat{n}(s, a) - \tilde{n}(s, a)| \leq C$, and the last inequality holds because $\sum_{s'} |\tilde{n}(s, a, s') - \hat{n}(s, a, s')| \leq C$. Combining (8), (9), (10) and (11) and via a union bound over k, h, s, a , we conclude that (7) holds with probability $1 - 3S^2AH(\log_2(KH) + 1)\delta$. \square

In the rest of this section, we let $\beta_h^k(s, a)$ be the shorthand of RHS of (7), i.e.,

$$\beta_h^k(s, a) := \min \left\{ 2\tilde{b}_h^k(s, a) + \frac{2C}{\tilde{n}(s, a) + C} + \sqrt{\frac{2\mathbb{V}(P_{s,a}, V_{h+1}^* - V_{h+1}^*)\iota}{\hat{n}^k(s, a)}} + \sqrt{\frac{2S\mathbb{V}(P_{s,a}, V_{h+1}^k - V_{h+1}^*)}{\hat{n}^k(s, a)}} + \frac{2S\iota}{3\hat{n}^k(s, a)}, 1 \right\}.$$

A.3. Regret Analysis

Let \mathcal{K} be the set of indice of episodes in which no update is triggered. By the update rule, it is easy to see that $|\mathcal{K}^C| \leq SA(\log_2(KH) + 1)$. Denote $h_0(k)$ to be the first time an update is triggered in the k -th episode if there is and otherwise $H + 1$. Define $\mathcal{X}_0 = \{(k, h_0(k)) | k \in \mathcal{K}^C\}$ and $\mathcal{X} = \{(k, h) | k \in \mathcal{K}^C, h_0(k) + 1 \leq h \leq H\}$. We further define $\bar{V}_h^k(s_h^k, a_h^k) = \mathbb{I}[(k, h) \notin \mathcal{X}] \bar{V}_h^k(s_h^k, a_h^k)$. We also set $\bar{\beta}_h^k(s_h^k, a_h^k) = \mathbb{I}[(k, h) \notin \mathcal{X}] \beta_h^k(s_h^k, a_h^k)$ and $\bar{r}_h^k = \mathbb{I}[(k, h) \notin \mathcal{X}] r(s_h^k, a_h^k)$. By Lemma 5, we have that with probability $1 - 3S^2AH(\log_2(KH) + 1)\delta$,

$$\bar{V}_h^k(s_h^k, a_h^k) \leq \bar{r}_h^k + \bar{\beta}_h^k(s_h^k, a_h^k) + P_{s,a} \bar{V}_{h+1}^k,$$

for any $(k, h) \notin \mathcal{X}_0$ and

$$\bar{V}_h^k(s_h^k, a_h^k) \leq \bar{r}_h^k + \bar{\beta}_h^k(s_h^k, a_h^k) + P_{s,a} \bar{V}_{h+1}^k + 1,$$

for any $(k, h) \in \mathcal{X}_0$.

By Lemma 5, with probability at least $1 - 5S^2AH(\log_2(KH) + 1)\delta$, it holds that

$$\begin{aligned}
 \text{Regret}(K) &:= \sum_{k=1}^K (V_1^*(s_1^k) - \tilde{V}_1^{\pi^k}(s_1^k)) \\
 &\leq \sum_{k=1}^K (V_1^k(s_1^k) - \tilde{V}_1^{\pi^k}(s_1^k)) \\
 &= \sum_{k=1}^K (\bar{V}_1^k(s_1^k) - \tilde{V}_1^{\pi^k}(s_1^k)) \\
 &= \sum_{k=1}^K (\bar{V}_1^k(s_1^k) - \sum_{h=1}^H \bar{r}_h^k) + \sum_{k=1}^K (\sum_{h=1}^H \bar{r}_h^k - \tilde{V}_1^{\pi^k}(s_1^k)) \\
 &= \sum_{k=1}^K \sum_{h=1}^H (P_{s_h^k, a_h^k} \bar{V}_{h+1}^k - \bar{V}_{h+1}^k(s_{h+1}^k)) + \sum_{k=1}^K \sum_{h=1}^H (\bar{V}_h^k(s_h^k) - \bar{r}_h^k - P_{s_h^k, a_h^k} \bar{V}_{h+1}^k) + \sum_{k=1}^K (\sum_{h=1}^H \bar{r}_h^k - \tilde{V}_1^{\pi^k}(s_1^k)) \\
 &\leq \sum_{k=1}^K \sum_{h=1}^H (P_{s_h^k, a_h^k} \bar{V}_{h+1}^k - \bar{V}_{h+1}^k(s_{h+1}^k)) + \sum_{k=1}^K \sum_{h=1}^H \bar{\beta}_h^k(s_h^k, a_h^k) + \sum_{k=1}^K (\sum_{h=1}^H \bar{r}_h^k - \tilde{V}_1^{\pi^k}(s_1^k)) + |\mathcal{K}^C|. \tag{12}
 \end{aligned}$$

Here the first inequality is due to the optimism of Q -function, and the last inequality holds by Lemma 5. Define $M_1 = \sum_{k=1}^K \sum_{h=1}^H (P_{s_h^k, a_h^k} \bar{V}_{h+1}^k - \bar{V}_{h+1}^k(s_{h+1}^k))$, $M_2 = \sum_{k=1}^K \sum_{h=1}^H \bar{\beta}_h^k(s_h^k, a_h^k)$ and $M_3 = \sum_{k=1}^K (\sum_{h=1}^H \bar{r}_h^k - \tilde{V}_1^{\pi^k}(s_1^k))$.

Bounding M_1 : For the M_1 term, we note that it can be viewed as a martingale. Hence by using a variance dependent concentraion inequality (3), we have that

$$\mathbb{P}[|M_1| > 2\sqrt{2 \sum_{k=1}^K \sum_{h=1}^H \mathbb{V}(P_{s_h^k, a_h^k}, \bar{V}_{h+1}^k)\iota} + 6\iota] \leq 2(\log_2(KH) + 1)\delta.$$

So in order to bound M_1 , it suffices to bound $M_4 := \sum_{k=1}^K \sum_{h=1}^H \mathbb{V}(P_{s_h^k, a_h^k}, \bar{V}_{h+1}^k)$. We will deal with this term later.

Bounding M_3 : For the M_3 term, we have

$$\begin{aligned} M_3 &= \sum_{k=1}^K \left(\sum_{h=1}^H \bar{r}_h^k - \tilde{V}_1^{\pi^k}(s_1^k) \right) \\ &= \sum_{k=1}^K \sum_{h=1}^H (\bar{r}_h^k - r_h^k) + \sum_{k=1}^K \left(\sum_{h=1}^H r_h^k - \tilde{V}_1^{\pi^k}(s_1^k) \right) \\ &\leq \sum_{k=1}^K \sum_{h=1}^H (r(s_h^k, a_h^k) - r_h^k) + \sum_{k=1}^K \left(\sum_{h=1}^H r_h^k - \tilde{V}_1^{\pi^k}(s_1^k) \right). \end{aligned} \quad (13)$$

For the first term $\sum_{k=1}^K \sum_{h=1}^H (r(s_h^k, a_h^k) - r_h^k)$ in (13), by (3), we obtain that

$$\mathbb{P}\left[\left| \sum_{k=1}^K \sum_{h=1}^H (r(s_h^k, a_h^k) - r_h^k) \right| > 2\sqrt{2 \sum_{k=1}^K \sum_{h=1}^H \text{Var}(s, a)\iota} + 6\iota \right] \leq 2(\log_2(KH) + 1)\delta,$$

where $\text{Var}(s, a) := \mathbb{E}[(R(s, a) - \mathbb{E}[R(s, a)])^2]$. Moreover, since the random variable $R(s, a) \in [0, 1]$, we have

$$\sum_{k=1}^K \sum_{h=1}^H \text{Var}(s, a) \leq \sum_{k=1}^K \sum_{h=1}^H r(s, a) \leq \sum_{k=1}^K \sum_{h=1}^H (r(s_h^k, a_h^k) - r_h^k) + K.$$

So we have $\mathbb{P}\left[\left| \sum_{k=1}^K \sum_{h=1}^H (r(s_h^k, a_h^k) - r_h^k) \right| > 2\sqrt{2(\sum_{k=1}^K \sum_{h=1}^H (r(s_h^k, a_h^k) - r_h^k) + K)\iota} + 6\iota \right] \leq 2(\log_2(KH) + 1)\delta$, which implies that with probability $1 - 2(\log_2(KH) + 1)\delta$ we have

$$\sum_{k=1}^K \sum_{h=1}^H (r(s_h^k, a_h^k) - r_h^k) \leq 6\sqrt{K\iota} + 21\iota. \quad (14)$$

For the second term in (13), because $\sum_{h=1}^H r_h^k - \tilde{V}_1^{\pi^k}(s_1^k)$, $k \in [K]$ is a martingale. Hence by Azuma's inequality, we have with probability $1 - \delta$ it holds that

$$\left| \sum_{k=1}^K \left(\sum_{h=1}^H r_h^k - \tilde{V}_1^{\pi^k}(s_1^k) \right) \right| \leq \sqrt{2K\iota}. \quad (15)$$

Combining (14) and (15), we have that $\mathbb{P}[|M_3| > 8\sqrt{K\iota} + 6\iota] \leq 2(\log_2(KH) + 1)\delta$.

Bounding M_2 : For the M_2 term, recall that

$$\begin{aligned} \beta_h^k(s, a) &= O(\min\left\{ \sqrt{\frac{\mathbb{V}(\hat{P}_{s,a}, V_{h+1}^k)\iota}{|\tilde{n}^k(s, a) - C|}}, 1 \right\} + \min\left\{ \sqrt{\frac{\hat{r}_h^k(s, a)\iota}{|\tilde{n}^k(s, a) - C|}}, 1 \right\} + \min\left\{ \frac{\iota}{|\tilde{n}^k(s, a) - C|}, 1 \right\} + \min\left\{ \frac{2C}{|\tilde{n}^k(s, a) - C|}, 1 \right\} \\ &\quad + \min\left\{ \frac{\sqrt{C}\iota}{|\tilde{n}^k(s, a) - C|}, 1 \right\} + \frac{2C}{\tilde{n}^k(s, a) + C} + \sqrt{\frac{\mathbb{V}(P_{s,a}, V_{h+1}^*)}{\hat{n}^k(s, a)}} + \sqrt{\frac{S\mathbb{V}(P_{s,a}, V_{h+1}^k - V_{h+1}^*)\iota}{\hat{n}^k(s, a)}} + \frac{S\iota}{\hat{n}^k(s, a)}. \end{aligned}$$

First note that

$$O(\min\{\frac{\iota}{|\tilde{n}^k(s, a) - C|}, 1\} + \min\{\frac{2C}{|\tilde{n}^k(s, a) - C|}, 1\} + \min\{\frac{\sqrt{C\iota}}{|\tilde{n}^k(s, a) - C|}, 1\} + \frac{2C}{\tilde{n}^k(s, a) + C}) = O(\min\{\frac{C\iota}{|\tilde{n}^k(s, a) - C|}, 1\}).$$

By Lemma 2, we have

$$\mathbb{P}[\hat{P}_{s,a,s'} > \frac{3}{2}P_{s,a,s'} + \frac{4\iota}{3\hat{n}^k(s, a)}] \leq \mathbb{P}[\hat{P}_{s,a,s'} - P_{s,a,s'} > \sqrt{\frac{2P_{s,a,s'}\iota}{\hat{n}^k(s, a)}} + \frac{\iota}{3\hat{n}^k(s, a)}] \leq \delta,$$

which implies that, with probability $1 - 2S^2AH(\log_2(KH) + 1)\delta$, it holds that

$$\begin{aligned} \mathbb{V}(\hat{P}_{s,a,s'}, V_{h+1}^k) &= \sum_{s'} \hat{P}_{s,a,s'} (V_{h+1}^k(s') - \hat{P}_{s,a}^k V_{h+1}^k)^2 \\ &\leq \sum_{s'} \hat{P}_{s,a,s'}^k (V_{h+1}^k(s') - P_{s,a} V_{h+1}^k)^2 \\ &\leq \sum_{s'} (\frac{3}{2}P_{s,a,s'} + \frac{4\iota}{3\hat{n}^k(s, a)}) (V_{h+1}^k(s') - P_{s,a} V_{h+1}^k)^2 \\ &\leq \frac{3}{2}\mathbb{V}(P_{s,a}, V_{h+1}^k) + \frac{4S\iota}{3\hat{n}^k(s, a)}. \end{aligned}$$

Note that $\mathbb{V}(P, X + Y) \leq 2(\mathbb{V}(P, X) + \mathbb{V}(P, Y))$ for any P, X, Y , so we conclude that

$$\begin{aligned} \beta_h^k(s, a) &\leq O(\min\{\sqrt{\frac{\mathbb{V}(P_{s,a}, V_{h+1}^k)\iota}{|\tilde{n}^k(s, a) - C|}}, 1\} + \min\{\sqrt{\frac{\hat{r}_h^k(s, a)\iota}{|\tilde{n}^k(s, a) - C|}}, 1\} + \sqrt{\frac{S\mathbb{V}(P_{s,a}, V_{h+1}^k - V_{h+1}^*)\iota}{\hat{n}^k(s, a)}} \\ &\quad + \min\{\frac{C\iota}{|\tilde{n}^k(s, a) - C|}, 1\} + \frac{S\iota}{\hat{n}^k(s, a)}). \end{aligned}$$

According to the update rule, despite those episodes in which an update is triggered, the number of visit of (s, a) between the i -th update of $\hat{P}_{s,a}$ and the $i + 1$ -th update of $\hat{P}_{s,a}$ do not exceeds 2^{i-1} , i.e., for any (s, a) and any $i \geq 3$, we have

$$\sum_{k=1}^K \sum_{h=1}^H \mathbb{I}[(s_h^k, a_h^k) = (s, a), \tilde{n}^k(s, a) - C = 2^{i-1}] \mathbb{I}[(k, h) \notin \mathcal{X}] \leq 2^{i-1}. \quad (16)$$

Let $l = \max\{i | 2^{i-1} + C \leq KH\}$. We calculate that

$$\begin{aligned} &\sum_{k=1}^K \sum_{h=1}^H \min\{\frac{C\iota}{|\tilde{n}^k(s, a) - C|}, 1\} \mathbb{I}[(k, h) \notin \mathcal{X}] \\ &= \sum_{k=1}^K \sum_{h=1}^H \sum_{s, a} \sum_{i=3}^l (\mathbb{I}[(s_h^k, a_h^k) = (s, a), \tilde{n}^k(s, a) - C = 2^{i-1}] \min\{\frac{C\iota}{|\tilde{n}^k(s, a) - C|}, 1\} \mathbb{I}[(k, h) \notin \mathcal{X}] \\ &\quad + \mathbb{I}[(s_h^k, a_h^k) = (s, a), \tilde{n}^k(s, a) - C < 4]) \\ &\leq \sum_{k=1}^K \sum_{h=1}^H \sum_{s, a} \sum_{i=3}^l (\mathbb{I}[(s_h^k, a_h^k) = (s, a), \tilde{n}^k(s, a) - C = 2^{i-1}] \mathbb{I}[(k, h) \notin \mathcal{X}] \frac{C\iota}{2^{i-1}}) + (C + 4) \\ &\leq \sum_{s, a} C\iota + (C + 4)SA \end{aligned} \quad (17)$$

$$= \tilde{O}(CSA). \quad (18)$$

Here (17) is by using (16). Let $\omega = \{\omega_h^k \geq 0 | 1 \leq h \leq H, 1 \leq k \leq K\}$ be a group of non-negative weights such that

$w_h^k \in [0, 1]$ for any k, h and $w_h^k = 0$ if $(k, h) \in \mathcal{X}$. We prove the following useful inequality:

$$\begin{aligned}
 & \sum_{k=1}^K \sum_{h=1}^H \min\left\{\sqrt{\frac{\omega_h^k}{\tilde{n}^k(s, a) - C}}, 1\right\} \\
 \leq & \sum_{k=1}^K \sum_{h=1}^H \sum_{s, a} \sum_{i=3}^l \mathbb{I}[(s_h^k, a_h^k) = (s, a), \tilde{n}^k(s, a) = 2^{i-1}] \sqrt{\frac{\omega_h^k}{2^{i-1}}} + CSA + 8SA(\log_2(KH) + 4)\iota \\
 = & \sum_{s, a} \sum_{i=3}^l \frac{1}{\sqrt{2^{i-1}}} \sum_{k=1}^K \sum_{h=1}^H \mathbb{I}[(s_h^k, a_h^k) = (s, a), \tilde{n}^k(s, a) = 2^{i-1}] \sqrt{\omega_h^k} + CSA + 8SA(\log_2(KH) + 4)\iota \\
 \leq & \sum_{s, a} \sum_{i=3}^l \sqrt{\frac{\sum_{k=1}^K \sum_{h=1}^H \mathbb{I}[(s_h^k, a_h^k) = (s, a), \tilde{n}^k(s, a) - C = 2^{i-1}]\iota}{2^{i-1}}} \sqrt{\sum_{k=1}^K \sum_{h=1}^H \mathbb{I}[(s_h^k, a_h^k) = (s, a), \tilde{n}^k(s, a) - C = 2^{i-1}]\omega_h^k} \\
 & + CSA + 8SA(\log_2(KH) + 4)\iota
 \end{aligned} \tag{19}$$

$$\leq \sqrt{SA\iota \sum_{k=1}^K \sum_{h=1}^H \omega_h^k} + CSA + 8SA(\log_2(KH) + 4)\iota. \tag{20}$$

Here (19) is by Cauchy-Schwarz inequality and (20) is due to (16). By using the same technique (or see (29) in (Zhang et al., 2020a)), we can prove that

$$\sum_{k=1}^K \sum_{h=1}^H \sqrt{\frac{\omega_h^k}{\hat{n}^k(s_h^k, a_h^k)}} \leq \sqrt{SA\iota \sum_{k=1}^K \sum_{h=1}^H \omega_h^k} + 8SA(\log_2(KH) + 4). \tag{21}$$

Let $I(k, h)$ be the shorthand of $\mathbb{I}[(k, h) \notin \mathcal{X}]$. By plugging in respectively $\omega_h^k = I(k, h)\hat{r}_h^k(s_h^k, a_h^k)$, $I(k, h)\mathbb{V}(P_{s_h^k, a_h^k}, V_{h+1}^k)$ into (20) and $\omega_h^k = I(k, h)\mathbb{V}(V_{h+1}^k - V_{h+1}^*)$ into (21), we obtain that

$$\begin{aligned}
 M_2 &= \sum_{k=1}^K \sum_{h=1}^H \tilde{\beta}_h^k(s_h^k, a_h^k) \\
 &= \sum_{k=1}^K \sum_{h=1}^H I(k, h)\beta_h^k(s_h^k, a_h^k) \\
 &\leq O\left(\sqrt{SA\iota \sum_{k=1}^K \sum_{h=1}^H \mathbb{V}(P_{s_h^k, a_h^k}, V_{h+1}^k)I(k, h)} + \sqrt{S^2A\iota \sum_{k=1}^K \sum_{h=1}^H \mathbb{V}(P_{s_h^k, a_h^k}, V_{h+1}^k - V_{h+1}^*)I(k, h)}\right) \\
 &\quad + O\left(\sqrt{SA\iota \sum_{k=1}^K \sum_{h=1}^H \hat{r}_h^k(s_h^k, a_h^k)I(k, h)} + \tilde{O}(CSA + S^2A)\right)
 \end{aligned} \tag{22}$$

We define $M_5 := \sum_{k=1}^K \sum_{h=1}^H \mathbb{V}(P_{s_h^k, a_h^k}, V_{h+1}^k - V_{h+1}^*)I(k, h + 1)$. We state the following Lemmas in (Zhang et al., 2020a) to complete the proof of Theorem 1, which we omit the detailed proof.

Lemma 6 (Lemma 5 in (Zhang et al., 2020a)). $\sum_{k=1}^K \sum_{h=1}^H \hat{r}_h^k(s_h^k, a_h^k)I(k, h) \leq 2 \sum_{k=1}^K \sum_{h=1}^H r_h^k + 4SA \leq 2K + 4SA$.

Lemma 7 (Lemma 6 in (Zhang et al., 2020a)). *With probability $1 - 2(\log_2(KH) + 1)\log_2(KH)\delta$, it holds that*

$$M_4 \leq 2M_2 + 2|\mathcal{K}^C| + 2K + \max\{46\iota, 8\sqrt{(M_2 + |\mathcal{K}^C| + K)\iota} + 6\iota\}.$$

Lemma 8 (Lemma 7 in (Zhang et al., 2020a)). *With probability $1 - 2(\log_2(KH) + 1)\log_2(KH)\delta$, it holds that*

$$M_5 \leq 2 \max\{M_2, 1\} + 2|\mathcal{K}^C| + \max\{46\iota, 8\sqrt{(M_2 + |\mathcal{K}^C|)\iota} + 6\iota\}.$$

Combining the above lemmas with (22), we have that with probability $1 - (6S^2AH(\log_2(KH) + 1) + 6(\log_2(KH) + 1)\log_2(H))\delta$,

$$M_2 \leq O(\sqrt{SAl\ell(M_4 + |\mathcal{K}^C|)} + \sqrt{S^2Al\ell(M_5 + |\mathcal{K}^C|)} + \sqrt{SAl\ell K} + CSA + S^2A\ell \log_2(KH)),$$

$$M_4 \leq 2M_2 + 2|\mathcal{K}^C| + \max\{46\ell, 8\sqrt{(M_2 + 2K)\ell} + 6\ell\},$$

$$M_5 \leq 2\max\{M_2, 1\} + 2|\mathcal{K}^C| + \max\{46\ell, \sqrt{M_2\ell} + 6\ell\}.$$

Which implies that

$$M_2 \leq O(\sqrt{SAK\ell} + CSA + \sqrt{S^2Al\sqrt{M_2\ell^{3/2}} + \sqrt{SAl\ell K} + S^2A\ell \log_2(KH)}) \quad (23)$$

$$\leq O(\sqrt{SAK\ell} + CSA + S^2A\ell \log_2(KH)). \quad (24)$$

Recalling (12) and (23), we conclude that, with probability $1 - (10S^2AH(\log_2(KH) + 2) + 6(\log_2(KH) + 1)\log_2(KH) + 1)\delta$,

$$\begin{aligned} \text{Regret}(K) &\leq M_1 + M_2 + M_3 \\ &\leq O(\sqrt{SAK\ell} + S^2A\ell \log_2(KH) + CSA + \sqrt{K\ell}) \\ &= O(\sqrt{SAK\ell} + S^2A\ell \log_2(KH) + CSA). \end{aligned}$$

Hence by rescaling δ , we finish the proof of Theorem 1.

B. Proof of Theorem 3

Proof. First, since we run Algorithm 4 to learn a $\frac{\mu_{\min}}{2}$ -policy cover, from Theorem 4, this requires

$$H(N'_g + N'_\phi + N'_p) = \tilde{\Omega}\left(\frac{M^4 A^4 H^2 \log |\mathcal{G}|}{\mu_{\min}^4 \gamma^2} + \frac{HMA}{\mu_{\min}} + \frac{M^2 AH^2}{\mu_{\min}^3}\right)$$

trajectories.

Next, we prove that in phase 2, Algorithm 2 learns a decoding function with ϵ_f decoding error, where ϵ_f is to be defined later. Formally, we prove that the following condition holds with high probability:

Condition 1 (Bijection between learned and true states). *There exists $\epsilon_f < \frac{1}{2}$ such that there is a bijective mapping $\alpha_h : \hat{\mathcal{S}}_h \rightarrow \mathcal{S}_h$ for which*

$$\mathbb{P}_{x \sim q(\cdot | \alpha_h(s))} [\hat{f}_h(x) = \hat{s}] \geq 1 - \epsilon_f.$$

In other words, this condition states that every estimated latent state \hat{s} roughly corresponds to a true latent state $\alpha_h(\hat{s})$, when we use the decoding function \hat{f}_h . This is because all but an ϵ_f fraction of contexts drawn from $\alpha_h(\hat{s})$ are decoded to their true latent state, and for each latent state s , there is a distinct estimated state $\alpha_h^{-1}(s)$ as the map α_h is a bijection. For simplicity, we define $\mathbf{p}(s, a) \in \mathbb{R}^M$ to be the forward transition distribution over \mathcal{S}_h for $s \in \mathcal{S}_{h-1}$ and $a \in \mathcal{A}$. We abuse notation to similarly use $\mathbf{p}(\hat{s}, a) \in \mathbb{R}^M$ to be the vector $\{\mathbb{P}(s | \hat{s}, a)\}_{s \in \mathcal{S}_h}$ of conditional probabilities \mathcal{S}_h for $\hat{s} \in \hat{\mathcal{S}}_{h-1}$ and $a \in \mathcal{A}$. Note that unlike $s \in \mathcal{S}_{h-1}$, $\hat{s} \in \hat{\mathcal{S}}_{h-1}$ is not a Markovian state and hence the conditional probability vector $\mathbf{p}(\hat{s}, a)$ depends on the specific distribution over $\hat{\mathcal{S}}_{h-1} \times \mathcal{A}$. In the following we will use $\mathbf{p}^\nu(\hat{s}, a)$ to emphasize this dependency where ν is the distribution, where ν is a distribution over $\hat{\mathcal{S}}_{h-1} \times \mathcal{A}$.

In the proof, we often compare two vectors indexed by \mathcal{S}_h and $\hat{\mathcal{S}}_h$. We will assume the order of the indices of these two vectors are matched according to α_h .

Establishing Condition 1. In order to establish the condition, we need to show that our decoding function \hat{f}_h predicts the underlying latent state correctly almost always. We do this in two steps. Since the functions \hat{f}_h are derived based on \hat{g}_h and $\hat{\phi}_h$, we analyze the properties of these two objects in the following two lemmas in (Du et al., 2019). In order to state the first lemma, we need some additional notation. Note that η_h and \hat{f}_{h-1} induce a distribution over $\mathcal{S}_{h-1} \times \hat{\mathcal{S}}_{h-1} \times \mathcal{A} \times \mathcal{S}_h$. We denote

this distribution as ν_h . With this distribution, we define the conditional backward probability $\hat{\mathbf{b}}_{\nu_h} : \mathcal{S}_h \rightarrow \Delta(\widehat{\mathcal{S}}_{h-1} \times \mathcal{A})$ as

$$\hat{\mathbf{b}}_{\nu_h}(\hat{s}, a | s'_1) = \frac{p_{h-1}^{\nu_h}(s'_1 | \hat{s}, a) \mathbb{P}^{\nu_h}(\hat{s}, a)}{\sum_{\hat{s}_1, a_1} p_{h-1}^{\nu_h}(s'_1 | \hat{s}_1, a_1) \mathbb{P}^{\nu_h}(\hat{s}, a_1)} \quad (25)$$

Recall that $\mathbf{p}_{h-1}^{\nu_h}$ above refers to the distribution over s'_1 according the transition dynamics, when \hat{s}, a are induced by ν_h . With this notation, we have the following lemma.

Lemma 9 (Lemma G.2 in (Du et al., 2019)). *Assume $\epsilon_f \leq \frac{\mu_{\min}^3 \gamma}{100M^4 A^3}$. Then the distributions $\hat{\mathbf{b}}_{\nu_h}(\hat{s}, a | s')$ are well separated for any pair $s'_1, s'_2 \in \mathcal{S}_h$:*

$$\left\| \hat{\mathbf{b}}_{\nu_h}(s'_1) - \hat{\mathbf{b}}_{\nu_h}(s'_2) \right\|_1 \geq \frac{\mu_{\min} \gamma}{3MA} \quad (26)$$

Furthermore, if $N_g = \Omega\left(\frac{M^3 A^3}{\epsilon_f \mu_{\min}^3 \gamma^2} \log\left(\frac{|G|H}{\delta}\right)\right)$, with probability at least $1 - \delta/H$, for every $s' \in \mathcal{S}_h$, \hat{g}_h satisfies

$$\mathbb{P}_{x' \sim q(\cdot | s')} \left[\left\| \hat{g}_h(x') - \hat{\mathbf{b}}_{\nu_h}(s') \right\|_1 \geq \frac{\gamma \mu_{\min}}{100MA} \right] \leq \epsilon_f \quad (27)$$

Proof. See (Du et al., 2019) for details. \square

The first part of Lemma 9 tell us that the latent states at level h are well separated if we embed them using $\phi(s') = \hat{\mathbf{b}}_{\nu_h}(s')$ as the state embedding. The second part guarantees that our regression procedure estimates this representation accurately. Together, these assertions imply that any two contexts from the same latent state (up to an ϵ_f fraction) are close to each other, while contexts from two different latent states are well-separated. Formally, with probability at least $1 - \frac{\delta}{H}$ over the N_g training data:

1. For any $s' \in \mathcal{S}_h$ and $x'_1, x'_2 \sim q(\cdot | s')$, we have with probability at least $1 - 2\epsilon_f$ over the emission process

$$\left\| \hat{g}_h(x'_1) - \hat{g}_h(x'_2) \right\|_1 \leq \frac{\mu_{\min} \gamma}{50MA}$$

2. For any $s'_1, s'_2 \in \mathcal{S}_h$ such that $s'_1 \neq s'_2, x'_1 \sim q(\cdot | s'_1)$ and $x'_2 \sim q(\cdot | s'_2)$, we have with probability at least $1 - 2\epsilon_f$ over the emission process

$$\left\| \hat{g}_h(x'_1) - \hat{g}_h(x'_2) \right\|_1 \geq \frac{\mu_{\min} \gamma}{4MA}$$

In other words, the mapping of contexts, as performed through the functions \hat{g}_h should be easy to cluster with each cluster roughly corresponding to a true latent state. Our next lemma guarantees that with enough samples for clustering, this is indeed the case.

Lemma 10 (Lemma G.3 in (Du et al., 2019)). *If $N_\phi = \Theta\left(\frac{MA}{\mu_{\min}} \log\left(\frac{MH}{\delta}\right)\right)$ and $\epsilon_f \leq \frac{\delta}{100HN_\phi}$ we have with probability at least $1 - \frac{\delta}{H}$, (1) for every $s' \in \mathcal{S}_h$, there exists at least one point $\mathbf{z} \in \mathcal{Z}$ such that $\mathbf{z} = \hat{g}_h(x')$ with $x' \sim q(\cdot | s')$ and $\left\| \hat{g}_h(x') - \hat{\mathbf{b}}_{\nu_h}(s') \right\|_1 \leq \frac{\mu_{\min} \gamma}{100MA}$ and (2) for every $\mathbf{z} = \hat{g}_h(x') \in \mathcal{Z}$ with $x' \sim q(\cdot | s')$, $\left\| \hat{g}_h(x') - \hat{\mathbf{b}}_{\nu_h}(s') \right\|_1 \leq \frac{\mu_{\min} \gamma}{100MA}$*

Proof. See (Du et al., 2019) for details. \square

Based on Lemma 9 and 10, we can establish that Condition 1 holds with high probability. Note that Condition 1 consists of two parts. The first part states that there exists a bijective map $\alpha_h : \widehat{\mathcal{S}}_h \rightarrow \mathcal{S}_h$. The second part states that the decoding error is small. To prove the first part, we explicitly construct the map α_h and show it is bijective. We define $\alpha_h : \widehat{\mathcal{S}}_h \rightarrow \mathcal{S}_h$ as

$$\alpha_h(\hat{s}') = \operatorname{argmin}_{s \in \mathcal{S}_h} \left\| \phi(s') - \hat{\phi}(\hat{s}') \right\|_1 \quad (28)$$

First observe that for any $\hat{s}' \in \widehat{\mathcal{S}}_h$, by the second conclusion of Lemma 10 we know there exists $s' \in \mathcal{S}_h$ such that

$$\left\| \hat{\phi}(\hat{s}') - \phi(s') \right\| \leq \frac{\gamma \mu_{\min}}{100MA}$$

This also implies for any $s'' \neq s'$

$$\left\| \widehat{\phi}(\hat{s}') - \phi(s'') \right\| \geq \|\phi(s'') - \phi(s')\| - \left\| \widehat{\phi}(\hat{s}') - \phi(s') \right\| \geq \frac{\gamma\mu_{\min}}{4MA}$$

Therefore we know $\alpha_h(\hat{s}') = s'$, i.e., α_h always maps the learned state to the correct original state.

We now prove α_h is injective, i.e., $\alpha(\hat{s}') \neq \alpha_h(\hat{s}'')$ for $\hat{s}' \neq \hat{s}'' \in \widehat{\mathcal{S}}_h$. Suppose there are $\hat{s}', \hat{s}'' \in \widehat{\mathcal{S}}_h$ such that $\alpha_h(\hat{s}') = \alpha_h(\hat{s}'') = s'$ for some $s' \in \mathcal{S}_h$. Then using the second conclusion of Lemma 10, we know

$$\left\| \widehat{\phi}(\hat{s}') - \widehat{\phi}(\hat{s}'') \right\|_1 \leq \left\| \widehat{\phi}(\hat{s}') - \phi(s') \right\|_1 + \left\| \phi(s') - \widehat{\phi}(\hat{s}'') \right\|_1 \leq \frac{\gamma\mu_{\min}}{50MA} \quad (29)$$

However, we know by Algorithm 2, every $\hat{s}' \neq \hat{s}'' \in \widehat{\mathcal{S}}_h$ must satisfy

$$\left\| \widehat{\phi}(\hat{s}') - \widehat{\phi}(\hat{s}'') \right\|_1 > \tau = \frac{\gamma\mu_{\min}}{30MA}$$

This leads to a contradiction and thus α_h is injective.

Next we prove α_h is surjective, i.e., for every $s' \in \mathcal{S}_h$, there exists $\hat{s}' \in \widehat{\mathcal{S}}_h$ such that $\alpha_h(\hat{s}') = s'$. The first conclusion in Lemma G. 3 guarantees that for each latent state $s' \in \mathcal{S}_h$, there exists $\mathbf{z} = \hat{g}(x') \in \mathcal{Z}$ with $x' \sim q(\cdot | s')$. The second conclusion of Lemma 10 guarantees that

$$\|\mathbf{z} - \phi(s')\|_1 \leq \frac{\gamma\mu_{\min}}{100MA}$$

Now we first assert that all points in a cluster are emitted from the same latent state by combining Equation (10), the second part of Lemma G. 3 and our setting of τ . Now the second part of Lemma G. 3 implies that there exists $\hat{s}' \in \widehat{\mathcal{S}}_h$ such that $\left\| \mathbf{z} - \widehat{\phi}(\hat{s}') \right\|_1 \leq \frac{\mu_{\min}\gamma}{50MA}$, since \mathbf{z} and $\widehat{\phi}(\hat{s}')$ correspond to \hat{g} evaluated on two different contexts in the same cluster. Therefore we have

$$\left\| \phi(s') - \widehat{\phi}(\hat{s}') \right\|_1 \leq \|\phi(s') - \mathbf{z}\|_1 + \left\| \mathbf{z} - \widehat{\phi}(\hat{s}') \right\|_1 \leq \frac{\mu_{\min}\gamma}{30MA}$$

Now we can show that $\alpha_h(\hat{s}') = s'$. To do this, we show that $\widehat{\phi}(\hat{s}')$ is closer to $\phi(s')$ than the embedding of any state in \mathcal{S}_h . Using the second conclusion of Lemma G. 3 and Equation 10 we know for any $s'' \neq s'$

$$\left\| \widehat{\phi}(\hat{s}') - \phi(s'') \right\|_1 \geq \|\phi(s') - \phi(s'')\|_1 - \left\| \widehat{\phi}(\hat{s}') - \phi(s') \right\|_1 \geq \frac{\gamma\mu_{\min}}{4MA}$$

We know $s' = \operatorname{argmin}_{s_1 \in \mathcal{S}_h} \left\| \widehat{\phi}(s_1) - \widehat{\phi}(\hat{s}') \right\|_1$. Therefore, by the definition of α_h we know $\alpha_h(\hat{s}') = s'$. Now we have finished the proof of the first part of Condition G.1.

For the second part of Condition G. 1, note for any $s' \in \mathcal{S}_h$ and $x' \sim q(\cdot | s')$, by Lemma G.2, we know with probability at least $1 - \epsilon_f$ over the emission process we have

$$\|\hat{g}_h(x') - \phi(s')\|_1 \leq \frac{\gamma\mu_{\min}}{100MA} \quad (30)$$

For $\hat{s}' = \alpha_h^{-1}(s')$, we have

$$\left\| \hat{g}_h(x') - \widehat{\phi}(\hat{s}') \right\|_1 \leq \|\hat{g}_h(x') - \phi(s')\|_1 + \left\| \phi(s') - \widehat{\phi}(\hat{s}') \right\|_1 \leq \frac{\gamma\mu_{\min}}{50MA}$$

On the other hand, for $\hat{s}'' \in \widehat{\mathcal{S}}_h$ with $\hat{s}'' \neq \alpha_h^{-1}(s')$, we have

$$\left\| \hat{g}_h(x') - \widehat{\phi}(\hat{s}'') \right\|_1 \geq -\|\hat{g}_h(x') - \phi(s')\|_1 + \|\phi(s') - \phi(\alpha_h(\hat{s}''))\|_1 - \left\| \phi(\alpha_h(\hat{s}'')) - \widehat{\phi}(\hat{s}'') \right\|_1 \geq \frac{\gamma\mu_{\min}}{4MA}$$

Therefore we have with probability at least $1 - \epsilon_f$

$$\hat{f}_h(x') = \operatorname{argmin}_{\hat{s}' \in \widehat{\mathcal{S}}_h} \left\| \widehat{\phi}(\hat{s}') - \hat{g}_h(x') \right\|_1 = \alpha_h^{-1}(s')$$

which is equivalent to the second part of Condition G.1.

Combine the analysis above, we have the following lemma.

Lemma 11. Assume $N_\phi = \Theta\left(\frac{MA}{\mu_{\min}} \log\left(\frac{MH}{\delta}\right)\right)$. For any $\epsilon_f < \min\left\{\frac{\mu_{\min}^3 \gamma}{100M^4 A^3}, \frac{\delta}{100HN_\phi}\right\}$, set $N_g = \Omega\left(\frac{M^3 A^3}{\epsilon_f \mu_{\min}^3 \gamma^2} \log\left(\frac{|\mathcal{G}|H}{\delta}\right)\right)$, then Algorithm 2 learns a decoding function \hat{f} such that there is a bijective mapping $\alpha_h : \hat{\mathcal{S}}_h \rightarrow \mathcal{S}_h$ for which

$$\mathbb{P}_{x \sim q(\cdot | \alpha_h(\hat{s}))} \left[\hat{f}_h(x) = \hat{s} \right] \geq 1 - \epsilon_f.$$

Finally, we prove that with an ϵ_f -error decoding function, with high probability there will be at most $2\epsilon_f T' + \sqrt{2T' \ln \frac{\delta}{2}}$ states in phase 3 where the agent makes a mistake. Here T' is the total number of steps in phase 3. This allows us to call CR-MVP with $C = 2\epsilon_f T' + \sqrt{2T' \ln \frac{\delta}{2}}$.

Lemma 12. Assume that the decoding function has ϵ_f decoding error, that is, $\mathbb{P}(f(s) \neq s) \leq \epsilon_f$. Define $\tilde{N}(s, a, s')$ as the counter without error, and $\hat{N}(s, a, s')$ as the counter with error. we have

$$\mathbb{P}\left(\sum_{s, a, s'} \left| \hat{N}(s, a, s') - \tilde{N}(s, a, s') \right| \leq 2\epsilon_f T' + \sqrt{2T' \ln \frac{2}{\delta}}\right) \leq \frac{\delta}{2} \quad (31)$$

Proof. This is a simple corollary of Chernoff-Hoeffding inequality. \square

Combine the three parts of the algorithm together, we can bound the total regret as

$$\begin{aligned} & H(N'_g + N'_\phi + N'_p) + H(N_g + N_\phi) + \tilde{O}\left(\sqrt{HMAK} + H^2 M^2 A + HMA \left(\epsilon_f T' + \sqrt{T' \ln \frac{\delta}{2}}\right)\right) \\ &= \tilde{O}\left(\frac{M^4 A^4 H^2 \log |\mathcal{G}|}{\mu_{\min}^4 \gamma^2} + \frac{HMA}{\mu_{\min}} + \frac{M^2 A H^2}{\mu_{\min}^3} + \frac{HMA}{\mu_{\min}} \log\left(\frac{MH}{\delta}\right) + \right. \\ & \quad \left. \frac{HM^3 A^3}{\epsilon_f \mu_{\min}^3 \gamma^2} \log\left(\frac{|\mathcal{G}|H}{\delta}\right) + \sqrt{HMAK} + H^2 M^2 A + HMA \left(\epsilon_f T' + \sqrt{T' \ln \frac{\delta}{2}}\right)\right) \end{aligned}$$

Finally, by setting $\epsilon_f = \sqrt{\frac{M^2 A^2}{\mu_{\min}^3 \gamma^2 H K} \log\left(\frac{|\mathcal{G}|H}{\delta}\right)}$, the regret is upper bounded by $\tilde{O}\left(\frac{H^{3/2} M^2 A^2 \sqrt{K}}{\mu_{\min}^{3/2} \gamma} + \text{poly}(H, M, A, \mu_{\min}^{-1}, \gamma^{-1})\right)$, which finishes the proof. \square

C. Proof of Lower Bounds

C.1. Proof of Theorem 2

Proof. First we prove the lower bound in the multi-armed bandit case, *i.e.* when there is only one state and A actions (called A arms in previous literature). Let \mathcal{E}_1 and \mathcal{E}_2 be two bandit instances. In \mathcal{E}_1 , the reward of arm 1 is $Ber(\frac{1}{2})$, and the reward of the other $A - 1$ arms is $Ber(0)$. There is no adversarial corruption in \mathcal{E}_1 . We denote $\text{Regret}(\mathcal{A}, \mathcal{E}_1, T)$ as the regret \mathcal{A} incurs under environment \mathcal{E}_1 after T steps.

Let $T_a(n)$ denote the number of times arm a is chosen in the first n steps. Under environment \mathcal{E}_1 . The agent incurs regret only when it makes a suboptimal action, which can only happen when it pulls arm $a \neq 1$. So the regret under \mathcal{E}_1 can be represented as

$$\text{Regret}(\mathcal{A}, \mathcal{E}_1, T) = \frac{1}{2} \sum_{i=2}^A \mathbb{E}[T_i(T)] \quad (32)$$

Now we consider two cases: **Case 1:** $\forall i \geq 2, \mathbb{P}_{\mathcal{A}, \mathcal{E}_1} (T_i(T) > \frac{C}{2}) \geq \frac{1}{2}$. Then by Lemma 13, we have

$$\text{Regret}(\mathcal{A}, \mathcal{E}_1, T) = \frac{1}{2} \sum_{i=2}^A \mathbb{E}[T_{s_1, i}(T)] \quad (33)$$

$$\geq \frac{C(A-1)}{16} \mathbb{P} \left(\sum_{i=2}^A T_{s_1, i}(T) \geq \frac{C(A-1)}{8} \right) \quad (34)$$

$$\geq \frac{C(A-1)}{48} \quad (35)$$

Case 2: $\exists i \geq 2$, such that $\mathbb{P}_{\mathcal{A}, \mathcal{E}_1} (T_i(T) \leq \frac{C}{2}) \geq \frac{1}{2}$. In this case we define another environment \mathcal{E}_2 identical to \mathcal{E}_1 except that the reward of arm i is $Ber(1)$. And for the first $\frac{C}{2}$ times the agent pulls arm i , the adversary will corrupt the reward to be $Ber(0)$. In such environment, the corruption level is C and the agent incurs regret at least $\frac{1}{2}$ each time the agent takes action $a \neq i$.

Now since the rewards \mathcal{E}_1 and \mathcal{E}_2 generate are the same before $T_i(T)$ exceeds $\frac{C}{2}$, the agent will pull arms identical to that in \mathcal{E}_1 when it runs algorithm \mathcal{A} , until it pulls arm i for more than $\frac{C}{2}$ times. So $\mathbb{P}_{\mathcal{A}, \mathcal{E}_2} (T_i(T) \leq \frac{C}{2}) \geq \frac{1}{2}$. Then the regret under \mathcal{E}_2 is

$$\text{Regret}(\mathcal{A}, \mathcal{E}_2, T) \geq \frac{1}{2} \sum_{j \neq i} \mathbb{E}[T_j(T)] \quad (36)$$

$$\geq \frac{1}{2} \left(T - \frac{C}{2} \right) \mathbb{P}_{\mathcal{A}, \mathcal{E}_1} \left(\sum_{j \neq i} T_j(T) \geq T - \frac{C}{2} \right) \quad (37)$$

$$= \frac{1}{2} \left(T - \frac{C}{2} \right) \mathbb{P}_{\mathcal{A}, \mathcal{E}_1} \left(T_i(T) \leq \frac{C}{2} \right) \quad (38)$$

$$\geq \frac{1}{4} \left(T - \frac{C}{2} \right) \quad (39)$$

Since $T \geq 2AC$, in case 2 the regret is at least $O(AC)$, which completes the proof.

Now we consider the MDP case. Without loss of generality, assume that $S = A^H$, where H is a positive integer. We construct the following MDP with H horizon: In each episode the MDP starts from a fixed state s_1 . And for any $h \in [H]$, in step h , each state s_h can be represented by a sequence of actions with length $h - 1$. And when the environment is in state $s_h = (a_1, \dots, a_{h-1})$ and the agent takes action a_h , it will transit to state $s_{h+1} = (a_1, \dots, a_{h-1}, a_h)$. In other words, the transition of this MDP can be represented by an A -nary tree, which is shown in Figure 1. The reward for the first $H - 1$ episodes are all 0, so the agent receives a reward signal only when it reaches a leaf node in step H .

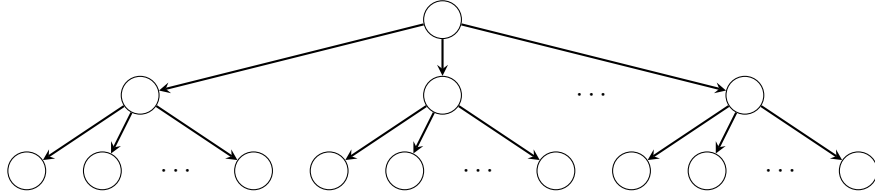


Figure 1. An A -nary tree with three layers.

The MDP is equivalent to a multi-armed bandit problem with SA arms. So we can use the result on multi-armed bandits and obtain a $\Omega(CSA)$ lower bound on the regret. \square

Lemma 13. Assume X_i are nonnegative random variables and $\mathbb{P}(X_i \geq C) \geq \frac{1}{2}$ for all $i \in [n]$. Then $\mathbb{P}(\sum_{i=1}^n X_i \geq \frac{nC}{4}) \geq \frac{1}{3}$.

Proof. Define $Y_i = \min\{C, X_i\}$, then $\mathbb{P}(\sum_{i=1}^n X_i \geq \frac{nC}{4}) \geq \mathbb{P}(\sum_{i=1}^n Y_i \geq \frac{nC}{4})$. Denote $z = \mathbb{P}(\sum_{i=1}^n X_i \geq \frac{nC}{4})$, then

$$\frac{nC}{2} \leq \mathbb{E}(\sum_{i=1}^n Y_i) \leq \mathbb{P}(\sum_{i=1}^n Y_i \geq \frac{nC}{4})nC + \mathbb{P}(\sum_{i=1}^n Y_i < \frac{nC}{4})\frac{nC}{4} \quad (40)$$

This yields

$$\frac{nC}{2} \leq znC + (1-z)\frac{nC}{4} \Rightarrow z \geq \frac{1}{3} \quad (41)$$

which finishes the proof. \square

C.2. Proof of Proposition 1

Proof. As in the proof of Theorem 3, we define \mathcal{E}_1 as a A -arm bandit instance without corruptions, and the reward vector of the A arms is $(\frac{1}{2}, 0, \dots, 0)$. Consider running algorithm \mathcal{A} on \mathcal{E}_1 . Suppose $i = \operatorname{argmin}_{j>1} \mathbb{E}_1[T_i]$, where T_i is the number of times arm i is pulled and \mathbb{E}_1 denotes taking expectation in \mathcal{E}_1 . We define environment \mathcal{E}_2 as follows: let the corruption level C be $\lceil 2\mathbb{E}[T_i] \rceil + 1$, and the reward vector is $(\frac{1}{2}, 0, \dots, 0, 1, 0, \dots, 0)$, where the i -th element is 1. Moreover, the adversary create a false 0 reward of the i -th arm whenever the i -th arm is pulled less than C times.

First we consider the regret incurred in \mathcal{E}_1 . Note that by the definition of T_i , the algorithm will incur at least $O(\mathbb{E}[T_i])$ regret in \mathcal{E}_1 , since arm i is suboptimal in \mathcal{E}_1 and is pulled for at least $O(\mathbb{E}[T_i])$ times. In \mathcal{E}_1 the corruption level is 0, so the regret is upper bounded by $O(\sqrt{K})$. This implies that

$$O(\sqrt{K}) \geq \operatorname{Regret}(\mathcal{A}, \mathcal{E}_1, K) \geq \Omega(\mathbb{E}[T_i]) = \Omega(C). \quad (42)$$

On the other hand, consider running the algorithm \mathcal{A} on \mathcal{E}_2 . Since the agent will receive a reward of 0 for the first C times it pulls arm i , the agent will act just as it were in \mathcal{E}_1 , unless it pulls arm i for more than C times. Note that the probability that $T_i \leq C$ is at least $\frac{1}{2}$ by the definition of T_i and C , which means that with high probability the algorithm will pull arm i for more than C times. As a result, the algorithm will incur regret at least $\frac{1}{2}(K - C) = \Omega(K)$ with probability $\frac{1}{2}$, since arm i is the optimal arm in \mathcal{E}_2 . This implies that

$$O(\sqrt{K} + K^\alpha C^\beta) \geq \operatorname{Regret}(\mathcal{A}, \mathcal{E}_2, K) \geq \Omega(K). \quad (43)$$

Hence by combining (42) and (43) we have,

$$K^{\frac{1}{2}} \geq C \geq K^{\frac{1-\alpha}{\beta}},$$

which is equal to $\alpha + \frac{\beta}{2} \geq 1$. \square

D. Detailed description of PCID algorithm

Here we present the original PCID algorithm in (Du et al., 2019) for convenience.

The following theorem shows that with a polynomial sample complexity, PCID returns an ϵ -policy cover of all latent states with high probability.

Theorem 4 (Theorem 4.1 in (Du et al., 2019)). *Fix any $\epsilon = O\left(\frac{\mu_{\min}^3 \gamma}{M^4 A^3 H}\right)$ and a failure probability $\delta > 0$. Set $N_g = \tilde{\Omega}\left(\frac{M^4 A^4 H \log |\mathcal{G}|}{\epsilon \mu_{\min}^3 \gamma^2}\right)$, $N_\phi = \tilde{\Theta}\left(\frac{MA}{\mu_{\min}}\right)$, $N_p = \tilde{\Omega}\left(\frac{M^2 AH^2}{\mu_{\min} \epsilon^2}\right)$, $\tau = \frac{\gamma}{30MA}$. Then with probability at least $1 - \delta$, Algorithm 4 returns an ϵ -policy cover of \mathcal{S} , with size at most MH .*

Algorithm 4 PCID

Input: N_g : sample size for learning context embeddings, N_ϕ : sample size for learning state embeddings, N_p : sample size for estimating transition probabilities, $\tau > 0$: a clustering threshold for learning latent states

Output: policy cover $\Pi = \Pi_1 \cup \dots \cup \Pi_{H+1}$

Let $\widehat{\mathcal{S}}_1 = \{s_1\}$.

Let $\hat{f}_1(x) = s_1$ for all $x \in \mathcal{X}$

Let $\Pi_1 = \{\pi_0\}$ where π_0 is the trivial 0-step policy. Initialize \hat{p} to an empty mapping.

for $h = 2, \dots, H + 1$ **do**

Let $\eta_h = U(\Pi_{h-1}) \odot U(\mathcal{A})$

Execute η_h for N_g times. $D_g = \{\hat{s}_{h-1}^i, a_{h-1}^i, x_h^i\}_{i=1}^{N_g}$ for $\hat{s}_{h-1} = \hat{f}_{h-1}(x_{h-1})$

Learn \hat{g}_h by calling ERM oracle on input D_g :

$$\hat{g}_h = \operatorname{argmin}_{g \in \mathcal{G}} \sum_{(\hat{s}, a, x') \in D_g} \|g(x') - \mathbf{e}_{(\hat{s}, a)}\|^2$$

Execute η_h for N_ϕ times. $\mathcal{Z} = \{\hat{\mathbf{z}}_i = \hat{\mathbf{g}}_h(x_h^i)\}_{i=1}^{N_\phi}$

Learn $\widehat{\mathcal{S}}_h$ and the state embedding map $\hat{\phi}_h : \widehat{\mathcal{S}}_h \rightarrow \mathcal{Z}$ by clustering \mathcal{Z} with threshold τ (see Algorithm 3).

Define $\hat{f}_h(x') = \operatorname{argmin}_{\hat{s} \in \widehat{\mathcal{S}}_h} \|\hat{\phi}(\hat{s}) - \hat{\mathbf{g}}_h(x')\|_1$

Execute η_h for N_p times. $D_p = \{\hat{s}_{h-1}^i, a_{h-1}^i, \hat{s}_h^i\}_{i=1}^{N_p}$ for $\hat{s}_{h-1} = \hat{f}_{h-1}(x_{h-1})$, $\hat{s}_h = \hat{f}_h(x_h)$

Define $\hat{p}(\hat{s}_h | \hat{s}_{h-1}, a_{h-1})$ equal to empirical conditional probabilities in D_p .

end for

for $\hat{s}' \in \widehat{\mathcal{S}}_h$ **do**

Run Algorithm 5 with inputs \hat{p} and \hat{s}' to obtain $(h-1)$ -step policy $\psi_{\hat{s}'} : \widehat{\mathcal{S}}_{[h-1]} \rightarrow \mathcal{A}$.

Set $\pi_{\hat{s}'}(x_\ell) = \psi_{\hat{s}'}(\hat{f}_\ell(x_\ell))$, $\ell \in [h-1]$, $x_\ell \in \mathcal{X}_\ell$

end for

Let $\Pi_h = (\pi_{\hat{s}})_{\hat{s} \in \widehat{\mathcal{S}}_h}$.

Algorithm 5 Dynamic Programming for Reaching a State

Input: target state $\hat{s}^* \in \widehat{\mathcal{S}}_h$, transition probabilities $\hat{p}(\hat{s}' | \hat{s}, a)$ for all $\hat{s} \in \widehat{\mathcal{S}}_\ell$, $a \in \mathcal{A}$, $\hat{s}' \in \widehat{\mathcal{S}}_{\ell+1}$, $\ell \in [h-1]$

Output: policy $\psi : \widehat{\mathcal{S}}_{[h-1]} \rightarrow \mathcal{A}$ maximizing $\hat{\mathbb{P}}^\psi(\hat{s}^*)$.

Let $v(\hat{s}^*) = 1$ and let $v(\hat{s}) = 0$ for all other $\hat{s} \in \widehat{\mathcal{S}}_h$.

for $l = h-1, h-2, \dots, 1$ **do**

for $\hat{s} \in \widehat{\mathcal{S}}_\ell$ **do**

$$\psi(\hat{s}) = \max_{a \in \mathcal{A}} \left[\sum_{\hat{s}' \in \widehat{\mathcal{S}}_{\ell+1}} v(\hat{s}') \hat{p}(\hat{s}' | \hat{s}, a) \right]$$

$$v(\hat{s}) = \sum_{\hat{s}' \in \widehat{\mathcal{S}}_{\ell+1}} v(\hat{s}') \hat{p}(\hat{s}' | \hat{s}, a = \psi(\hat{s}))$$

end for

end for