

# Appendix

## A Technical lemmas

**Lemma A.1** (Multiplicative Chernoff bound [Chernoff et al. \(1952\)](#)). *Let  $X$  be a Binomial random variable with parameter  $p, n$ . For any  $\delta > 0$ , we have that*

$$\mathbb{P}[X < (1 - \delta)pn] < \left( \frac{e^{-\delta}}{(1 - \delta)^{1-\delta}} \right)^{np}.$$

A slightly looser bound that suffices for our propose:

$$\mathbb{P}[X < (1 - \delta)pn] < e^{-\frac{\delta^2 pn}{2}}.$$

**Lemma A.2** (Hoeffdings Inequality [Sridharan \(2002\)](#)). *Let  $x_1, \dots, x_n$  be independent bounded random variables such that  $\mathbb{E}[x_i] = 0$  and  $|x_i| \leq \xi_i$  with probability 1. Then for any  $\epsilon > 0$  we have*

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n x_i \geq \epsilon\right) \leq e^{-\frac{2n^2 \epsilon^2}{\sum_{i=1}^n \xi_i^2}}.$$

**Lemma A.3** (Bernsteins Inequality). *Let  $x_1, \dots, x_n$  be independent bounded random variables such that  $\mathbb{E}[x_i] = 0$  and  $|x_i| \leq \xi$  with probability 1. Let  $\sigma^2 = \frac{1}{n} \sum_{i=1}^n \text{Var}[x_i]$ , then with probability  $1 - \delta$  we have*

$$\frac{1}{n} \sum_{i=1}^n x_i \leq \sqrt{\frac{2\sigma^2 \cdot \log(1/\delta)}{n}} + \frac{2\xi}{3n} \log(1/\delta)$$

**Lemma A.4** (McDiarmids Inequality ([Sridharan, 2002](#))). *Let  $x_1, \dots, x_n$  be independent random variables and  $S : X^n \rightarrow \mathbb{R}$  be a measurable function which is invariant under permutation and let the random variable  $Z$  be given by  $Z = S(x_1, x_2, \dots, x_n)$ . Assume  $S$  has bounded difference: i.e.*

$$\sup_{x_1, \dots, x_n, x'_i} |S(x_1, \dots, x_i, \dots, x_n) - S(x_1, \dots, x'_i, \dots, x_n)| \leq \xi_i,$$

then for any  $\epsilon > 0$  we have

$$\mathbb{P}(|Z - \mathbb{E}[Z]| \geq \epsilon) \leq 2e^{-\frac{2\epsilon^2}{\sum_{i=1}^n \xi_i^2}}.$$

**Lemma A.5** (Azuma-Hoeffding inequality). *Suppose  $X_k, k = 1, 2, 3, \dots$  is a martingale and  $|X_k - X_{k-1}| \leq c_k$  almost surely. Then for all positive integers  $N$  and any  $\epsilon > 0$ ,*

$$\mathbb{P}(|X_N - X_0| \geq \epsilon) \leq 2e^{-\frac{\epsilon^2}{2 \sum_{i=1}^N c_i^2}}.$$

**Lemma A.6** (Freedman's inequality [Tropp et al. \(2011\)](#)). *Let  $X$  be the martingale associated with a filter  $\mathcal{F}$  (i.e.  $X_i = \mathbb{E}[X|\mathcal{F}_i]$ ) satisfying  $|X_i - X_{i-1}| \leq M$  for  $i = 1, \dots, n$ . Denote  $W := \sum_{i=1}^n \text{Var}(X_i|\mathcal{F}_{i-1})$  then we have*

$$\mathbb{P}(|X - \mathbb{E}[X]| \geq \epsilon, W \leq \sigma^2) \leq 2e^{-\frac{\epsilon^2}{2(\sigma^2 + M\epsilon/3)}}.$$

Or in other words, with probability  $1 - \delta$ ,

$$|X - \mathbb{E}[X]| \leq \sqrt{8\sigma^2 \cdot \log(1/\delta)} + \frac{2M}{3} \cdot \log(1/\delta), \quad \text{Or } W \geq \sigma^2.$$

**Lemma A.7** (Best arm identification lower bound [Krishnamurthy et al. \(2016\)](#)). *For any  $A \geq 2$  and  $\tau \leq \sqrt{1/8}$  and any best arm identification algorithm that produces an estimate  $\hat{a}$ , there exists a multi-arm bandit problem for which the best arm  $a^*$  is  $\tau$  better than all others, but  $\mathbb{P}[\hat{a} \neq a^*] \geq 1/3$  unless the number of samples  $T$  is at least  $\frac{A}{72\tau^2}$ .*

## B On error metric for OPE

In this section, we discuss the metric considered in this work. Traditionally, most works directly use *Mean Square Error* (MSE)  $\mathbb{E}[(\hat{v}^\pi - v^\pi)^2]$  as the criterion for measuring OPE methods *e.g.* Thomas and Brunskill (2016); Thomas (2015); Thomas et al. (2017); Farajtabar et al. (2018), or equivalently, by proposing unbiased estimators and discussing its variance *e.g.* Jiang and Li (2016). Alternately, one can consider bounding the absolute difference between  $v^\pi$  and  $\hat{v}^\pi$  with high probability (*e.g.* Duan et al. (2020)), *i.e.*  $|\hat{v}^\pi - v^\pi| \leq \epsilon_{\text{prob}}$  *w.h.p.* Generally speaking, high probability bound can be seen as a stricter criterion compared to MSE since

$$\begin{aligned}\mathbb{E}[(\hat{v}^\pi - v^\pi)^2] &= \mathbb{E}[(\hat{v}^\pi - v^\pi)^2 \mathbf{1}_E] + \mathbb{E}[(\hat{v}^\pi - v^\pi)^2 \mathbf{1}_{E^c}] \\ &\leq \epsilon_{\text{prob}}(\delta)^2 \cdot (1 - \delta) + H^2 \cdot \delta,\end{aligned}$$

where  $E$  is the event that  $\epsilon_{\text{prob}}$  error holds and  $\delta$  is the failure probability. As a result, if both  $\delta$  and  $\epsilon_{\text{prob}}(\delta)$  can be controlled small, then the high probability bound implies a result for MSE bound. This is realistic, since  $\delta$  mostly appears inside the logarithmic term of  $\epsilon_{\text{prob}}(\delta)$  so the second term can be scaled to sufficiently small without affecting the polynomial dependence for the first term.

Table 2: Summary of Uniform OPE results for  $H$ -horizon non-stationary setting

Method/Analysis	Policy class	Guarantee	Sample complexity
Simulation Lemma	All policies	$\epsilon$ -uniform convergence	$O(H^4 S^2 / d_m \epsilon^2)$
Theorem 3.3	All policies	$\epsilon$ -uniform convergence	$O(H^4 S / d_m \epsilon^2)$
Theorem 3.5	All deterministic policies	$\epsilon$ -uniform convergence	$O(H^3 S / d_m \epsilon^2)$
Theorem 3.7	local policies	$\epsilon$ -uniform convergence	$O(H^3 / d_m \epsilon^2)$
Minimax lower bound (Theorem 3.8)	————	over class $\mathcal{M}_{d_m}$	$\Omega(H^3 / d_m \epsilon^2)$

## C Some preparations

In this section we present some results that are critical for proving the main theorems.

**Lemma C.1.** *For any  $0 < \delta < 1$ , there exists an absolute constant  $c_1$  such that when total episode  $n > c_1 \cdot 1/d_m \cdot \log(HSA/\delta)$ , then with probability  $1 - \delta$ ,*

$$n_{s_t, a_t} \geq n \cdot d_t^\mu(s_t, a_t)/2, \quad \forall s_t, a_t.$$

If state  $s_t$  is not accessible, then  $n_{s_t, a_t} = d_t^\mu(s_t, a_t) = 0$  so the lemma holds trivially.<sup>6</sup>

*Proof of Lemma C.1.* Define  $E := \{\exists t, s_t, a_t \text{ s.t. } n_{s_t, a_t} < nd_t^\mu(s_t, a_t)/2\}$ . Then combining the multiplicative Chernoff bound (Lemma A.1 in the Appendix) and a union bound over each  $t, s_t$  and  $a_t$ , we obtain

$$\begin{aligned}\mathbb{P}[E] &\leq \sum_t \sum_{s_t} \sum_{a_t} \mathbb{P}[n_{s_t, a_t} < nd_t^\mu(s_t, a_t)/2] \\ &\leq HSA \cdot e^{-\frac{n \cdot \min_{t, s_t, a_t} d_t^\mu(s_t, a_t)}{8}} = HSA \cdot e^{-\frac{n \cdot d_m}{8}} := \delta\end{aligned}$$

solving this for  $n$  then provides the stated result.  $\square$

Now we define:  $N := \min_{t, s_t, a_t} n_{s_t, a_t}$ , then above implies  $N \geq nd_m/2$  (recall  $d_m$  in Assumption 2.2). Now we aggregate only the first  $N$  pieces of data in each state-action  $(s_t, a_t)$ <sup>7</sup> of off-policy data  $\mathcal{D}$  and they consist of a new dataset  $\mathcal{D}' = \{(s_t, a_t, s_{t+1}^{(i)}, r_t^{(i)}) : i = 1, \dots, N; t \in [H]; s_t \in \mathcal{S}, a_t \in \mathcal{A}\}$ , and is a subset of  $\mathcal{D}$ . For the rest of paper, we will use either  $\mathcal{D}'$  or the original  $\mathcal{D}$  to create OPEMA  $\hat{v}^\pi$  (only for theoretical analysis purpose). Whether  $\mathcal{D}$  or  $\mathcal{D}'$  is used will be stated clearly in each context.

<sup>6</sup>In general, non-accessible state will not affect our results so to make our presentation succinct we will not mention non-accessible state for the rest of paper unless necessary.

<sup>7</sup>Note we can do this since by definition  $N \leq n_{s_t, a_t}$  for all  $s_t, a_t$ .

**Remark C.2.** It is worth mentioning that when use  $\mathcal{D}'$  to construct  $\hat{v}^\pi$ ,  $n_{s_t, a_t}^{\mathcal{D}'} = N$  for all  $s_t, a_t$ . Also,  $N := \min n_{s_t, a_t}^{\mathcal{D}}$  (note  $n_{s_t, a_t}^{\mathcal{D}}$  is the count from  $\mathcal{D}$ ) itself is a random variable and in the extreme case we could have  $N = 0$  and if that happens  $\hat{v}^\pi = 0$  (since in that case  $\hat{P}_t \equiv 0$  and  $\hat{d}_t^\pi$  is degenerated). However, there is only tiny probability  $N$  will be small, as guaranteed by Lemma C.1.

We wanted to point out that this technique of dropping certain amount of data, is not uncommon for analyzing model-based method in RL: e.g. Rmax exploration (Brafman and Tenenbholz, 2002) for online episodic setting (see [Jiang (2018), Notes on Rmax exploration] Section 2 Algorithm for tabular MDP. The data they use is the known set  $K$  with parameter  $m$ , in step3 data pairs observed more than  $m$  times are not recorded).

### C.1 Fictitious OPEMA estimator.

Similar to Xie et al. (2019); Yin and Wang (2020), we introduce an unbiased version of  $\hat{v}^\pi$  to fill in the gap at  $(s_t, a_t)$  where  $n_{s_t, a_t}$  is small. Concretely, every component in  $\hat{v}^\pi$  is substituted by the fictitious counterpart, i.e.  $\tilde{v}^\pi := \sum_{t=1}^H \langle \tilde{d}_t^\pi, \tilde{r}_t \rangle$ , with  $\tilde{d}_t^\pi = \tilde{P}_t^\pi \tilde{d}_{t-1}^\pi$  and  $\tilde{P}_t^\pi(s_t|s_{t-1}) = \sum_{a_{t-1}} \tilde{P}_t(s_t|s_{t-1}, a_{t-1})\pi(a_{t-1}|s_{t-1})$ . In particular, consider the high probability event in Lemma C.1, i.e. let  $E_t$  denotes the event  $\{n_{s_t, a_t} \geq nd_t^\mu(s_t, a_t)/2\}$ <sup>8</sup>, then we define

$$\begin{aligned}\tilde{r}_t(s_t, a_t) &= \hat{r}_t(s_t, a_t)\mathbf{1}(E_t) + r_t(s_t, a_t)\mathbf{1}(E_t^c) \\ \tilde{P}_{t+1}(\cdot|s_t, a_t) &= \hat{P}_{t+1}(\cdot|s_t, a_t)\mathbf{1}(E_t) + P_{t+1}(\cdot|s_t, a_t)\mathbf{1}(E_t^c).\end{aligned}$$

Similarly, for the OPEMA estimator uses data  $\mathcal{D}'$ , the fictitious estimator is set to be

$$\begin{aligned}\tilde{r}_t(s_t, a_t) &= \hat{r}_t(s_t, a_t)\mathbf{1}(E) + r_t(s_t, a_t)\mathbf{1}(E^c) \\ \tilde{P}_{t+1}(\cdot|s_t, a_t) &= \hat{P}_{t+1}(\cdot|s_t, a_t)\mathbf{1}(E) + P_{t+1}(\cdot|s_t, a_t)\mathbf{1}(E^c)\end{aligned}$$

where  $E$  denote the event  $\{N \geq nd_m/2\}$ .

$\tilde{v}^\pi$  creates a bridge between  $\hat{v}^\pi$  and  $v^\pi$  because of its unbiasedness and it is also bounded by  $H$  (see Lemma B.3 and Lemma B.5 in Yin and Wang (2020) for those preliminary results). Also,  $\tilde{v}^\pi$  is identical to  $\hat{v}^\pi$  with high probability, as stated by the following lemma.

**Lemma C.3.** For any  $0 < \delta < 1$ , there exists an absolute constant  $c_1$  such that when total episode  $n > c_1 d_m \cdot \log(HSA/\delta)$ , then with probability  $1 - \delta$ ,

$$\sup_{\pi \in \Pi} |\hat{v}^\pi - \tilde{v}^\pi| = 0.$$

*Proof.* This Lemma is a direct corollary of Lemma C.1 by considering the event  $E_1 := \{\exists t, s_t, a_t \text{ s.t. } n_{s_t, a_t} < nd_t^\mu(s_t, a_t)/2\}$  or  $\{N < nd_m/2\}$  since  $\hat{v}^\pi$  and  $\tilde{v}^\pi$  are identical on  $E_1^c$ .  $\square$

Note  $\hat{v}^\pi$  and  $\tilde{v}^\pi$  even equal to each other uniformly over all  $\pi$  in  $\Pi$ . This is not surprising since only logging policy  $\mu$  will decide if they are equal or not. This lemma shows how close  $\hat{v}^\pi$  and  $\tilde{v}^\pi$  are. Therefore in the following it suffices to consider the uniform convergence of  $\sup_{\pi \in \Pi} |\hat{v}^\pi - v^\pi|$ .

Next by using a fictitious analogy of state-action expression as in equation (1), we have:

$$\begin{aligned}\sup_{\pi \in \Pi} |\hat{v}^\pi - v^\pi| &= \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi, \tilde{r}_t \rangle - \sum_{t=1}^H \langle d_t^\pi, r_t \rangle \right| \\ &= \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi, \tilde{r}_t \rangle - \sum_{t=1}^H \langle \tilde{d}_t^\pi, r_t \rangle + \sum_{t=1}^H \langle \tilde{d}_t^\pi, r_t \rangle - \sum_{t=1}^H \langle d_t^\pi, r_t \rangle \right| \\ &\leq \underbrace{\sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right|}_{(*)} + \underbrace{\sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi, \tilde{r}_t - r_t \rangle \right|}_{(**)}\end{aligned}\tag{3}$$

<sup>8</sup>More rigorously,  $E_t$  depends on the specific pair  $s_t, a_t$  and should be written as  $E_t(s_t, a_t)$ . However, for brevity we just use  $E_t$  and this notation should be clear in each context.

We first deal with  $(**)$  by the following lemma.

**Lemma C.4.** *We have with probability  $1 - \delta$ :*

$$\sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi, \tilde{r}_t - r_t \rangle \right| \leq O\left(\sqrt{\frac{H^2 \log(HSA/\delta)}{n \cdot d_m}}\right)$$

*Proof of Lemma C.4.* Since  $|\langle \tilde{d}_t^\pi, \tilde{r}_t - r_t \rangle| \leq \|\tilde{d}_t^\pi\|_1 \cdot \|\tilde{r}_t - r_t\|_\infty$ , we obtain

$$\left| \sum_{t=1}^H \langle \tilde{d}_t^\pi, \tilde{r}_t - r_t \rangle \right| \leq \sum_{t=1}^H \|\tilde{d}_t^\pi\|_1 \cdot \|\tilde{r}_t - r_t\|_\infty = \sum_{t=1}^H \|\tilde{r}_t - r_t\|_\infty,$$

where we used  $\tilde{d}_t^\pi(\cdot)$  is a probability distribution. Therefore above expression further indicates  $\sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi, \tilde{r}_t - r_t \rangle \right| \leq \sum_{t=1}^H \|\tilde{r}_t - r_t\|_\infty$ . Now by a union bound and Hoeffding inequality (Lemma A.2),

$$\begin{aligned} \mathbb{P}(\sup_t \|\tilde{r}_t - r_t\|_\infty > \epsilon) &= \mathbb{P}(\sup_{t, s_t, a_t} |\tilde{r}_t(s_t, a_t) - r_t(s_t, a_t)| > \epsilon) \\ &\leq HSA \cdot \mathbb{P}(|\tilde{r}_t(s_t, a_t) - r_t(s_t, a_t)| > \epsilon) \\ &= HSA \cdot \mathbb{P}(|\hat{r}_t(s_t, a_t) - r_t(s_t, a_t)| \mathbf{1}(E_t) > \epsilon) \\ &\leq 2HSA \cdot \mathbb{E}[\mathbb{E}[e^{-2n_{s_t, a_t} \epsilon^2} | E_t]] \\ &\leq 2HSA \cdot \mathbb{E}[\mathbb{E}[e^{-nd_m \epsilon^2} | E_t]] = 2HSA \cdot e^{-nd_m \epsilon^2} := \frac{\delta}{2}. \end{aligned}$$

where we use  $\mathbb{P}(A) = \mathbb{E}[\mathbf{1}_A] = \mathbb{E}[\mathbb{E}[\mathbf{1}_A | X]]$ . Solving for  $\epsilon$ , then it follows:

$$\sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi, \tilde{r}_t - r_t \rangle \right| \leq \sum_{t=1}^H \|\tilde{r}_t - r_t\|_\infty \leq O\left(\sqrt{\frac{H^2 \log(HSA/\delta)}{n \cdot d_m}}\right)$$

with probability  $1 - \delta$ . The case for  $E = \{N \geq nd_m/2\}$  can be proved easily in a similar way.  $\square$

Note that in order to measure the randomness in reward, sample complexity  $n$  only has dependence of order  $H^2$ , this result implies random reward will only cause error of lower order dependence in  $H$ . Therefore, in many RL literature deterministic reward is directly assumed. Next we consider  $(*)$  in (3) by decomposing  $\sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle$  into a martingale type representation. This is the key for our proof since with it we can use either uniform concentration inequalities or martingale concentration inequalities to prove efficiency.

## C.2 Decomposition of $\sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle$

Let  $\tilde{d}_t^\pi \in \mathbb{R}^{S \cdot A}$  denote the marginal state-action probability vector,  $\pi_t \in \mathbb{R}^{(S \cdot A) \times S}$  is the policy matrix with  $(\pi_t)_{(s_t, a_t), s_t} = \pi_t(a_t | s_t)$  and  $(\pi_t)_{(s_t, a_t), s} = 0$  for  $s \neq s_t$ . Moreover, let state-action transition matrix  $T_t \in \mathbb{R}^{S \times (S \cdot A)}$  to be  $(T_t)_{s_t, (s_{t-1}, a_{t-1})} = P_t(s_t | s_{t-1}, a_{t-1})$ , then we have

$$\tilde{d}_t^\pi = \pi_t \tilde{T}_t \tilde{d}_{t-1}^\pi \tag{4}$$

$$d_t^\pi = \pi_t T_t d_{t-1}^\pi. \tag{5}$$

Therefore we have

$$\tilde{d}_t^\pi - d_t^\pi = \pi_t (\tilde{T}_t - T_t) \tilde{d}_{t-1}^\pi + \pi_t T_t (\tilde{d}_{t-1}^\pi - d_{t-1}^\pi) \tag{6}$$

recursively apply this formula, we have

$$\tilde{d}_t^\pi - d_t^\pi = \sum_{h=2}^t \Gamma_{h+1:t} \pi_h (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi + \Gamma_{1:t} (\tilde{d}_1^\pi - d_1^\pi) \tag{7}$$

where  $\Gamma_{h:t} = \prod_{v=h}^t \pi_v T_v$  and  $\Gamma_{t+1:t} := 1$ . Now let  $X = \sum_{t=1}^H \langle r_t, \tilde{d}_t^\pi - d_t^\pi \rangle$ , then we have the following:

**Theorem C.5** (martingale decomposition of  $X$ : Restate of the fictitious version of Lemma 3.1). *We have:*

$$X = \sum_{h=2}^H \langle V_h^\pi(s), ((\tilde{T}_h - T_h)\tilde{d}_{h-1}^\pi)(s) \rangle + \langle V_1^\pi(s), (\tilde{d}_1^\pi - d_1^\pi)(s) \rangle,$$

where the inner product is taken w.r.t states.

*Proof of Theorem C.5.* By applying (7) and the change of summation, we have

$$\begin{aligned} X &= \sum_{t=1}^H \left( \sum_{h=2}^t \langle r_t, \Gamma_{h+1:t} \pi_h (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle + \langle r_t, \Gamma_{1:t} (\tilde{d}_1^\pi - d_1^\pi) \rangle \right) \\ &= \sum_{t=1}^H \left( \sum_{h=2}^t \langle r_t, \Gamma_{h+1:t} \pi_h (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle \right) + \sum_{h=1}^H \langle r_t, \Gamma_{1:t} (\tilde{d}_1^\pi - d_1^\pi) \rangle \\ &= \sum_{t=2}^H \left( \sum_{h=2}^t \langle r_t, \Gamma_{h+1:t} \pi_h (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle \right) + \sum_{h=1}^H \langle r_t, \Gamma_{1:t} (\tilde{d}_1^\pi - d_1^\pi) \rangle \\ &= \sum_{h=2}^H \left( \sum_{t=h}^H \langle r_t, \Gamma_{h+1:t} \pi_h (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle \right) + \sum_{h=1}^H \langle (\pi_1^T \Gamma_{1:t}^T r_t)(s), (\tilde{d}_1^\pi - d_1^\pi)(s) \rangle \\ &= \sum_{h=2}^H \left( \underbrace{\left\langle \sum_{t=h}^H \pi_h^T \Gamma_{h+1:t}^T r_t, (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \right\rangle}_{V_h^\pi(s)} \right) + \underbrace{\left\langle \left( \sum_{h=1}^H \pi_1^T \Gamma_{1:t}^T r_t \right)(s), (\tilde{d}_1^\pi - d_1^\pi)(s) \right\rangle}_{V_1^\pi(s)} \end{aligned}$$

□

## D Proof of uniform convergence in OPE with full policies using standard uniform concentration tools: Theorem 3.3

As a reminder for the reader, the OPEMA estimator used in this section is with data subset  $\mathcal{D}'$ . Also, by Lemma C.4 we only need to consider  $\sup_{\pi \in \Pi} |\sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle|$ .

**Theorem D.1.** *There exists an absolute constant  $c$  such that if  $n > c \cdot \frac{1}{d_m} \cdot \log(HSA/\delta)$ , then with probability  $1 - \delta$ , we have:*

$$\sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right| \leq O\left(\sqrt{\frac{H^4 \log(HSA/\delta)}{nd_m}}\right) + \mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right| \right]$$

*Proof of Theorem D.1.* Note in data  $\mathcal{D}' = \{(s_t, a_t, s_{t+1}^{(i)}) : i = 1, \dots, N; t = 1, \dots, H; s_t \in \mathcal{S}, a_t \in \mathcal{A}\}$ <sup>9</sup>, not only  $s_{t+1}^{(i)}$  but also  $N$  are random variables.

We first conditional on  $N$ , then  $(s_t, a_t, s_{t+1}^{(i)})$ 's are independent samples for all  $i, s_t, a_t$  since any sample will not contain information about other samples. Therefore we can regroup  $\mathcal{D}'$  into  $N$  independent samples with  $\mathcal{D}' = \{X^{(i)} : i = 1, \dots, N\}$  where  $X^{(i)} = \{(s_t, a_t, s_{t+1}^{(i)}), t = 1, \dots, H; s_t \in \mathcal{S}, a_t \in \mathcal{A}\}$ . Now for any  $i_0$ , change  $X^{(i_0)}$  to  $X'^{(i_0)} = \{(s_t, a_t, s_{t+1}'^{(i_0)}), t = 1, \dots, H; s_t \in \mathcal{S}, a_t \in \mathcal{A}\}$  and keep the rest  $N - 1$  data the same, use this data to

<sup>9</sup>Here we do not include  $r_t^{(i)}$  since the quantity  $\sup_{\pi \in \Pi} |\sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle|$  only contains the mean reward function  $r_t$ .

create new estimator with state-action transition  $\tilde{d}'^\pi$ , then we have

$$\begin{aligned}
 & \left| \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right| - \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}'^\pi - d_t^\pi, r_t \rangle \right| \right| \\
 & \leq \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle - \sum_{t=1}^H \langle \tilde{d}'^\pi - d_t^\pi, r_t \rangle \right| \\
 & \leq \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle - \sum_{t=1}^H \langle \tilde{d}'^\pi - d_t^\pi, r_t \rangle \right| \\
 & = \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - \tilde{d}'^\pi, r_t \rangle \right| \\
 & = \sup_{\pi \in \Pi} \left| \sum_{h=2}^H \langle \tilde{V}_h'^\pi, (\tilde{T}_h - \tilde{T}_h') \tilde{d}_{h-1}^\pi \rangle + \langle \tilde{V}_1'^\pi, \tilde{d}_1^\pi - \tilde{d}'^\pi \rangle \right|,
 \end{aligned}$$

where the last equation comes from the trick that substitutes  $d_t^\pi$  by  $\tilde{d}'^\pi$  in Theorem C.5. By definition, the above equals to

$$\begin{aligned}
 & = \sup_{\pi \in \Pi} \left| \sum_{h=2}^H \langle \tilde{V}_h'^\pi, (\hat{T}_h - \hat{T}_h') \hat{d}_{h-1}^\pi \rangle + \langle \tilde{V}_1'^\pi, \hat{d}_1^\pi - \hat{d}'^\pi \rangle \right| \cdot \mathbf{1}(E) \\
 & \leq \sup_{\pi \in \Pi} \left( \sum_{h=2}^H \|(\hat{T}_h - \hat{T}_h')^T \tilde{V}_h'^\pi\|_\infty \|\hat{d}_{h-1}^\pi\|_1 + |\langle \tilde{V}_1'^\pi, \hat{d}_1^\pi - \hat{d}'^\pi \rangle| \right) \cdot \mathbf{1}(E) \\
 & \leq \sup_{\pi \in \Pi} \left( \sum_{h=2}^H \|(\hat{T}_h - \hat{T}_h')^T \tilde{V}_h'^\pi\|_\infty + |\langle \tilde{V}_1'^\pi, \hat{d}_1^\pi - \hat{d}'^\pi \rangle| \right) \cdot \mathbf{1}(E)
 \end{aligned}$$

Note the change of a single  $X^{(i_0)}$  will only change two entries of each row of  $(\hat{T}_h - \hat{T}_h')^T$  by  $1/N$  since with data  $\mathcal{D}'$ ,  $n_{s_t, a_t} = N$  for all  $s_t, a_t$ . Or in other words, given  $E$ ,

$$\hat{T}_h^T - \hat{T}_h'^T = \begin{bmatrix} 0 & \dots & 0 & \frac{1}{N} & 0 & \dots & -\frac{1}{N} & \dots & 0 \\ 0 & \frac{1}{N} & 0 & \dots & -\frac{1}{N} & \dots & \dots & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ -\frac{1}{N} & 0 & \dots & 0 & \dots & \dots & 0 & \dots & \frac{1}{N} \end{bmatrix},$$

where the locations of  $1/N, -1/N$  in each row are random as it depends on how different is  $X^{(i_0)}$  from  $X^{(i_0)}$ . However, based on this fact, it is enough for us to guarantee

$$\|(\hat{T}_h - \hat{T}_h')^T \tilde{V}_h'^\pi\|_\infty \leq \frac{2}{N} \|\tilde{V}_h'^\pi\|_\infty \leq \frac{2}{N} (H - h + 1) \leq \frac{2}{N} H$$

and same result holds for  $|\langle \tilde{V}_1'^\pi, \hat{d}_1^\pi - \hat{d}'^\pi \rangle| \leq 2H/N$  given  $N$ .

Combine all the results above, for a single change of  $X^{(i_0)}$  we have

$$\left| \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right| - \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}'^\pi - d_t^\pi, r_t \rangle \right| \right| \leq 2 \frac{H^2}{N} \mathbf{1}(E) \leq 2 \frac{H^2}{N}$$

for any fixed  $N$ . If we let  $Z = S(X^{(1)}, \dots, X^{(N)}) = \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right|$ , then for a given  $N$  by independence and above bounded difference condition we can apply Mcdiarmid inequality Lemma A.4 (where  $\xi_i = 2H^2/N$ ) to obtain

$$\mathbb{P}(|Z - \mathbb{E}[Z]| \geq \epsilon | N) \leq 2e^{-\frac{N\epsilon^2}{2H^4}} := \frac{\delta}{2} \tag{8}$$

Now note when  $n > O(\frac{1}{d_m} \cdot \log(HSA/\delta))$ , by Lemma C.1 we can obtain  $N > nd_m/2$  with probability  $1 - \delta/2$ , combining this result and solving  $\epsilon$  in (8), we have

$$\sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right| \leq O\left(\sqrt{\frac{H^4 \log(HSA/\delta)}{n \cdot d_m}}\right) + \mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right| \right]$$

with probability  $1 - \delta$ . □

Next before bounding  $\mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right| \right]$ , we first give a useful lemma.

Let  $\gamma \in (0, 1)$  to be any threshold parameter. Then we first have the following lemma:

**Lemma D.2.** *Recall by definition  $P_h(s_h, |s_{h-1}, a_{h-1}) = T_h(s_h, |s_{h-1}, a_{h-1})$ . It holds that with probability  $1 - \delta$ , for all  $t, s_t, a_t \in [H], \mathcal{S}, \mathcal{A}$ : if  $P_h(s_h | s_{h-1}, a_{h-1}) \leq \gamma$ , then*

$$\left| \tilde{T}_h(s_h | s_{h-1}, a_{h-1}) - T_h(s_h | s_{h-1}, a_{h-1}) \right| \leq \sqrt{\frac{\gamma \log(HSA/\delta)}{2nd_m}} + \frac{2 \log(HSA/\delta)}{3nd_m};$$

if  $P_h(s_h, |s_{h-1}, a_{h-1}) > \gamma$ , then

$$\left| \frac{\tilde{T}_h(s_h | s_{h-1}, a_{h-1}) - T_h(s_h | s_{h-1}, a_{h-1})}{T_h(s_h | s_{h-1}, a_{h-1})} \right| \leq \sqrt{\frac{\log(HSA/\delta)}{2nd_m \gamma}} + \frac{2 \log(HSA/\delta)}{3nd_m \gamma};$$

*Proof.* First consider the case where  $P_h(s_h | s_{h-1}, a_{h-1}) \leq \gamma$ .

$$\tilde{T}_h(s_h | s_{h-1}, a_{h-1}) - T_h(s_h | s_{h-1}, a_{h-1}) = \frac{1}{n_{s_{h-1}, a_{h-1}}} \sum_{i=1}^{n_{s_{h-1}, a_{h-1}}} \left( \mathbb{1}[s_h^{(i)}] - T_h(s_h | s_{h-1}, a_{h-1}) \right) \mathbb{1}(E_h),$$

since  $\text{Var}[\mathbb{1}[s_h^{(i)}] | s_{h-1}, a_{h-1}] = P_h(s_h | s_{h-1}, a_{h-1})(1 - P_h(s_h | s_{h-1}, a_{h-1})) \leq P_h(s_h | s_{h-1}, a_{h-1}) \leq \gamma$ , therefore by Lemma A.3,

$$\left| \tilde{T}_h(s_h | s_{h-1}, a_{h-1}) - T_h(s_h | s_{h-1}, a_{h-1}) \right| \leq \mathbb{1}(E_h) \left( \sqrt{\frac{\gamma \log(1/\delta)}{n_{s_{h-1}, a_{h-1}}}} + \frac{2 \log(1/\delta)}{n_{s_{h-1}, a_{h-1}}} \right) \leq \sqrt{\frac{\gamma \log(1/\delta)}{2nd_m}} + \frac{2 \log(1/\delta)}{3nd_m};$$

Second, when  $P_h(s_h | s_{h-1}, a_{h-1}) > \gamma$ .

$$\frac{\tilde{T}_h(s_h | s_{h-1}, a_{h-1}) - T_h(s_h | s_{h-1}, a_{h-1})}{T_h(s_h | s_{h-1}, a_{h-1})} = \frac{1}{n_{s_{h-1}, a_{h-1}}} \sum_{i=1}^{n_{s_{h-1}, a_{h-1}}} \left( \frac{\mathbb{1}[s_h^{(i)}]}{T_h(s_h | s_{h-1}, a_{h-1})} - 1 \right) \mathbb{1}(E_h),$$

since

$$\text{Var} \left[ \frac{\mathbb{1}[s_h^{(i)}]}{T_h(s_h | s_{h-1}, a_{h-1})} \middle| s_{h-1}, a_{h-1} \right] \leq \frac{1}{T_h(s_h | s_{h-1}, a_{h-1})^2} \text{Var} [\mathbb{1}[s_h^{(i)}] | s_{h-1}, a_{h-1}] \leq \frac{1}{T_h(s_h | s_{h-1}, a_{h-1})} \leq \frac{1}{\gamma},$$

and since  $\frac{\mathbb{1}[s_h^{(i)}]}{T_h(s_h | s_{h-1}, a_{h-1})} \leq 1/\gamma$ , again by Bernstein inequality we have

$$\left| \frac{\tilde{T}_h(s_h | s_{h-1}, a_{h-1}) - T_h(s_h | s_{h-1}, a_{h-1})}{T_h(s_h | s_{h-1}, a_{h-1})} \right| \leq \sqrt{\frac{\log(1/\delta)}{2nd_m \gamma}} + \frac{2 \log(1/\delta)}{3nd_m \gamma};$$

apply the union bound over  $t, s_t, a_t$  we obtain the stated result. □

**Bounding**  $\mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right| \right]$ . First note by Theorem C.5:

$$\mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right| \right] \leq \sum_{h=2}^H \mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \langle v_h^\pi(s), ((\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi(s)) \rangle \right| \right] + \mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \langle V_1^\pi(s), (\tilde{d}_1^\pi - d_1^\pi(s)) \rangle \right| \right],$$

so it suffices to bound each  $\mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \langle V_h^\pi(s), ((\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi(s)) \rangle \right| \right]$ . First of all,

$$\begin{aligned} & \mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \langle V_h^\pi(s), ((\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi(s)) \rangle \right| \right] \\ &= \mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \sum_{s_h, s_{h-1}, a_{h-1}} V_h^\pi(s_h) (\tilde{T}_h - T_h)(s_h | s_{h-1}, a_{h-1}) \tilde{d}_{h-1}^\pi(s_{h-1}, a_{h-1}) \right| \right] \\ &\leq \mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \sum_{s_h, s_{h-1}, a_{h-1}} V_h^\pi(s_h) (\tilde{T}_h - T_h)(s_h | s_{h-1}, a_{h-1}) \tilde{d}_{h-1}^\pi(s_{h-1}, a_{h-1}) \right| \cdot \mathbb{1}[T_h(s_h | s_{h-1}, a_{h-1}) > \gamma] \right] \\ &+ \mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \sum_{s_h, s_{h-1}, a_{h-1}} V_h^\pi(s_h) (\tilde{T}_h - T_h)(s_h | s_{h-1}, a_{h-1}) \tilde{d}_{h-1}^\pi(s_{h-1}, a_{h-1}) \right| \cdot \mathbb{1}[T_h(s_h | s_{h-1}, a_{h-1}) \leq \gamma] \right] \\ &= \underbrace{\mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \sum_{s_h, s_{h-1}, a_{h-1}} V_h^\pi(s_h) T_h(s_h | s_{h-1}, a_{h-1}) \tilde{d}_{h-1}^\pi(s_{h-1}, a_{h-1}) \frac{\tilde{T}_h - T_h}{T_h}(s_h | s_{h-1}, a_{h-1}) \right| \cdot \mathbb{1}[T_h > \gamma] \right]}_{(a)} \\ &+ \underbrace{\mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \sum_{s_h, s_{h-1}, a_{h-1}} V_h^\pi(s_h) \tilde{d}_{h-1}^\pi(s_{h-1}, a_{h-1}) (\tilde{T}_h - T_h)(s_h | s_{h-1}, a_{h-1}) \right| \cdot \mathbb{1}[T_h(s_h | s_{h-1}, a_{h-1}) \leq \gamma] \right]}_{(b)}, \end{aligned}$$

Apply Lemma D.2 with  $\delta'/2$  where  $\delta' = \delta/H$ , then

$$\begin{aligned} (a) &\leq \sup_{\pi \in \Pi} \left| \sum_{s_h, s_{h-1}, a_{h-1}} V_h^\pi(s_h) T_h(s_h | s_{h-1}, a_{h-1}) \tilde{d}_{h-1}^\pi(s_{h-1}, a_{h-1}) \left( \sqrt{\frac{\log(2HSA/\delta')}{2nd_m\gamma}} + \frac{2\log(2HSA/\delta')}{3nd_m\gamma} \right) \right| \left(1 - \frac{\delta'}{2}\right) \\ &\quad + H\delta'/2 \\ &\leq \sup_{\pi \in \Pi} \left| \sum_{s_h, s_{h-1}, a_{h-1}} V_h^\pi(s_h) T_h(s_h | s_{h-1}, a_{h-1}) \tilde{d}_{h-1}^\pi(s_{h-1}, a_{h-1}) \left( \sqrt{\frac{\log(2H^2SA/\delta)}{2nd_m\gamma}} + \frac{2\log(2H^2SA/\delta)}{3nd_m\gamma} \right) \right| \\ &\quad + \delta/2 \\ &\leq \sup_{\pi \in \Pi} \left| H \left( \sqrt{\frac{2\log(H^2SA/\delta)}{2nd_m\gamma}} + \frac{2\log(2H^2SA/\delta)}{3nd_m\gamma} \right) \right| + \delta/2 = H \left( \sqrt{\frac{\log(2H^2SA/\delta)}{2nd_m\gamma}} + \frac{2\log(2H^2SA/\delta)}{3nd_m\gamma} \right) + \delta/2, \\ (b) &\leq \sup_{\pi \in \Pi} \left| \sum_{s_h, s_{h-1}, a_{h-1}} V_h^\pi(s_h) \tilde{d}_{h-1}^\pi(s_{h-1}, a_{h-1}) \left( \sqrt{\frac{\gamma \log(2HSA/\delta)}{2nd_m}} + \frac{2\log(2HSA/\delta)}{3nd_m} \right) \right| \left(1 - \frac{\delta'}{2}\right) + H \frac{\delta'}{2} \\ &\leq \sup_{\pi \in \Pi} \left| HS \left( \sqrt{\frac{\gamma \log(2H^2SA/\delta)}{2nd_m}} + \frac{2\log(2H^2SA/\delta)}{3nd_m} \right) \right| + \frac{\delta}{2} = HS \left( \sqrt{\frac{\gamma \log(2H^2SA/\delta)}{2nd_m}} + \frac{2\log(2H^2SA/\delta)}{3nd_m} \right) + \frac{\delta}{2}, \end{aligned}$$



Hence we have for any  $\gamma$ ,

$$\begin{aligned} & \mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \langle V_h^\pi(s), ((\tilde{T}_h - T_h)\tilde{d}_{h-1}^\pi)(s) \rangle \right| \right] \\ & \leq H \left( \sqrt{\frac{\log(2H^2SA/\delta)}{2nd_m\gamma}} + \frac{2\log(2H^2SA/\delta)}{3nd_m\gamma} \right) + HS \left( \sqrt{\frac{\gamma\log(2H^2SA/\delta)}{2nd_m}} + \frac{2\log(2H^2SA/\delta)}{3nd_m} \right) + \delta \end{aligned}$$

In particular, choose  $\gamma = 1/S < 1$ , then above becomes

$$\mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \langle V_h^\pi(s), ((\tilde{T}_h - T_h)\tilde{d}_{h-1}^\pi)(s) \rangle \right| \right] \leq \sqrt{\frac{2H^2S\log(2H^2SA/\delta)}{nd_m}} + \frac{4HS\log(2H^2SA/\delta)}{3nd_m} + \delta$$

Critically, above holds for any  $\forall 1 > \delta > 0$ . Based on theorem condition  $n > c \cdot 1/d_m \log(HSA/\theta) > c \cdot 1/d_m$ <sup>10</sup>, choose  $\delta = \frac{c}{nd_m}$ , then above is further less equal to

$$\sqrt{\frac{2H^2S\log(2nH^2SA)}{nd_m}} + \frac{4HS\log(2nH^2SA)}{3nd_m} + \frac{c}{nd_m} \leq \sqrt{\frac{2H^2S\log(2nH^2SA)}{nd_m}} + C \cdot \frac{HS\log(2nH^2SA)}{3nd_m}$$

where  $C$  is a new constant absorbs  $1/nd_m$ . If we further reducing it to

Finally, summing over all  $H$ , and again using new constant  $C'$  to absorb higher order term, we obtain

$$\mathbb{E} \left[ \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right| \right] \leq C' \sqrt{\frac{H^4S\log(nHSA)}{nd_m}}$$

Combing this with Theorem D.1 and Lemma C.4, we have proved Theorem 3.3.

**Remark D.3.** The key for proving this uniform convergence bound is that applying concentration inequality only to terms that are independent of the policies, i.e.  $\tilde{T}_h(s_h|s_{h-1}, a_{h-1}) - T_h(s_h|s_{h-1}, a_{h-1})$ . Therefore when taking supremum over policies, high probability event holds with same probability without decay.

## E Proof of uniform convergence in OPE with deterministic policies using martingale concentration inequalities: Theorem 3.5

A reminder that all results in this section use data  $\mathcal{D}$  for OPEMA estimator  $\hat{v}^\pi$ .

### E.1 Martingale concentration result on $\sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle$ .

Let  $X = \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle$  and  $\mathcal{D}_h := \{s_t^{(i)}, a_t^{(i)} : t = 1, \dots, h\}_{i=1}^n$ . Since  $\mathcal{D}_h$  forms a filtration, then by law of total expectation we have  $X_t = \mathbb{E}[X|\mathcal{D}_t]$  is martingale. Moreover, we have

**Lemma E.1.**

$$X_t := \mathbb{E}[X|\mathcal{D}_t] = \sum_{h=2}^t \langle V_h^\pi, (\tilde{T}_h - T_h)\tilde{d}_{h-1}^\pi \rangle + \langle V_1^\pi, \tilde{d}_1^\pi - d_1^\pi \rangle.$$

*Proof of Lemma E.1.* By martingale decomposition Theorem C.5 and note  $\tilde{T}_i, \tilde{d}_i^\pi$  are measurable *w.r.t.*  $\mathcal{D}_t$  for  $i = 1, \dots, t$ , so we have

$$\mathbb{E}[X|\mathcal{D}_t] = \sum_{h=t+1}^H \mathbb{E} \left[ \langle V_h^\pi, (\tilde{T}_h - T_h)\tilde{d}_{h-1}^\pi \rangle \middle| \mathcal{D}_t \right] + \sum_{h=2}^t \langle V_h^\pi, (\tilde{T}_h - T_h)\tilde{d}_{h-1}^\pi \rangle + \langle V_1^\pi, (\tilde{d}_1^\pi - d_1^\pi) \rangle.$$

<sup>10</sup>Note the  $\theta$  in  $\log(HSA/\theta)$  is identical to the failure probability in Theorem D.1

Note for  $h \geq t+1$ ,  $\mathcal{D}_t \subset \mathcal{D}_{h-1}$ , so by total law of expectation (tower property) we have

$$\begin{aligned} & \mathbb{E} \left[ \langle V_h^\pi, (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle \middle| \mathcal{D}_t \right] \\ &= \mathbb{E} \left[ \mathbb{E} \left[ \langle V_h^\pi, (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle \middle| \mathcal{D}_{h-1} \right] \middle| \mathcal{D}_t \right] \\ &= \mathbb{E} \left[ \langle V_h^\pi, \mathbb{E} \left[ (\tilde{T}_h - T_h) \middle| \mathcal{D}_{h-1} \right] \tilde{d}_{h-1}^\pi \rangle \middle| \mathcal{D}_t \right] = 0 \end{aligned}$$

where the last equality uses  $\tilde{T}_h$  is unbiased of  $T_h$  given  $\mathcal{D}_{h-1}$ . This gives the desired result.  $\square$

Next we show martingale difference  $|X_t - X_{t-1}|$  is bounded with high probability.

**Lemma E.2.** *With probability  $1 - \delta$ ,*

$$\sup_t |X_t - X_{t-1}| \leq O\left(\sqrt{\frac{H^2 \log(HSA/\delta)}{n \cdot d_m}}\right).$$

*Proof.*

$$|X_t - X_{t-1}| = \langle V_t^\pi, (\tilde{T}_t - T_t) \tilde{d}_{t-1}^\pi \rangle \leq \|(\tilde{T}_t - T_t)^T V_t^\pi\|_\infty \|\tilde{d}_{t-1}^\pi\|_1 = \|(\tilde{T}_t - T_t)^T V_t^\pi\|_\infty.$$

For any fixed pair  $(s_t, a_t)$ , we have

$$\begin{aligned} & ((\tilde{T}_t - T_t)^T V_t^\pi)(s_{t-1}, a_{t-1}) \\ &= \mathbb{1}(E_{t-1}) \cdot ((\hat{T}_t - T_t)^T V_t^\pi)(s_{t-1}, a_{t-1}) \\ &= \mathbb{1}(E_{t-1}) \cdot \sum_{s_t} V_t^\pi(s_t) (\hat{T}_t - T_t)(s_t | s_{t-1}, a_{t-1}) \\ &= \mathbb{1}(E_{t-1}) \cdot \left( \sum_{s_t} V_t^\pi(s_t) \hat{T}_t(s_t | s_{t-1}, a_{t-1}) - \mathbb{E}[V_t^\pi | s_{t-1}, a_{t-1}] \right) \\ &= \mathbb{1}(E_{t-1}) \cdot \left( \sum_{s_t} V_t^\pi(s_t) \frac{1}{n_{s_{t-1}, a_{t-1}}} \sum_{i=1}^n \mathbb{1}(s_t^{(i)} = s_t, s_{t-1}^{(i)} = s_{t-1}, a_{t-1}^{(i)} = a_{t-1}) - \mathbb{E}[V_t^\pi | s_{t-1}, a_{t-1}] \right) \\ &= \mathbb{1}(E_{t-1}) \cdot \left( \frac{1}{n_{s_{t-1}, a_{t-1}}} \sum_{i=1}^n V_t^\pi(s_t^{(i)}) \mathbb{1}(s_t^{(i)} = s_t, s_{t-1}^{(i)} = s_{t-1}, a_{t-1}^{(i)} = a_{t-1}) - \mathbb{E}[V_t^\pi | s_{t-1}, a_{t-1}] \right) \\ &= \mathbb{1}(E_{t-1}) \cdot \left( \frac{1}{n_{s_{t-1}, a_{t-1}}} \sum_{i: s_{t-1}^{(i)} = s_{t-1}, a_{t-1}^{(i)} = a_{t-1}} V_t^\pi(s_t^{(i)}) - \mathbb{E}[V_t^\pi | s_{t-1}, a_{t-1}] \right), \end{aligned}$$

where the fourth line uses the definition of  $\hat{T}_t$  and the fifth line uses the fact  $\sum_{s_t} V_t^\pi(s_t) \mathbb{1}(s_t^{(i)} = s_t, s_{t-1}^{(i)} = s_{t-1}, a_{t-1}^{(i)} = a_{t-1}) = V_t^\pi(s_t^{(i)}) \mathbb{1}(s_t^{(i)} = s_t, s_{t-1}^{(i)} = s_{t-1}, a_{t-1}^{(i)} = a_{t-1})$ .

Note  $\|V_t^\pi(\cdot)\|_\infty \leq H$  and also conditional on  $E_t$ ,  $n_{s_t, a_t} \geq nd_t^\mu(s_t, a_t)/2$ , therefore by Hoeffding's inequality and a Union bound we obtain with probability  $1 - \delta$

$$\sup_t |X_t - X_{t-1}| \leq O\left(\sqrt{\frac{H^2 \log(HSA/\delta)}{n \cdot \min_{t, s_t, a_t} d_t^\mu(s_t, a_t)}}\right) = O\left(\sqrt{\frac{H^2 \log(HSA/\delta)}{n \cdot d_m}}\right).$$

$\square$

Next we calculate the conditional variance of  $\text{Var}[X_{t+1} | \mathcal{D}_t]$ .

**Lemma E.3.** *We have the following decomposition of conditional variance:*

$$\text{Var}[X_{t+1} | \mathcal{D}_t] = \sum_{s_t, a_t} \frac{\tilde{d}_t^\pi(s_t, a_t)^2 \cdot \mathbb{1}(E_t)}{n_{s_t, a_t}} \cdot \text{Var}[V_{t+1}^\pi(s_{t+1}^{(1)}) | s_t^{(1)} = s_t, a_t^{(1)} = a_t]$$

*Proof.* Indeed,

$$\begin{aligned}
 \text{Var}[X_{t+1}|\mathcal{D}_t] &= \text{Var} \left[ \sum_{s_t, a_t} \sum_{s_{t+1}} V_{t+1}^\pi(s_{t+1}) (\tilde{T} - T)(s_{t+1}|s_t, a_t) \tilde{d}_t^\pi(s_t, a_t) \middle| \mathcal{D}_t \right] \\
 &= \sum_{s_t, a_t} \text{Var} \left[ \sum_{s_{t+1}} V_{t+1}^\pi(s_{t+1}) (\tilde{T} - T)(s_{t+1}|s_t, a_t) \middle| \mathcal{D}_t \right] \tilde{d}_t^\pi(s_t, a_t)^2 \\
 &= \sum_{s_t, a_t} \mathbb{1}(E_t) \cdot \text{Var} \left[ \sum_{s_{t+1}} V_{t+1}^\pi(s_{t+1}) \hat{T}(s_{t+1}|s_t, a_t) \middle| \mathcal{D}_t \right] \tilde{d}_t^\pi(s_t, a_t)^2 \\
 &= \sum_{s_t, a_t} \mathbb{1}(E_t) \cdot \text{Var} \left[ \sum_{s_{t+1}} V_{t+1}^\pi(s_{t+1}) \frac{1}{n_{s_t, a_t}} \sum_{i=1}^n \mathbb{1}(s_{t+1}^{(i)} = s_{t+1}, s_t^{(i)} = s_t, a_t^{(i)} = a_t) \middle| \mathcal{D}_t \right] \tilde{d}_t^\pi(s_t, a_t)^2 \\
 &= \sum_{s_t, a_t} \frac{\mathbb{1}(E_t)}{n_{s_t, a_t}^2} \cdot \text{Var} \left[ \sum_{i: s_t^{(i)} = s_t, a_t^{(i)} = a_t} V_{t+1}^\pi(s_{t+1}^{(i)}) \middle| \mathcal{D}_t \right] \tilde{d}_t^\pi(s_t, a_t)^2 \\
 &= \sum_{s_t, a_t} \frac{\tilde{d}_t^\pi(s_t, a_t)^2 \cdot \mathbb{1}(E_t)}{n_{s_t, a_t}} \cdot \text{Var}[V_{t+1}^\pi(s_{t+1}^{(1)}) | s_t^{(1)} = s_t, a_t^{(1)} = a_t]
 \end{aligned} \tag{9}$$

where the second equal sign comes from the fact that when conditional on  $\mathcal{D}_t$ , we can separate  $n$  episodes into  $SA$  groups and episodes from different groups are independent of each other. The third uses  $\mathbb{1}(E_t)$  is measurable w.r.t  $\mathcal{D}_t$ . Similarly, the last equal sign again uses  $n_{s_t, a_t}$  episodes are independent given  $\mathcal{D}_t$ .  $\square$

**Lemma E.4** (Yin and Wang (2020) Lemma 3.4). *For any policy  $\pi$  and any MDP.*

$$\begin{aligned}
 \text{Var}_\pi \left[ \sum_{t=1}^H r_t^{(1)} \right] &= \sum_{t=1}^H \left( \mathbb{E}_\pi \left[ \text{Var} \left[ r_t^{(1)} + V_{t+1}^\pi(s_{t+1}^{(1)}) \middle| s_t^{(1)}, a_t^{(1)} \right] \right] \right. \\
 &\quad \left. + \mathbb{E}_\pi \left[ \text{Var} \left[ \mathbb{E}[r_t^{(1)} + V_{t+1}^\pi(s_{t+1}^{(1)}) | s_t^{(1)}, a_t^{(1)}] \middle| s_t^{(1)} \right] \right] \right).
 \end{aligned}$$

This Lemma suggests if we can bound  $\tilde{d}_t^\pi$  by  $O(d_t^\pi)$  with high probability, then by Lemma E.3 we have w.h.p

$$\sum_{t=1}^H \text{Var}[X_{t+1}|\mathcal{D}_t] \leq O\left(\frac{1}{nd_m} \cdot \sum_{t=1}^H \mathbb{E}[\text{Var}[V_{t+1}^\pi(s_{t+1}^{(1)}) | s_t^{(1)}, a_t^{(1)}]]\right) \leq O\left(\frac{H^2}{nd_m}\right)$$

Note this gives only  $H^2$  dependence for  $\sum_{t=1}^H \text{Var}[X_{t+1}|\mathcal{D}_t]$  instead of a naive bound with  $H^3$  and helps us to save a  $H$  factor.

Next we show how to bound  $\tilde{d}_t^\pi$ .

## E.2 Bounding $\tilde{d}_t^\pi(s_t, a_t) - d_t^\pi(s_t, a_t)$

Our analysis is based on using martingale structure to derive bound on  $\tilde{d}_t^\pi(s_t, a_t) - d_t^\pi(s_t, a_t)$  for fixed  $t, s_t, a_t$  with probability  $1 - \delta/HSA$ , then use a union bound to get a bound for all  $\tilde{d}_t^\pi(s_t, a_t) - d_t^\pi(s_t, a_t)$  with probability  $1 - \delta$ .

Concretely, in (7) if we only extract the specific  $(s_t, a_t)$ , then we have

$$\tilde{d}_t^\pi(s_t, a_t) - d_t^\pi(s_t, a_t) = \sum_{h=2}^t (\Gamma_{h+1:t} \pi_h (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi)(s_t, a_t) + (\Gamma_{1:t} (\tilde{d}_1^\pi - d_1^\pi))(s_t, a_t),$$

here  $\tilde{d}_t^\pi(s_t, a_t) - d_t^\pi(s_t, a_t)$  already forms a martingale with filtration  $\mathcal{F}_t = \sigma(\mathcal{D}_t)$  and  $(\Gamma_{h+1:t} \pi_h (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi)(s_t, a_t)$  is the corresponding martingale difference since

$$\mathbb{E}[(\Gamma_{h+1:t} \pi_h (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi)(s_t, a_t) | \mathcal{F}_{h-1}] = (\Gamma_{h+1:t} \pi_h \mathbb{E}[(\tilde{T}_h - T_h) | \mathcal{F}_{h-1}] \tilde{d}_{h-1}^\pi)(s_t, a_t) = 0.$$

Now we fix specific  $(s_t, a_t)$ . Then denote  $(\Gamma_{h+1:t}\pi_h)(s_t, a_t) := \Gamma'_{h:t} \in \mathbb{R}^{1 \times S}$ , then we have

$$|(\Gamma_{h+1:t}\pi_h(\tilde{T}_h - T_h)\tilde{d}_{h-1}^\pi)(s_t, a_t)| = |\Gamma'_{h:t}(\tilde{T}_h - T_h)\tilde{d}_{h-1}^\pi| = |\langle (\tilde{T}_h - T_h)^T \Gamma'_{h:t}, \tilde{d}_{h-1}^\pi \rangle| \leq \|\Gamma'_{h:t}(\tilde{T}_h - T_h)\|_\infty \cdot 1.$$

Note here  $\Gamma'_{h:t}(\tilde{T}_h - T_h)$  is a row vector with dimension  $SA$ .

**Bounding  $\|\Gamma'_{h:t}(\tilde{T}_h - T_h)\|_\infty$**

In fact, for any given  $(s_{h-1}, a_{h-1})$ , we have

$$\begin{aligned} & \Gamma'_{h:t}(\tilde{T}_h - T_h)(s_{h-1}, a_{h-1}) = \mathbf{1}(E_t) \cdot \Gamma'_{h:t}(\hat{T}_h - T_h)(s_{h-1}, a_{h-1}) \\ &= \mathbf{1}(E_t) \cdot \Gamma'_{h:t} \left( \frac{1}{n_{s_{t-1}, a_{t-1}}} \sum_{i: s_{h-1}^{(i)} = s_{h-1}, a_{h-1}^{(i)} = a_{h-1}} \mathbf{e}_{s_h^{(i)}} - \mathbb{E}[\mathbf{e}_{s_h^{(1)}} | s_{h-1}^{(1)} = s_{h-1}, a_{h-1}^{(1)} = a_{h-1}] \right) \\ &= \mathbf{1}(E_t) \left( \frac{1}{n_{s_{t-1}, a_{t-1}}} \sum_{i: s_{h-1}^{(i)} = s_{h-1}, a_{h-1}^{(i)} = a_{h-1}} \Gamma'_{h:t}(s_h^{(i)}) - \mathbb{E}[\Gamma'_{h:t}(s_h^{(1)}) | s_{h-1}^{(1)} = s_{h-1}, a_{h-1}^{(1)} = a_{h-1}] \right) \end{aligned}$$

Note by definition  $\Gamma'_{h:t}(s_h^{(i)}) \leq 1$ , since  $(\Gamma_{h+1:t}\pi_h)(s_t, a_t) := \Gamma'_{h:t} \in \mathbb{R}^{1 \times S}$  and  $\Gamma_{h+1:t}, \pi_h$  are just probability transitions. Therefore by Hoeffding's inequality and law of total expectation, we have

$$\begin{aligned} & \mathbb{P}(|\Gamma'_{h:t}(\tilde{T}_h - T_h)(s_{h-1}, a_{h-1})| > \epsilon) = \mathbb{P}(|\Gamma'_{h:t}(\hat{T}_h - T_h)(s_{h-1}, a_{h-1})| > \epsilon | E_t) \\ & \leq \mathbb{E} \left[ \exp\left(-\frac{2n_{s_{h-1}, a_{h-1}} \cdot \epsilon^2}{1}\right) \middle| E_t \right] \leq \exp\left(-\frac{nd_{h-1}^\mu(s_{h-1}, a_{h-1}) \cdot \epsilon^2}{1}\right) \end{aligned}$$

and apply a union bound to get

$$\begin{aligned} & P(\sup_h \|\Gamma'_{h:t}(\tilde{T}_h - T_h)\|_\infty > \epsilon) \leq H \cdot \sup_h P(\|\Gamma'_{h:t}(\tilde{T}_h - T_h)\|_\infty > \epsilon) \\ & \leq HSA \cdot \sup_{h, s_{h-1}, a_{h-1}} \mathbb{P}(|\Gamma'_{h:t}(\tilde{T}_h - T_h)(s_{h-1}, a_{h-1})| > \epsilon) \\ & \leq HSA \cdot \exp\left(-\frac{n \min d_{h-1}^\mu(s_{h-1}, a_{h-1}) \cdot \epsilon^2}{1}\right) := \frac{\delta}{HSA} \end{aligned} \tag{10}$$

Let the right hand side of (10) to be  $\delta/HSA$ , then we have w.p.  $1 - \delta/HSA$ ,

$$\sup_h \|\Gamma'_{h:t}(\tilde{T}_h - T_h)\|_\infty \leq O\left(\sqrt{\frac{1}{n \cdot d_m} \log \frac{H^2 S^2 A^2}{\delta}}\right). \tag{11}$$

**Go back to bounding  $\tilde{d}_t^\pi(s_t, a_t) - d_t^\pi(s_t, a_t)$ .** By Azuma-Hoeffding's inequality (Lemma A.5), we have<sup>11</sup>

$$\mathbb{P}(|\tilde{d}_t^\pi(s_t, a_t) - d_t^\pi(s_t, a_t)| > \epsilon) \leq \exp\left(-\frac{\epsilon^2}{\sum_{i=1}^t (\sup_h \|\Gamma'_{h:t}(\tilde{T}_h - T_h)\|_\infty)^2}\right) := \delta/HSA,$$

where  $\sum_{i=1}^t (\sup_h \|\Gamma'_{h:t}(\tilde{T}_h - T_h)\|_\infty)^2$  is the sum of bounded square differences in Azuma-Hoeffding's inequality. Therefore we have w.p.  $1 - \delta/HSA$ ,

$$|\tilde{d}_t^\pi(s_t, a_t) - d_t^\pi(s_t, a_t)| \leq O\left(\sqrt{t \cdot (\sup_h \|\Gamma'_{h:t}(\tilde{T}_h - T_h)\|_\infty)^2 \log \frac{HSA}{\delta}}\right), \tag{12}$$

<sup>11</sup>To be more precise here we actually use a weaker version of Azuma-Hoeffding's inequality, see Remark E.7.

combining (11) with above we further have that w.p.  $1 - 2\delta/HSA$ ,

$$|\tilde{d}_t^\pi(s_t, a_t) - d_t^\pi(s_t, a_t)| \leq O\left(\sqrt{\frac{t}{nd_m} \log \frac{H^2 S^2 A^2}{\delta} \log \frac{HSA}{\delta}}\right)$$

**Lastly**, by a union bound and simple scaling (from  $2\delta$  to  $\delta$ ) we have w.p.  $1 - \delta$

$$\sup_t \|\tilde{d}_t^\pi - d_t^\pi\|_\infty \leq O\left(\sqrt{\frac{H}{nd_m} \log \frac{H^2 S^2 A^2}{\delta} \log \frac{HSA}{\delta}}\right).$$

This implies that w.p.  $1 - \delta$ ,  $\forall t, s_t, a_t$ ,

$$\tilde{d}_t^\pi(s_t, a_t)^2 \leq 2d_t^\pi(s_t, a_t)^2 + O\left(\frac{H}{nd_m} \log \frac{H^2 S^2 A^2}{\delta} \log \frac{HSA}{\delta}\right). \quad (13)$$

Combining (13) with Lemma E.4 and Lemma E.3, we obtain:

**Lemma E.5.** *With probability  $1 - \delta$ ,*

$$\sum_{t=1}^H \text{Var}[X_{t+1}|\mathcal{D}_t] \leq O\left(\frac{H^2}{nd_m}\right) + O\left(\frac{H^4 SA}{n^2 d_m^2} \cdot \log \frac{H^2 S^2 A^2}{\delta} \log \frac{HSA}{\delta}\right) \quad (14)$$

*Proof of Lemma E.5.* By (13) and Lemma E.3, we have  $\forall t$ , with probability at least  $1 - \delta$ ,

$$\begin{aligned} \text{Var}[X_{t+1}|\mathcal{D}_t] &\leq \sum_{s_t, a_t} O\left(\frac{\tilde{d}_t^\pi(s_t, a_t)^2}{nd_m}\right) \cdot \text{Var}[V_{t+1}^\pi(s_{t+1}^{(1)})|s_t^{(1)} = s_t, a_t^{(1)} = a_t] \\ &\leq \sum_{s_t, a_t} O\left(\frac{1}{nd_m}\right) \left(2d_t^\pi(s_t, a_t)^2 + O\left(\frac{H}{nd_m} \log \frac{H^2 S^2 A^2}{\delta} \log \frac{HSA}{\delta}\right)\right) \cdot \text{Var}[V_{t+1}^\pi(s_{t+1}^{(1)})|s_t^{(1)} = s_t, a_t^{(1)} = a_t] \\ &\leq \sum_{s_t, a_t} O\left(\frac{1}{nd_m}\right) \left(2d_t^\pi(s_t, a_t) + O\left(\frac{H}{nd_m} \log \frac{H^2 S^2 A^2}{\delta} \log \frac{HSA}{\delta}\right)\right) \cdot \text{Var}[V_{t+1}^\pi(s_{t+1}^{(1)})|s_t^{(1)} = s_t, a_t^{(1)} = a_t] \\ &\leq O\left(\frac{1}{nd_m}\right) \mathbb{E} \left[ \text{Var}[V_{t+1}^\pi(s_{t+1}^{(1)})|s_t^{(1)}, a_t^{(1)}] \right] + O\left(\frac{1}{nd_m} \cdot \frac{H}{nd_m} \log \frac{H^2 S^2 A^2}{\delta} \log \frac{HSA}{\delta} \cdot H^2 SA\right) \\ &= O\left(\frac{1}{nd_m}\right) \mathbb{E} \left[ \text{Var}[V_{t+1}^\pi(s_{t+1}^{(1)})|s_t^{(1)}, a_t^{(1)}] \right] + O\left(\frac{H^3 SA}{n^2 d_m^2} \cdot \log \frac{H^2 S^2 A^2}{\delta} \log \frac{HSA}{\delta}\right) \end{aligned}$$

then sum over  $t$  and apply Lemma E.4 gives the stated result.  $\square$

Combining all the results, we are able to prove:

**Theorem E.6.** *With probability  $1 - \delta$ , we have*

$$\left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right| \leq O\left(\sqrt{\frac{H^2 \log(HSA/\delta)}{nd_m}}\right) + \sqrt{\frac{H^4 SA \cdot \log(H^2 S^2 A^2/\delta) \log(HSA/\delta)}{n^2 d_m^2}}$$

where  $O(\cdot)$  absorbs only the absolute constants.

*Proof of Theorem E.6.* Recall  $X = \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle$  and by law of total expectation it is easy to show  $E[X] = 0$ . Next denote  $\sigma^2 = O\left(\frac{H^2}{nd_m}\right) + O\left(\frac{H^4 SA}{n^2 d_m^2} \cdot \log \frac{H^2 S^2 A^2}{\delta} \log \frac{HSA}{\delta}\right)$  as in Lemma E.5 and also let  $M = \sup_t |X_t - X_{t-1}|$ . Then by Freedman inequality (Lemma A.6), we have with probability  $1 - \delta/3$ ,

$$|X - \mathbb{E}[X]| \leq \sqrt{8\sigma^2 \cdot \log(3/\delta)} + \frac{2M}{3} \cdot \log(3/\delta), \quad \text{Or } W \geq \sigma^2.$$

where  $W = \sum_{t=1}^H \text{Var}[X_{t+1}|\mathcal{D}_t]$ . Next by Lemma E.5, we have  $\mathbb{P}(W \geq \sigma^2) \leq 1/3\delta$ , this implies with probability  $1 - 2\delta/3$ ,

$$|X - \mathbb{E}[X]| \leq \sqrt{8\sigma^2 \cdot \log(3/\delta)} + \frac{2M}{3} \cdot \log(3/\delta).$$

Finally, by Lemma E.2, we have  $\mathbb{P}(M \geq O(\sqrt{\frac{H^2 \log(HSA/\delta)}{nd_m}})) \leq \delta/3$ . Also use  $\mathbb{E}[X] = 0$ , we have with probability  $1 - \delta$ ,

$$|X| \leq \sqrt{8\sigma^2 \cdot \log(3/\delta)} + O(\sqrt{\frac{H^2 \cdot \log(HSA/\delta)}{nd_m}} \log(3/\delta)).$$

Plugging back the expression of  $\sigma^2 = O(\frac{H^2}{nd_m}) + O(\frac{H^4 SA}{n^2 d_m^2} \cdot \log \frac{H^2 S^2 A^2}{\delta} \log \frac{HSA}{\delta})$  and assimilating the same order terms give the desired result.  $\square$

**Remark E.7.** Rigorously, standard Azuma-Hoeffding's inequality Lemma A.5 does not apply to (12) since  $\sup_h \|\Gamma'_{h:t}(\tilde{T}_h - T_h)\|_\infty$  is not a deterministic upper bound, we only have the difference bound with high probability sense, see (11). Therefore, strictly speaking, we need to apply Theorem 32 in Chung and Lu (2006) which is a weaker Azuma-Hoeffding's inequality allowing bounded difference with high probability. The same logic applies for a weaker freedman's inequality consisting of Theorem 34 and Theorem 37 in Chung and Lu (2006) since our martingale difference  $M = \sup_t |X_t - X_{t-1}|$  in the proof of Theorem E.6 is bounded with high probability. We avoid explicitly using them in order to make our proofs more readable for our readers.

We end this section by giving the proofs of Theorem 3.4 and Theorem 3.5.

*Proof of Lemma 3.4 and Theorem 3.5.* The proof of Lemma 3.4 comes from Lemma C.3, Lemma C.4 and Theorem E.6. The proof of Theorem 3.5 relies on applying a union bound over  $\Pi$  in Theorem 3.4 (recall all non-stationary deterministic policies have  $|\Pi| = A^{HS}$ ), then extra dependence of  $\sqrt{\log(|\Pi|)} = \sqrt{HS \log(A)}$  pops out. Note that the higher order term has two trailing log terms (see the right hand side of (14)), so when replacing  $\delta$  by  $\delta/|\Pi|$  with a union bound, both terms will give extra  $\sqrt{HS}$  dependence so in higher order term we have extra  $HS$  dependence but not just  $\sqrt{HS}$ .  $\square$

## F Proof of uniform convergence problem with local policy class.

In this section, we consider using OPEMA estimator with data  $\mathcal{D}'$ . Also, WLOG we only consider deterministic reward (as implied by Lemma C.4 random reward only causes lower order dependence). Also, we fix  $N > 0$  for the moment. First recall for all  $t = 1, \dots, H$

$$\begin{aligned} V_t^\pi(s_t) &= \mathbb{E}_\pi \left[ \sum_{t'=t}^H r_{t'}(s_{t'}^{(1)}, a_{t'}^{(1)}) \middle| s_t^{(1)} = s_t \right] \\ Q_t^\pi(s_t, a_t) &= \mathbb{E}_\pi \left[ \sum_{t'=t}^H r_{t'}(s_{t'}^{(1)}, a_{t'}^{(1)}) \middle| s_t^{(1)} = s_t, a_t^{(1)} = a_t \right] \end{aligned}$$

where  $r_t(s, a)$  are deterministic rewards and  $s_t^{(1)}, a_t^{(1)}$  are random variables. Consider  $V_t^\pi, Q_t^\pi$  as vectors, then by standard Bellman equations we have for all  $t = 1, \dots, H$  (define  $V_{H+1} = Q_{H+1} = 0$ )

$$Q_t^\pi = r_t + P_{t+1}^\pi Q_{t+1}^\pi = r_t + P_{t+1} V_{t+1}^\pi, \quad (15)$$

where  $P_t^\pi \in \mathbb{R}^{(SA) \times (SA)}$  is the state-action transition and  $P_t(\cdot, \cdot) \in \mathbb{R}^{(SA) \times S}$  is the transition probabilities defined in Section 2. Also, we have bellman optimality equations:

$$Q_t^* = r_t + P_{t+1} V_{t+1}^*, \quad V_t^*(s_t) := \max_{a_t} Q_t^*(s_t, a_t), \quad \pi_t^*(s_t) := \operatorname{argmax}_{a_t} Q_t^*(s_t, a_t) \quad \forall s_t \quad (16)$$

where  $\pi^*$  is one optimal deterministic policy. The corresponding Bellman equations and Bellman optimality equations for empirical MDP  $\widehat{M}$  are defined similarly. Since we consider deterministic rewards, by Bellman equations we have

$$\widehat{Q}_t^\pi - Q_t^\pi = \widehat{P}_{t+1}^\pi \widehat{Q}_{t+1}^\pi - P_{t+1}^\pi Q_{t+1}^\pi = (\widehat{P}_{t+1}^\pi - P_{t+1}^\pi) \widehat{Q}_{t+1}^\pi + P_{t+1}^\pi (\widehat{Q}_{t+1}^\pi - Q_{t+1}^\pi)$$

for  $t = 1, \dots, H$ . By writing it recursively, we have  $\forall t = 1, \dots, H - 1$

$$\begin{aligned} \widehat{Q}_t^\pi - Q_t^\pi &= \sum_{h=t+1}^H \Gamma_{t+1:h-1}^\pi (\widehat{P}_h^\pi - P_h^\pi) \widehat{Q}_h^\pi \\ &= \sum_{h=t+1}^H \Gamma_{t+1:h-1}^\pi (\widehat{P}_h - P_h) \widehat{V}_h^\pi \end{aligned}$$

where  $\Gamma_{t:h}^\pi = \prod_{i=t}^h P_i^\pi$  is the multi-step state-action transition and  $\Gamma_{t+1:t}^\pi := I$ .

Note  $\widehat{\pi}^*$  to be the empirical optimal policy over  $\widehat{M}$ , we are interested in how to obtain uniform convergence for any policy  $\pi$  that is close to  $\widehat{\pi}^*$ . More precisely, in this section we consider the policy class  $\Pi_1$  to be:

$$\Pi_1 := \{\pi : s.t. \|\widehat{V}_t^\pi - \widehat{V}_t^{\widehat{\pi}^*}\|_\infty \leq \epsilon_{\text{opt}}, \forall t = 1, \dots, H\}$$

where  $\epsilon_{\text{opt}} \geq 0$  is a parameter decides how large the policy class is. We now assume  $\widehat{\pi}$  to be any policy within  $\Pi_1$  throughout this section. **Also,  $\widehat{\pi}$  may be a policy learned from a learning algorithm using the data  $\mathcal{D}$ . In this case,  $\widehat{\pi}$  may not be independent of  $\widehat{P}$ .**

We start with the following simple calculation:<sup>12</sup>

$$\begin{aligned} \left| \widehat{Q}_t^{\widehat{\pi}} - Q_t^{\widehat{\pi}} \right| &\leq \sum_{h=t+1}^H \Gamma_{t+1:h-1}^\pi \left| (\widehat{P}_h - P_h) \widehat{V}_h^{\widehat{\pi}} \right| \\ &\leq \underbrace{\sum_{h=t+1}^H \Gamma_{t+1:h-1}^\pi \left| (\widehat{P}_h - P_h) \widehat{V}_h^{\widehat{\pi}^*} \right|}_{(***)} + \underbrace{\sum_{h=t+1}^H \Gamma_{t+1:h-1}^\pi \left| (\widehat{P}_h - P_h) (\widehat{V}_h^{\widehat{\pi}^*} - \widehat{V}_h^{\widehat{\pi}}) \right|}_{(****)} \end{aligned} \quad (17)$$

We now analyze  $(***)$  and  $(****)$ .

### F.1 Analyzing $\sum_{h=t+1}^H \Gamma_{t+1:h-1}^\pi \left| (\widehat{P}_h - P_h) (\widehat{V}_h^{\widehat{\pi}^*} - \widehat{V}_h^{\widehat{\pi}}) \right|$

First, by vector induced matrix norm<sup>13</sup> we have

$$\begin{aligned} \left\| \sum_{h=t+1}^H \Gamma_{t+1:h-1}^\pi \cdot \left| (\widehat{P}_h - P_h) (\widehat{V}_h^{\widehat{\pi}^*} - \widehat{V}_h^{\widehat{\pi}}) \right| \right\|_\infty &\leq H \cdot \sup_h \left\| \Gamma_{t+1:h-1}^\pi \right\|_\infty \left\| |(\widehat{P}_h - P_h) (\widehat{V}_h^{\widehat{\pi}^*} - \widehat{V}_h^{\widehat{\pi}})| \right\|_\infty \\ &\leq H \cdot \sup_h \left\| |(\widehat{P}_h - P_h) (\widehat{V}_h^{\widehat{\pi}^*} - \widehat{V}_h^{\widehat{\pi}})| \right\|_\infty \end{aligned}$$

where the last equal sign uses multi-step transition  $\Gamma_{t+1:h-1}^\pi$  is row-stochastic. Note given  $N$ ,  $\widehat{P}_t(\cdot|\cdot, \cdot)$  all have  $N$  in the denominator. Therefore, by Hoeffding inequality and a union bound we have with probability  $1 - \delta$ ,

$$\sup_{t, s_t, s_{t-1}, a_{t-1}} |\widehat{P}_t(s_t|s_{t-1}, a_{t-1}) - P_t(s_t|s_{t-1}, a_{t-1})| \leq O\left(\sqrt{\frac{\log(HSA/\delta)}{N}}\right),$$

this indicates

$$\sup_h \left\| |(\widehat{P}_h - P_h) (\widehat{V}_h^{\widehat{\pi}^*} - \widehat{V}_h^{\widehat{\pi}})| \right\|_\infty \leq \epsilon_{\text{opt}} \cdot \sup_h \left\| |\widehat{P}_h - P_h| \cdot \mathbf{1} \right\|_\infty \leq \epsilon_{\text{opt}} \cdot O\left(S\sqrt{\frac{\log(HSA/\delta)}{N}}\right),$$

where  $\mathbf{1} \in \mathbb{R}^S$  is all-one vector. To sum up, we have

<sup>12</sup>Since all quantities in the calculation are vectors, so the absolute value  $|\cdot|$  used is point-wise operator.

<sup>13</sup>For  $A$  a matrix and  $x$  a vector we have  $\|Ax\|_\infty \leq \|A\|_\infty \|x\|_\infty$ .

**Lemma F.1.** Fix  $N > 0$ , we have with probability  $1 - \delta$ , for all  $t = 1, \dots, H - 1$

$$\sum_{h=t+1}^H \Gamma_{t+1:h-1}^{\hat{\pi}} \left| (\hat{P}_h - P_h)(\hat{V}_h^{\hat{\pi}^*} - \hat{V}_h^{\hat{\pi}}) \right| \leq \epsilon_{\text{opt}} \cdot O \left( \sqrt{\frac{H^2 S^2 \log(HSA/\delta)}{N}} \cdot \mathbf{1} \right)$$

Now we consider  $(***)$ .

**F.2 Analyzing**  $\sum_{h=t+1}^H \Gamma_{t+1:h-1}^{\hat{\pi}} \left| (\hat{P}_h - P_h) \hat{V}_h^{\hat{\pi}^*} \right|$ .

**Lemma F.2.** Given  $N$ , we have with probability  $1 - \delta$ ,  $\forall t = 1, \dots, H - 1$

$$\sum_{h=t+1}^H \Gamma_{t+1:h-1}^{\hat{\pi}} \left| (\hat{P}_h - P_h) \hat{V}_h^{\hat{\pi}^*} \right| \leq \sum_{h=t+1}^H \Gamma_{t+1:h-1}^{\hat{\pi}} \left( 4 \sqrt{\frac{\log(HSA/\delta)}{N}} \sqrt{\text{Var}(\hat{V}_h^{\hat{\pi}^*})} + \frac{4(H-t)}{3N} \log\left(\frac{HSA}{\delta}\right) \cdot \mathbf{1} \right)$$

where  $\text{Var}(v_t^\pi) \in \mathbb{R}^{SA}$  and  $\text{Var}(V_t^\pi)(s_{t-1}, a_{t-1}) = \text{Var}_{s_t}[V_t^\pi(\cdot) | s_{t-1}, a_{t-1}]$  and  $|\cdot|, \sqrt{\cdot}$  are point-wise operator.

*Proof of Lemma F.2.* The key point is to guarantee  $\hat{P}_h$  is independent of  $\hat{V}_h^{\hat{\pi}^*}$  so that we can apply Bernstein inequality w.r.t the randomness in  $\hat{P}_h$ . In fact, note given  $N$  all data pairs in  $\mathcal{D}'$  are independent of each other, and  $\hat{P}_h$  only uses data from  $h-1$  to  $h$ . Moreover,  $\hat{V}_h^{\hat{\pi}^*}$  only uses data from time  $h$  to  $H$  since  $\hat{V}_h^\pi$  uses data from  $h$  to  $H$  by bellman equation (15) for any  $\pi$  and optimal policy  $\hat{\pi}_{h:H}^*$  also only uses data from  $h$  to  $H$  by bellman optimality equation (16).

Then by Bernstein inequality (Lemma A.3), with probability  $1 - \delta$

$$\left| (\hat{P}_h - P_h) \hat{V}_h^{\hat{\pi}^*} \right| (s_{t-1}, a_{t-1}) \leq 4 \sqrt{\frac{\log(1/\delta)}{N}} \sqrt{\text{Var}(\hat{V}_h^{\hat{\pi}^*})(s_{t-1}, a_{t-1})} + \frac{4(H-t)}{3N} \log\left(\frac{1}{\delta}\right)$$

apply a union bound and take the sum we get the stated result.  $\square$

Now combine Lemma F.1 and Lemma F.2 we obtain with probability  $1 - \delta$ , for all  $t = 1, \dots, H - 1$

$$\begin{aligned} \left| \hat{Q}_t^{\hat{\pi}} - Q_t^{\hat{\pi}} \right| &\leq \sum_{h=t+1}^H \Gamma_{t+1:h-1}^{\hat{\pi}} \left( 4 \sqrt{\frac{\log(HSA/\delta)}{N}} \sqrt{\text{Var}(\hat{V}_h^{\hat{\pi}^*})} + \frac{4(H-t)}{3N} \log\left(\frac{HSA}{\delta}\right) \cdot \mathbf{1} \right) \\ &\quad + c_1 \epsilon_{\text{opt}} \cdot \sqrt{\frac{H^2 S^2 \log(HSA/\delta)}{N}} \cdot \mathbf{1} \\ &\leq 4 \sqrt{\frac{\log(HSA/\delta)}{N}} \sum_{h=t+1}^H \Gamma_{t+1:h-1}^{\hat{\pi}} \sqrt{\text{Var}(\hat{V}_h^{\hat{\pi}^*})} + \frac{4H^2}{3N} \log\left(\frac{HSA}{\delta}\right) \cdot \mathbf{1} \\ &\quad + c_1 \epsilon_{\text{opt}} \cdot \sqrt{\frac{H^2 S^2 \log(HSA/\delta)}{N}} \cdot \mathbf{1}, \end{aligned} \tag{18}$$

Next note  $\sqrt{\text{Var}(\cdot)}$  is a norm, therefore by norm triangle inequality we have

$$\begin{aligned} \sqrt{\text{Var}(\hat{V}_h^{\hat{\pi}^*})} &\leq \sqrt{\text{Var}(\hat{V}_h^{\hat{\pi}^*} - \hat{V}_h^{\hat{\pi}})} + \sqrt{\text{Var}(\hat{V}_h^{\hat{\pi}} - V_h^{\hat{\pi}})} + \sqrt{\text{Var}(V_h^{\hat{\pi}})} \\ &\leq \left\| \hat{V}_h^{\hat{\pi}^*} - \hat{V}_h^{\hat{\pi}} \right\|_\infty \cdot \mathbf{1} + \left\| \hat{V}_h^{\hat{\pi}} - V_h^{\hat{\pi}} \right\|_\infty \cdot \mathbf{1} + \sqrt{\text{Var}(V_h^{\hat{\pi}})} \\ &\leq \epsilon_{\text{opt}} \cdot \mathbf{1} + \left\| \hat{Q}_h^{\hat{\pi}} - Q_h^{\hat{\pi}} \right\|_\infty \cdot \mathbf{1} + \sqrt{\text{Var}(V_h^{\hat{\pi}})} \end{aligned}$$



Plug this into (18) to obtain

$$\begin{aligned} \left| \widehat{Q}_t^{\widehat{\pi}} - Q_t^{\widehat{\pi}} \right| &\leq 4 \sqrt{\frac{\log(HSA/\delta)}{N}} \sum_{h=t+1}^H \left( \Gamma_{t+1:h-1}^{\widehat{\pi}} \sqrt{\text{Var}(V_h^{\widehat{\pi}})} + \left\| \widehat{Q}_h^{\widehat{\pi}} - Q_h^{\widehat{\pi}} \right\|_{\infty} \cdot \mathbf{1} \right) + \frac{4H^2}{3N} \log\left(\frac{HSA}{\delta}\right) \cdot \mathbf{1} \\ &\quad + c_2 \epsilon_{\text{opt}} \cdot \sqrt{\frac{H^2 S^2 \log(HSA/\delta)}{N}} \cdot \mathbf{1}. \end{aligned} \quad (19)$$

Next lemma helps us to bound  $\sum_{h=t+1}^H \Gamma_{t+1:h-1}^{\widehat{\pi}} \sqrt{\text{Var}(V_h^{\widehat{\pi}})}$ .

**Lemma F.3.** *A conditional version of Lemma E.4 holds:*

$$\begin{aligned} \text{Var}_{\pi} \left[ \sum_{t=h}^H r_t^{(1)} \middle| s_h^{(1)} = s_h, a_h^{(1)} = a_h \right] &= \sum_{t=h}^H \left( \mathbb{E}_{\pi} \left[ \text{Var} \left[ r_t^{(1)} + V_{t+1}^{\pi}(s_{t+1}^{(1)}) \middle| s_t^{(1)}, a_t^{(1)} \right] \middle| s_h^{(1)} = s_h, a_h^{(1)} = a_h \right] \right. \\ &\quad \left. + \mathbb{E}_{\pi} \left[ \text{Var} \left[ \mathbb{E}[r_t^{(1)} + V_{t+1}^{\pi}(s_{t+1}^{(1)}) \middle| s_t^{(1)}, a_t^{(1)}] \middle| s_t^{(1)} \right] \middle| s_h^{(1)} = s_h, a_h^{(1)} = a_h \right] \right). \end{aligned} \quad (20)$$

and by using (20) we can show

$$\sum_{h=t+1}^H \Gamma_{t+1:h-1}^{\widehat{\pi}} \sqrt{\text{Var}(V_h^{\widehat{\pi}})} \leq \sqrt{(H-t)^3} \cdot \mathbf{1}.$$

*Proof.* The proof of (20) uses the identical trick as Lemma E.4 except the total law of variance is replaced by the total law of conditional variance.

Moreover, recall  $\Gamma_{t+1:h-1}^{\widehat{\pi}} = \prod_{i=t+1}^{h-1} P_i^{\widehat{\pi}}$  is the multi-step transition, so for any pair  $(s_t, a_t)$ ,

$$\begin{aligned} &\sum_{h=t+1}^H \left( \Gamma_{t+1:h-1}^{\widehat{\pi}} \sqrt{\text{Var}(V_h^{\widehat{\pi}})} \right) (s_t, a_t) \\ &= \sum_{h=t+1}^H \sum_{s_{h-1}, a_{h-1}} \sqrt{\text{Var}[V_h^{\widehat{\pi}} | s_{h-1}, a_{h-1}] d_t^{\widehat{\pi}}(s_{h-1}, a_{h-1} | s_t, a_t)} \\ &= \sum_{h=t+1}^H \sum_{s_{h-1}, a_{h-1}} \sqrt{\text{Var}[V_h^{\widehat{\pi}} | s_{h-1}, a_{h-1}] d_t^{\widehat{\pi}}(s_{h-1}, a_{h-1} | s_t, a_t)} \cdot \sqrt{d_t^{\widehat{\pi}}(s_{h-1}, a_{h-1} | s_t, a_t)} \\ &\leq \sum_{h=t+1}^H \sqrt{\sum_{s_{h-1}, a_{h-1}} \text{Var}[V_h^{\widehat{\pi}} | s_{h-1}, a_{h-1}] d_t^{\widehat{\pi}}(s_{h-1}, a_{h-1} | s_t, a_t)} \cdot \sum_{s_{h-1}, a_{h-1}} d_t^{\widehat{\pi}}(s_{h-1}, a_{h-1} | s_t, a_t) \\ &= \sum_{h=t+1}^H \sqrt{\sum_{s_{h-1}, a_{h-1}} \text{Var}[V_h^{\widehat{\pi}} | s_{h-1}, a_{h-1}] d_t^{\widehat{\pi}}(s_{h-1}, a_{h-1} | s_t, a_t)} \\ &= \sum_{h=t+1}^H \sqrt{\mathbb{E}_{\widehat{\pi}} \left[ \text{Var}[V_h^{\widehat{\pi}} | s_{h-1}^{(1)}, a_{h-1}^{(1)}] \middle| s_t, a_t \right]} \\ &= \sum_{h=t+1}^H \sqrt{1} \cdot \sqrt{\mathbb{E}_{\widehat{\pi}} \left[ \text{Var}[V_h^{\widehat{\pi}} | s_{h-1}^{(1)}, a_{h-1}^{(1)}] \middle| s_t, a_t \right]} \\ &\leq \sqrt{(H-t) \sum_{h=t+1}^H \mathbb{E}_{\widehat{\pi}} \left[ \text{Var}[V_h^{\widehat{\pi}} | s_{h-1}^{(1)}, a_{h-1}^{(1)}] \middle| s_t, a_t \right]} \\ &\leq \sqrt{(H-t) \cdot \text{Var}_{\widehat{\pi}} \left[ \sum_{h=t+1}^H r_h^{(1)} \middle| s_t^{(1)} = s_t, a_t^{(1)} = a_t \right]} \leq \sqrt{(H-t)^3} \end{aligned}$$

where all the inequalities are Cauchy-Schwarz inequalities.  $\square$

Apply Lemma F.3 to bound (19), and use  $\infty$  norm on both sides, we obtain

**Theorem F.4.** *Conditional on  $N > 0$ , then with probability  $1 - \delta$ , we have for all  $t = 1, \dots, H - 1$*

$$\begin{aligned} \|\hat{Q}_t^{\hat{\pi}} - Q_t^{\hat{\pi}}\|_{\infty} &\leq 4\sqrt{\frac{H^3 \log(HSA/\delta)}{N}} + 4\sqrt{\frac{\log(HSA/\delta)}{N}} \sum_{h=t+1}^H \|\hat{Q}_h^{\hat{\pi}} - Q_h^{\hat{\pi}}\|_{\infty} + \frac{4H^2}{3N} \log\left(\frac{HSA}{\delta}\right) \\ &\quad + c_2 \epsilon_{opt} \cdot \sqrt{\frac{H^2 S^2 \log(HSA/\delta)}{N}}. \end{aligned}$$

Then by using backward induction and Theorem F.4, we have the following:

**Theorem F.5.** *Suppose  $N \geq 64H^2 \cdot \log(HSA/\delta)$  and  $\epsilon_{opt} \leq \sqrt{H}/S$ , then we have with probability  $1 - \delta$ ,*

$$\|\hat{Q}_1^{\hat{\pi}} - Q_1^{\hat{\pi}}\|_{\infty} \leq 2(9 + c_2) \sqrt{\frac{H^3 \log(HSA/\delta)}{N}}$$

where  $c_2$  is the same constant in Theorem F.4.

*Proof.* Under the condition, by Theorem F.4 it is easy to check for all  $t = 1, \dots, H - 1$  with probability  $1 - \delta$ ,

$$\|\hat{Q}_t^{\hat{\pi}} - Q_t^{\hat{\pi}}\|_{\infty} \leq (5 + c_2) \sqrt{\frac{H^3 \log(HSA/\delta)}{N}} + 4\sqrt{\frac{\log(HSA/\delta)}{N}} \sum_{h=t+1}^H \|\hat{Q}_h^{\hat{\pi}} - Q_h^{\hat{\pi}}\|_{\infty},$$

which we conditional on.

For  $t = H - 1$ , we have

$$\begin{aligned} \|\hat{Q}_{H-1}^{\hat{\pi}} - Q_{H-1}^{\hat{\pi}}\|_{\infty} &\leq (5 + c_2) \sqrt{\frac{H^3 \log(HSA/\delta)}{N}} + 4\sqrt{\frac{\log(HSA/\delta)}{N}} \|\hat{Q}_H^{\hat{\pi}} - Q_H^{\hat{\pi}}\|_{\infty} \\ &\leq (5 + c_2) \sqrt{\frac{H^3 \log(HSA/\delta)}{N}} + 4\sqrt{\frac{H^2 \log(HSA/\delta)}{N}} \\ &\leq (9 + c_2) \sqrt{\frac{H^3 \log(HSA/\delta)}{N}} \end{aligned}$$

Suppose  $\|\hat{Q}_h^{\hat{\pi}} - Q_h^{\hat{\pi}}\|_{\infty} \leq 2(9 + c_2) \sqrt{\frac{H^3 \log(HSA/\delta)}{N}}$  holds for all  $h = t + 1, \dots, H$ , then for  $h = t$ , we have

$$\begin{aligned} \|\hat{Q}_t^{\hat{\pi}} - Q_t^{\hat{\pi}}\|_{\infty} &\leq (5 + c_2) \sqrt{\frac{H^3 \log(HSA/\delta)}{N}} + 4\sqrt{\frac{\log(HSA/\delta)}{N}} \sum_{h=t+1}^H \|\hat{Q}_h^{\hat{\pi}} - Q_h^{\hat{\pi}}\|_{\infty} \\ &\leq (9 + c_2) \sqrt{\frac{H^3 \log(HSA/\delta)}{N}} + 4\sqrt{\frac{(H-1)^2 \log(HSA/\delta)}{N}} \cdot 2(9 + c_2) \sqrt{\frac{H^3 \log(HSA/\delta)}{N}} \\ &\leq 2(9 + c_2) \sqrt{\frac{H^3 \log(HSA/\delta)}{N}} \end{aligned}$$

where the last line uses the condition  $N \geq 64H^2 \cdot \log(HSA/\delta)$ . By induction, we have the result.  $\square$

*Proof of Theorem 3.7.* By Theorem F.5 we have for  $N \geq c \cdot H^2 \cdot \log(HSA/\delta)$ ,

$$\mathbb{P}\left(\left\|\hat{Q}_1^{\hat{\pi}} - Q_1^{\hat{\pi}}\right\|_{\infty} \geq 2(9 + c_2) \sqrt{\frac{H^3 \log(HSA/\delta)}{N}} \mid N\right) \leq \delta$$

The only thing left is to use Lemma C.1 to bound the event that  $\{N < nd_m/2\}$  has small probability.

Last but not least, the condition  $n > c_1 H^2 \log(HSA/\delta)/d_m$  is sufficient for applying Lemma C.1 and it also implies  $N \geq c \cdot H^2 \cdot \log(HSA/\delta)$  (the condition of Theorem F.5) when  $N \geq nd_m/2$  since:

$$n > c_1 H^2 \log(HSA/\delta)/d_m \Rightarrow nd_m/2 \geq c_2 H^2 \log(HSA/\delta)$$

which implies  $N \geq c_2 \cdot H^2 \cdot \log(HSA/\delta)$  when  $N \geq nd_m/2$ .  $\square$

## G Proof of uniform convergence lower bound.

In this section we prove a uniform convergence OPE lower bound of  $\Omega(H^3/d_m\epsilon^2)$ . Conceptually, uniform convergence lower bound can be derived by a reduction to the lower bound of identifying the  $\epsilon$ -optimal policy. There are quite a few literature that provide information theoretical lower bounds in different setting, *e.g.* Dann and Brunskill (2015); Jiang et al. (2017); Krishnamurthy et al. (2016); Jin et al. (2018); Sidford et al. (2018). However, to the best of our knowledge, there is no result proven for the non-stationary transition finite horizon episodic setting with bounded rewards. For example, Sidford et al. (2018) prove the result sample complexity lower bound of  $\Omega(H^3SA/\epsilon^2)$  with stationary MDP and their proof cannot be directly applied to non-stationary setting as they reduce the problem to infinite horizon discounted setting which always has stationary transitions. Dann and Brunskill (2015) prove the episode complexity of  $\tilde{\Omega}(H^2SA/\epsilon^2)$  for the stationary transition setting. Jin et al. (2018) prove the  $\Omega(\sqrt{H^2SAT})$  regret lower bound for non-stationary finite horizon online setting but it is not clear how to translate the regret to PAC-learning setting by keeping the same sample complexity optimality. Jiang et al. (2017) prove the  $\Omega(HSA/\epsilon^2)$  lower bound for the non-stationary finite horizon offline episodic setting where they assume  $\sum_{i=1}^H r_i \leq 1$  and this is also different from our setting since we have  $0 \leq r_t \leq 1$  for each time step.

Our proof consists of three steps. **1.** We will first show a minimax lower bound (**over all MDP instances**) for learning  $\epsilon$ -optimal policy is  $\Omega(H^3SA/\epsilon^2)$ ; **2.** Based on 1, we can further show a minimax lower bound (**over problem class  $\mathcal{M}_{d_m}$** ) for learning  $\epsilon$ -optimal policy is  $\Omega(H^3/d_m\epsilon^2)$ ; **3.** prove the uniform convergence OPE lower bound of the same rate.

### G.1 Information theoretical lower sample complexity bound over all MDP instances for identifying $\epsilon$ -optimal policy.

In fact, a modified construction of Theorem 5 in Jiang et al. (2017) is our tool for obtaining  $\Omega(H^3SA/\epsilon^2)$  lower bound. We can get the additional  $H^2$  factor by using  $\sum_{i=1}^H r_i$  can be of order  $O(H)$ .

**Theorem G.1.** *Given  $H \geq 2$ ,  $A \geq 2$ ,  $0 < \epsilon < \frac{1}{48\sqrt{8}}$  and  $S \geq c_1$  where  $c_1$  is a universal constant. Then there exists another universal constant  $c$  such that for any algorithm and any  $n \leq cH^3SA/\epsilon^2$ , there exists a non-stationary  $H$  horizon MDP with probability at least  $1/12$ , the algorithm outputs a policy  $\hat{\pi}$  with  $v^* - v^{\hat{\pi}} \geq \epsilon$ .*

Like in Jiang et al. (2017), the proof relies on embedding  $\Theta(HS)$  independent multi-arm bandit problems into a hard-to-learn MDP so that any algorithm that wants to output a near-optimal policy needs to identify the best action in  $\Omega(HS)$  problems. However, in our construction we make a further modification of Jiang et al. (2017) so that there is **no** waiting states, which is crucial for the reduction from offline family. We also double the length of the hard-to-learn MDP instance so that the latter half uses a “naive” copy construction which is uninformative. The uninformative extension will help to produce the additional  $H^2$  factor.

*Proof of Theorem G.2.* We construct a non-stationary MDP with  $S$  states per level,  $A$  actions per state and has horizon  $2H$ . At each time step, states are categorized into four types with two special states  $g_h, b_h$  and the remaining  $S - 2$  “bandit” states denoted by  $s_{h,i}$ ,  $i \in [S - 2]$ . Each bandit state has an unknown best action  $a_{h,i}^*$  that provides the highest expected reward comparing to other actions.

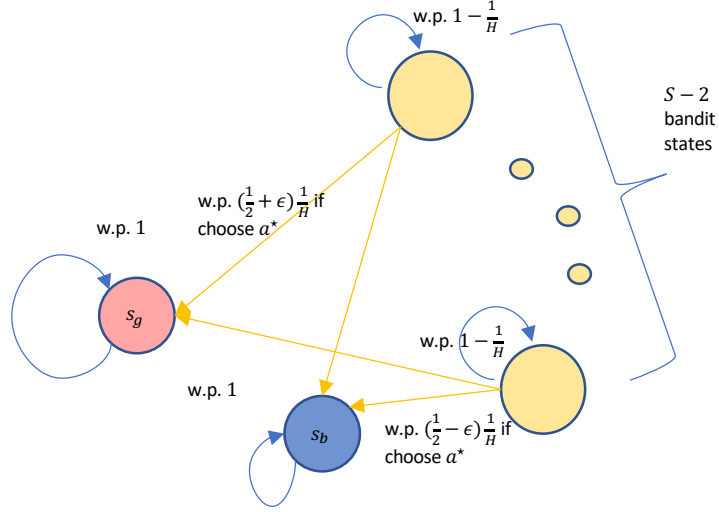


Figure 2: An illustration of the state-space transition diagram from our lower bound construction in Theorem G.2. The new construction eliminates the waiting states, thus making it possible to obtain a lower bound that explicitly depends on parameter  $d_m$  in Theorem G.2.

The transition dynamics are defined as follows:

- for  $h = 1, \dots, H - 1$ ,
  - For bandit states  $b_{h,i}$ , there is probability  $1 - \frac{1}{H}$  to transition to  $b_{h+1,i}$  regardless of the action chosen. For the rest of  $\frac{1}{H}$  probability, optimal action  $a_{h,i}^*$  will have probability  $\frac{1}{2} + \tau$  or  $\frac{1}{2} - \tau$  transition to  $g_{h+1}$  or  $b_{h+1}$  and all other actions  $a$  will have equal probability  $\frac{1}{2}$  for either  $g_{h+1}$  or  $b_{h+1}$ , where  $\tau$  is a parameter will be decided later. Or equivalently,

$$\mathbb{P}(\cdot | s_{h,i}, a_{h,i}^*) = \begin{cases} 1 - \frac{1}{H} & \text{if } \cdot = s_{h+1,i} \\ (\frac{1}{2} + \tau) \cdot \frac{1}{H} & \text{if } \cdot = g_{h+1} \\ (\frac{1}{2} - \tau) \cdot \frac{1}{H} & \text{if } \cdot = b_{h+1} \end{cases} \quad \mathbb{P}(\cdot | s_{h,i}, a) = \begin{cases} 1 - \frac{1}{H} & \text{if } \cdot = s_{h+1,i} \\ \frac{1}{2} \cdot \frac{1}{H} & \text{if } \cdot = g_{h+1} \\ \frac{1}{2} \cdot \frac{1}{H} & \text{if } \cdot = b_{h+1} \end{cases}$$

- $g_h$  always transitions to  $g_{h+1}$  and  $b_h$  always transitions to  $b_{h+1}$ , *i.e.* for all  $a \in \mathcal{A}$ , we have

$$\mathbb{P}(g_{h+1} | g_h, a) = 1, \quad \mathbb{P}(b_{h+1} | b_h, a) = 1.$$

We will determine parameter  $\tau$  at the end of the proof.

- for  $h = H, \dots, 2H - 1$ , all states will always transition to the same type of states for the next step, *i.e.*  $\forall a \in \mathcal{A}$ ,

$$\mathbb{P}(g_{h+1} | g_h, a) = \mathbb{P}(b_{h+1} | b_h, a) = \mathbb{P}(s_{h+1,i} | s_{h,i}, a) = 1, \quad \forall i \in [S - 2]. \quad (21)$$

- The initial distribution is decided by:

$$\mathbb{P}(s_{1,i}) = \frac{1}{S}, \quad \forall i \in [S - 2], \quad \mathbb{P}(g_1) = \frac{1}{S}, \quad \mathbb{P}(b_1) = \frac{1}{S} \quad (22)$$

- State  $s$  will receives reward 1 if and only if  $s = g_h$  and  $h \geq H$ . The reward at all other states is zero.

By this construction the optimal policy must take  $a_{h,i}^*$  for each bandit state  $s_{h,i}$  for at least the first half of the MDP, *i.e.* need to take  $a_{h,i}^*$  for  $h \leq H$ . In other words, this construction embeds at least  $H(S - 2)$  independent best arm identification problems that are identical to the stochastic multi-arm bandit problem in Lemma A.7

into the MDP. Note the key innovation here is that we can remove the waiting states used in Jiang et al. (2017) but still keep the multi-arm bandit problem independent!<sup>14</sup>

Notice in our construction, for any bandit state  $s_{h,i}$  with  $h \leq H$ , the difference of the expected reward between optimal action  $a_{h,i}^*$  and other actions is:

$$\begin{aligned}
 & \left(\frac{1}{2} + \tau\right) \cdot \frac{1}{H} \cdot \mathbb{E}[r_{(h+1):2H}|g_{h+1}] + \left(\frac{1}{2} - \tau\right) \cdot \frac{1}{H} \cdot \mathbb{E}[r_{(h+1):2H}|b_{h+1}] + \left(1 - \frac{1}{H}\right) \cdot \mathbb{E}[r_{(h+1):2H}|s_{h+1,i}] \\
 & - \frac{1}{2H} \cdot \mathbb{E}[r_{(h+1):2H}|g_{h+1}] - \frac{1}{2H} \cdot \mathbb{E}[r_{(h+1):2H}|b_{h+1}] - \left(1 - \frac{1}{H}\right) \cdot \mathbb{E}[r_{(h+1):2H}|s_{h+1,i}] \\
 = & \left(\frac{1}{2} + \tau\right) \cdot \frac{1}{H} \cdot \mathbb{E}[r_{(h+1):2H}|g_{h+1}] + \left(\frac{1}{2} - \tau\right) \cdot \frac{1}{H} \cdot \mathbb{E}[r_{(h+1):2H}|b_{h+1}] \\
 & - \frac{1}{2H} \cdot \mathbb{E}[r_{(h+1):2H}|g_{h+1}] - \frac{1}{2H} \cdot \mathbb{E}[r_{(h+1):2H}|b_{h+1}] \\
 = & \left(\frac{1}{2} + \tau\right) \frac{1}{H} \cdot H + \left(\frac{1}{2} - \tau\right) \frac{1}{H} \cdot 0 - \frac{1}{2H} \cdot H + \frac{1}{2H} \cdot 0 = \tau
 \end{aligned} \tag{23}$$

so it seems by Lemma A.7 one suffices to use the least possible  $\frac{A}{72(\tau)^2}$  samples to identify the best action  $a_{h,i}^*$ . However, note the construction of the latter half of the MDP (21) uses mindless reproduction of previous steps and therefore provides no additional information about the best action once the state at time  $H$  is known. In other words, observing  $\sum_{t=1}^{2H} r_t = H$  is equivalent as observing  $\sum_{t=1}^H r_t = 1$ . Therefore, for the bandit states in the first half the samples that provide information for identifying the best arm is up to time  $H$ . As a result, the difference of the expected reward between optimal action  $a_{h,i}^*$  and other action for identifying the best arm should be corrected as:

$$\begin{aligned}
 & \left(\frac{1}{2} + \tau\right) \cdot \frac{1}{H} \cdot \mathbb{E}[r_{(h+1):H}|g_{h+1}] + \left(\frac{1}{2} - \tau\right) \cdot \frac{1}{H} \cdot \mathbb{E}[r_{(h+1):H}|b_{h+1}] + \left(1 - \frac{1}{H}\right) \cdot \mathbb{E}[r_{(h+1):H}|s_{h+1,i}] \\
 & - \frac{1}{2H} \cdot \mathbb{E}[r_{(h+1):H}|g_{h+1}] - \frac{1}{2H} \cdot \mathbb{E}[r_{(h+1):H}|b_{h+1}] - \left(1 - \frac{1}{H}\right) \cdot \mathbb{E}[r_{(h+1):H}|s_{h+1,i}] \\
 = & \left(\frac{1}{2} + \tau\right) \cdot \frac{1}{H} \cdot \mathbb{E}[r_{(h+1):H}|g_{h+1}] + \left(\frac{1}{2} - \tau\right) \cdot \frac{1}{H} \cdot \mathbb{E}[r_{(h+1):H}|b_{h+1}] \\
 & - \frac{1}{2H} \cdot \mathbb{E}[r_{(h+1):H}|g_{h+1}] - \frac{1}{2H} \cdot \mathbb{E}[r_{(h+1):H}|b_{h+1}] \\
 = & \left(\frac{1}{2} + \tau\right) \frac{1}{H} \cdot 1 + \left(\frac{1}{2} - \tau\right) \frac{1}{H} \cdot 0 - \frac{1}{2H} \cdot 1 + \frac{1}{2H} \cdot 0 = \frac{\tau}{H}
 \end{aligned}$$

Now by Lemma A.7, for each bandit state  $s_{h,i}$  satisfying  $h \leq H$ , unless  $\frac{A}{72(\tau/H)^2}$  samples are collected from that state, the learning algorithm fails to identify the optimal action  $a_{h,i}^*$  with probability at least  $1/3$ .

After running any algorithm, let  $C$  be the set of  $(h, s)$  pairs for which the algorithm identifies the correct action. Let  $D$  be the set of  $(h, s)$  pairs for which the algorithm collects fewer than  $\frac{A}{72(\tau/H)^2}$  samples. Then by Lemma A.7 we have

$$\begin{aligned}
 \mathbb{E}[|C|] &= \mathbb{E} \left[ \sum_{(h,s)} \mathbf{1}[a_{h,s} = a_{h,s}^*] \right] \leq ((S-2)H - |D|) + \mathbb{E} \left[ \sum_{(h,s) \in D} \mathbf{1}[a_{h,s} = a_{h,s}^*] \right] \\
 &\leq ((S-2)H - |D|) + \frac{2}{3}|D| = (S-2)H - \frac{1}{3}|D|.
 \end{aligned}$$

If we have  $n \leq \frac{H(S-2)}{2} \times \frac{A}{72(\tau/H)^2}$ , by pigeonhole principle the algorithm can collect  $\frac{A}{72(\tau/H)^2}$  samples for at most half of the bandit problems, i.e.  $|D| \geq H(S-2)/2$ . Therefore we have

$$\mathbb{E}[|C|] \leq (S-2)H - \frac{1}{3}|D| \leq \frac{5}{6}(S-2)H.$$

Then by Markov inequality

$$\mathbb{P} \left[ |C| \geq \frac{11}{12}H(S-2) \right] \leq \frac{5/6}{11/12} = \frac{10}{11}$$

<sup>14</sup>Here independence means solving one bandit problem provides no information on other bandit problems.

so the algorithm failed to identify the optimal action on  $1/12$  fraction of the bandit problems with probability at least  $1/11$ . Note for each failure in identification, the reward is differ by  $\tau$  (see (23)), therefore under the event  $\{|C'| \geq \frac{1}{12}H(S-2)\}$ , following the similar calculation of Jiang et al. (2017) the suboptimality of the policy produced by the algorithm is

$$\begin{aligned} \epsilon &:= v^* - v^{\hat{\pi}} = \mathbb{P}[\text{visit } C'] \times \tau + \mathbb{P}[\text{visit } C] \times 0 = \mathbb{P}\left[\bigcup_{(h,i) \in C'} \text{visit}(h,i)\right] \times \tau \\ &= \sum_{(h,i) \in C'} \mathbb{P}[\text{visit}(h,i)] \times \tau = \sum_{(h,i) \in C'} \frac{1}{HS} (1 - 1/H)^{h-1} \tau \\ &\geq \sum_{(h,i) \in C'} \frac{1}{HS} (1 - 1/H)^H \tau \geq \sum_{(h,i) \in C'} \frac{1}{HS} \frac{1}{4} \tau \\ &\geq \frac{H(S-2)}{12} \frac{1}{HS} \frac{1}{4} \tau = c_1 \frac{\tau}{48}. \end{aligned}$$

where the third equal sign uses all best arm identification problems are independent. Now we set  $\tau = \min(\sqrt{1/8}, 48\epsilon/c_1)$  and under condition  $n \leq cH^3SA/\epsilon^2$ , we have

$$n \leq cH^3SA/\epsilon^2 \leq c48^2H^3SA/\tau^2 = c48^2 \cdot 72HS \cdot \frac{A}{72(\tau/H)^2} := c'HS \cdot \frac{A}{72\tau^2} \leq \frac{H(S-2)}{2} \cdot \frac{A}{72\tau^2},$$

the last inequality holds as long as  $S \geq 2/(1-2c')$ . Therefore in this situation, with probability at least  $1/11$ ,  $v^* - v^{\hat{\pi}} \geq \epsilon$ . Finally, we can use scaling to reduce the horizon from  $2H$  to  $H$ . □

## G.2 Information theoretical lower sample complexity bound over problems in $\mathcal{M}_{d_m}$ for identifying $\epsilon$ -optimal policy.

For all  $0 < d_m \leq \frac{1}{SA}$ , let the class of problems be

$$\mathcal{M}_{d_m} := \{(\mu, M) \mid \min_{t, s_t, a_t} d_t^\mu(s_t, a_t) \geq d_m\},$$

now we consider deriving minimax lower bound over this class.

**Theorem G.2.** *Under the same condition of Theorem G.1. In addition assume  $0 < d_m \leq \frac{1}{SA}$ . There exists another universal constant  $c$  such that when  $n \leq cH^3/d_m\epsilon^2$ , we always have*

$$\inf_{v^{\pi_{alg}}} \sup_{(\mu, M) \in \mathcal{M}_{d_m}} \mathbb{P}_{\mu, M}(v^* - v^{\pi_{alg}} \geq \epsilon) \geq p.$$

*Proof.* The hard instance  $(\mu, M)$  we used is based on Theorem G.1, which is described as follows.

- for the MDP  $M = (\mathcal{S}, \mathcal{A}, r, P, d_1, 2H+2)$ ,
  - Initial distribution  $d_1$  will always enter state  $s_0$ , and there are two actions with action  $a_1$  always transitions to  $s_{\text{yes}}$  and action  $a_2$  always transitions to  $s_{\text{no}}$ . The reward at the first time  $r_1(s, a) = 0$  for any  $s, a$ .
  - For state  $s_{\text{no}}$ , it will always transition back to itself regardless of the action and receive reward 0, i.e.

$$P_t(s_{\text{no}}|s_{\text{no}}, a) = 1, r_t(s_{\text{no}}, a) = 0, \forall t, \forall a.$$

- For state  $s_{\text{yes}}$ , it will transition to the MDP construction in Theorem G.1 with horizon  $2H$  and  $s_{\text{yes}}$  always receives reward zero.
- For  $t = 1$ , choose  $\mu(a_1|s_0) = \frac{1}{2}d_mSA$  and  $\mu(a_2|s_0) = 1 - \frac{1}{2}d_mSA$ . For  $t \geq 2$ , choose  $\mu$  to be uniform policy, i.e.  $\mu(a_t|s_t) = 1/A$ .

Based on this construction, the optimal policy has the form  $\pi^* = (a_1, \dots)$  and therefore the MDP branch that enters  $s_{\text{no}}$  is uninformative. Hence, data collected by that part is uninformed about the optimal policy and there is only  $\frac{1}{2}d_mSA$  proportion of data from  $s_{\text{yes}}$  are useful. Moreover, by Theorem G.1 the rest of Markov chain succeeded from  $s_{\text{yes}}$  requires  $\Omega(H^3SA/\epsilon^2)$  episodes (regardless of the exploration strategy/logging policy), so the actual data complexity needed for the whole construction  $(\mu, M)$  is  $\frac{\Omega(H^3SA/\epsilon^2)}{d_mSA} = \Omega(H^3/d_m\epsilon^2)$ .

It remains to check this construction  $\mu, M$  stays within  $\mathcal{M}_{d_m}$ .

- For  $t = 1$ , we have  $d_1(s_0, a_1) = \frac{1}{2}d_mSA \geq d_m$  (since  $S \geq 2$ ) and  $d_1(s_0, a_2) = 1 - \frac{1}{2}d_mSA \geq d_m$  (this is since  $d_m \leq \frac{1}{SA} \leq \frac{2}{2+SA}$ );
- For  $t = 2$ ,  $d_2(s_{\text{yes}}, a) = \frac{1}{2}d_mSA \cdot \frac{1}{A} = \frac{1}{2}d_mS \geq d_m$  (since  $S \geq 2$ ) and similar for  $s_{\text{no}}$ ;
- For  $t \geq 3$ , for  $g_h$  and  $b_h$  in the sub-chain inherited from  $s_{\text{yes}}$ , note  $d_h(g_h) \leq d_{h+1}(g_{h+1})$  (since  $g_h$  and  $b_h$  are absorbing states regardless of actions), therefore  $d_h(g_h) \geq d_1(g_1) = d_1(s_{\text{yes}}) \cdot \mathbb{P}(g_1|s_{\text{yes}}) = \frac{1}{2}d_mSA \cdot \frac{1}{S} = \frac{1}{2}d_mA$ , since  $\mu$  is uniform so  $d_h(g_h, a) \geq \Omega(d_mA) \cdot \frac{1}{A} = \Omega(d_m)$  for all  $a$ . Similar result can be derived for  $b_h$  in identical way.

For bandit state, we have for all  $i \in [S - 2]$ ,

$$\begin{aligned} d_{t+1}^\mu(s_{t+1,i}) &\geq \mathbb{P}^\mu(s_{t+1,i}, s_{t,i}, s_{t-1,i}, \dots, s_{2,i}, s_{1,i}, s_{\text{yes}}, s_0) \\ &= \prod_{u=1}^t \mathbb{P}^\mu(s_{u+1,i}|s_u) \mathbb{P}^\mu(s_{1,i}|s_{\text{yes}}) \mathbb{P}^\mu(s_{\text{yes}}|s_0) \\ &= \left(1 - \frac{1}{H}\right)^t \left(\frac{1}{S}\right) \left(\frac{1}{2}d_mSA\right) \geq cd_mA, \end{aligned}$$

now by  $\mu$  is uniform we have  $d_{t+1}^\mu(s_{t+1,i}, a) \geq \Omega(d_mA) \cdot \frac{1}{A} = \Omega(d_m)$  for all  $a$ . This concludes the proof.  $\square$

**Remark G.3.** A directly corollary is that the sample complexity in Theorem 4.1 part 3. is optimal. Indeed, for the case  $\epsilon_{\text{opt}} = 0$ , Theorem 4.1 implies  $\hat{\pi}$  is the  $\epsilon$ -optimal policy learned with sample complexity  $O(H^3 \log(HSA/\delta)/d_m\epsilon^2)$ . Theorem G.2 implies this sample complexity cannot be further reduced up to the logarithmic factor.

### G.3 Information theoretical lower sample complexity bound for uniform convergence in OPE.

By applying Theorem G.2, we can now prove Theorem 3.8.

*Proof of Theorem 3.8.* We prove it by contradiction. Suppose there is one off-policy evaluation method  $\hat{v}^\pi$  such that

$$\sup_{\pi \in \Pi} |\hat{v}^\pi - v^\pi| \leq o\left(\sqrt{\frac{H^3}{d_m n}}\right),$$

where  $o(\cdot)$  represents the standard small  $o$ -notation. Then by

$$\begin{aligned} 0 &\leq v^{\pi^*} - v^{\hat{\pi}^*} = v^{\pi^*} - \hat{v}^{\hat{\pi}^*} + \hat{v}^{\hat{\pi}^*} - v^{\hat{\pi}^*} \\ &\leq |v^{\pi^*} - \hat{v}^{\pi^*}| + |\hat{v}^{\hat{\pi}^*} - v^{\hat{\pi}^*}| \leq 2 \sup_{\pi} |v^\pi - \hat{v}^\pi|. \end{aligned}$$

this OPE method implies a  $\epsilon$ -optimal policy learning algorithm with sample complexity  $o(H^3/d_m\epsilon^2)$  which is smaller than the information theoretical lower bound obtained in Theorem G.2. Contradiction!  $\square$

## H Proofs of Theorem 4.1

*Proof of Theorem 4.1.* Part 1. and Part 2. are just direct corollaries. We only prove Part 3. here. Indeed, by definition of empirical optimal policy we have  $\hat{Q}^{\pi^*} \leq \hat{Q}^{\hat{\pi}^*}$ , so we have the following:

$$\begin{aligned} Q_1^{\pi^*} - Q_1^{\hat{\pi}} &= Q_1^{\pi^*} - \hat{Q}_1^{\hat{\pi}^*} + \hat{Q}_1^{\hat{\pi}^*} - \hat{Q}_1^{\hat{\pi}} + \hat{Q}_1^{\hat{\pi}} - Q_1^{\hat{\pi}} \\ &\leq Q_1^{\pi^*} - \hat{Q}_1^{\pi^*} + \hat{Q}_1^{\hat{\pi}^*} - \hat{Q}_1^{\hat{\pi}} + \hat{Q}_1^{\hat{\pi}} - Q_1^{\hat{\pi}} \\ &\leq Q_1^{\pi^*} - \hat{Q}_1^{\pi^*} + \epsilon_{\text{opt}} \cdot \mathbf{1} + \hat{Q}_1^{\hat{\pi}} - Q_1^{\hat{\pi}} \end{aligned}$$

and  $\hat{Q}_1^{\hat{\pi}} - Q_1^{\hat{\pi}}$  can be bounded by Theorem 3.7 using local uniform convergence.  $Q_1^{\pi^*} - \hat{Q}_1^{\pi^*}$  can be bounded by  $O(\sqrt{\frac{H^3 \log(HSA/\delta)}{nd_m}})$  using the similar technique in Section F even without introducing  $\epsilon_{\text{opt}}$  since  $\pi^*$  is a fixed policy. All these implies:

$$Q_1^{\pi^*} - Q_1^{\hat{\pi}} \leq \left( O(\sqrt{\frac{H^3 \log(HSA/\delta)}{nd_m}}) + \epsilon_{\text{opt}} \right) \cdot \mathbf{1}.$$

Especially when  $\epsilon_{\text{opt}} = 0$  then this is slightly stronger than the stated result since:

$$v_1^{\pi^*} - v_1^{\hat{\pi}^*} = Q_1^{\pi^*}(\cdot, \pi^*(\cdot)) - Q_1^{\hat{\pi}^*}(\cdot, \hat{\pi}^*(\cdot)) \leq Q_1^{\pi^*}(\cdot, \pi^*(\cdot)) - Q_1^{\hat{\pi}^*}(\cdot, \pi^*(\cdot)) \leq \|Q_1^{\pi^*} - Q_1^{\hat{\pi}^*}\|_{\infty} \leq O(\sqrt{\frac{H^3 \log(HSA/\delta)}{nd_m}}) \cdot \mathbf{1}$$

□

## I Simulation details

The non-stationary MDP with used for the experiments have 2 states  $s_0, s_1$  and 2 actions  $a_1, a_2$  where action  $a_1$  has probability 1 always going back the current state and for action  $a_2$ , there is one state s.t. after choosing  $a_2$  the dynamic transitions to both states with equal probability  $\frac{1}{2}$  and the other one has asymmetric probability assignment ( $\frac{1}{4}$  and  $\frac{3}{4}$ ). The transition after choosing  $a_2$  is changing over different time steps therefore the MDP is non-stationary and the change is decided by a sequence of pseudo-random numbers. More formally,  $P_t$  can be either

$$\mathbb{P}(s_0|s_0, a_1) = 1; \mathbb{P}(s_1|s_1, a_1) = 1; \mathbb{P}(\cdot|s_0, a_2) = \begin{cases} \frac{1}{2}, & \text{if } \cdot = s_1 \\ \frac{1}{2}, & \text{if } \cdot = s_0 \end{cases} \quad ; \quad \mathbb{P}(\cdot|s_1, a_2) = \begin{cases} \frac{3}{4}, & \text{if } \cdot = s_1 \\ \frac{1}{4}, & \text{if } \cdot = s_0 \end{cases}$$

or

$$\mathbb{P}(s_0|s_0, a_1) = 1; \mathbb{P}(s_1|s_1, a_1) = 1; \mathbb{P}(\cdot|s_0, a_2) = \begin{cases} \frac{1}{4}, & \text{if } \cdot = s_1 \\ \frac{3}{4}, & \text{if } \cdot = s_0 \end{cases} \quad ; \quad \mathbb{P}(\cdot|s_1, a_2) = \begin{cases} \frac{1}{2}, & \text{if } \cdot = s_1 \\ \frac{1}{2}, & \text{if } \cdot = s_0 \end{cases}$$

Moreover, to make the learning problem non-trivial we use non-stationary rewards with 4 categories, *i.e.*  $r_t(s, a) \in \{\frac{1}{4}, \frac{2}{4}, \frac{3}{4}, 1\}$  and assignment of  $r_t(s, a)$  for each value is changing over time. That means, one possible assignment can be

$$r_t(s_0, a_1) = 1/4, r_t(s_0, a_2) = 2/4, r_t(s_1, a_1) = 3/4, r_t(s_1, a_2) = 1/4.$$

Moreover, the logging policy in Figure 1(a) is uniform with  $\mu_t(a_1|s) = \mu_t(a_2|s) = \frac{1}{2}$  for both states. We implement the non-stationary MDP in the Python environment and pseudo-random numbers  $p_t, r_t$ 's are generated by keeping `numpy.random.seed(100)`.



We fix episodes  $n = 2048$  and run each algorithm under  $K = 100$  macro-replications with data  $\mathcal{D}_{(k)} = \left\{ (s_t^{(i)}, a_t^{(i)}, r_t^{(i)}) \right\}_{(k)}^{i \in [n], t \in [H]}$ , and use each  $\mathcal{D}_{(k)}$  ( $k = 1, \dots, K$ ) to construct an estimator  $\hat{v}_{[k]}^\pi$ , then the (empirical) RMSE for fixed policy is computed as:

$$\text{RMSE\_FIX} = \sqrt{\frac{\sum_{k=1}^K (\hat{v}_{[k]}^\pi - v_{\text{true}}^\pi)^2}{K}},$$

and RMSE for suboptimality gap is computed as

$$\text{RMSE\_SUB} = \sqrt{\frac{\sum_{k=1}^K (v_{[k]}^{\hat{\pi}^*} - v_{\text{true}}^{\pi^*})^2}{K}},$$

and RMSE for empirical optimal policy gap is computed as

$$\text{RMSE\_EMPIRICAL} = \sqrt{\frac{\sum_{k=1}^K (\hat{v}_{[k]}^{\hat{\pi}^*} - v_{\text{true}}^{\pi^*})^2}{K}},$$

where  $v_{\text{true}}^\pi$  is obtained by calculating  $P_{t+1,t}^\pi(s'|s) = \sum_a P_{t+1,t}(s'|s, a) \pi_t(a|s)$ , the marginal state distribution  $d_t^\pi = P_{t,t-1}^\pi d_{t-1}^\pi$ ,  $r_t^\pi(s_t) = \sum_{a_t} r_t(s_t, a_t) \pi_t(a_t|s_t)$  and  $v_{\text{true}}^\pi = \sum_{t=1}^H \sum_{s_t} d_t^\pi(s_t) r_t^\pi(s_t)$ .  $v_{\text{true}}^{\pi^*}$  is obtained by running Value Iteration exhaustively until the error converges to 0. The average relative error for suboptimality (average of  $|v_{[k]}^{\hat{\pi}^*} - v_{\text{true}}^{\pi^*}|/v_{\text{true}}^{\pi^*}$ ) at  $H = 1000$  is 0.0011. Lastly, we also show the scaling of  $|\hat{v}^{\hat{\pi}^*} - v^{\pi^*}|$  in Figure 3, which shares a similar pattern as the suboptimality plot as a whole.<sup>15</sup>

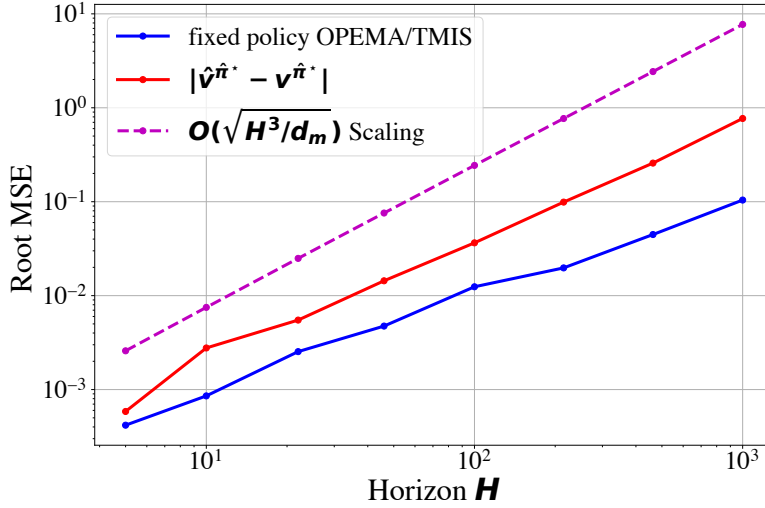


Figure 3: Log-log plot showing the dependence on horizon of uniform OPE and pointwise OPE via learning ( $|\hat{v}^{\hat{\pi}^*} - v^{\pi^*}|$ ) over a non-stationary MDP example.

## J On improvement over vanilla simulation lemma for fixed policy evaluation

**Vanilla simulation lemma, Lemma 1 of Jiang (2018).** Without loss of generality, assuming reward is deterministic function over state-action. By definition of Bellman equation, we have the following:

$$\hat{V}_t^\pi = r + \hat{P}_{t+1}^\pi \hat{V}_{t+1}^\pi, \quad V_t^\pi = r + P_{t+1}^\pi V_{t+1}^\pi,$$

<sup>15</sup>Here we do point out the empirical dependence on  $H$  for  $|\hat{v}^{\hat{\pi}^*} - v^{\pi^*}|$  in the Figure 3 is actually less than  $H^{1.5}$ , this comes from that the MDP example we choose is not the “hardest” example for quantity  $|\hat{v}^{\hat{\pi}^*} - v^{\pi^*}|$ , as opposed to quantity  $|v^* - v^{\hat{\pi}^*}|$  in Figure 1.

define  $\epsilon_P = \sup_{t, s_t, a_t} \|\hat{P}_t(\cdot|s_t, a_t) - P_t(\cdot|s_t, a_t)\|_1$ , then by Hoeffding's inequality and union bound, with probability  $1 - \delta$ ,

$$\epsilon_P \leq S \cdot \sup_{t, s_t, a_t} \|\hat{P}_t(\cdot|s_t, a_t) - P_t(\cdot|s_t, a_t)\|_\infty \leq S \cdot \sup_{t, s_t, a_t} O\left(\sqrt{\frac{\log(HSA/\delta)}{n_{s_t, a_t}}} \mathbf{1}(E_t)\right) = O\left(\sqrt{\frac{S^2 \log(HSA/\delta)}{n \cdot d_m}}\right)$$

then

$$\begin{aligned} \hat{V}_t^\pi - V_t^\pi &= \hat{P}_{t+1}^\pi \hat{V}_{t+1}^\pi - P_{t+1}^\pi V_{t+1}^\pi \\ &\leq \left( \|\hat{P}_{t+1}^\pi - P_{t+1}^\pi\|_1 \|\hat{V}_{t+1}^\pi\|_\infty + \|P_{t+1}^\pi\|_1 \|\hat{V}_{t+1}^\pi - V_{t+1}^\pi\|_\infty \right) \cdot \mathbf{1} \\ &\leq \left( H\epsilon_P + \|\hat{V}_{t+1}^\pi - V_{t+1}^\pi\|_\infty \right) \cdot \mathbf{1}, \end{aligned}$$

solving recursively, we have

$$\|\hat{V}_1^\pi - V_1^\pi\|_\infty \leq H^2 \epsilon_P \leq O\left(\sqrt{\frac{H^4 S^2 \log(HSA/\delta)}{n \cdot d_m}}\right).$$

This verifies SL has complexity  $\tilde{O}(H^4 S^2 / d_m \epsilon^2)$ . We do point out above standard analysis can be improved (e.g. Jiang (2018) Section 2.2) to  $\tilde{O}(H^4 S / d_m \epsilon^2)$ , then in this case our analysis (Lemma 3.4) has an improvement of  $H^2 S$  with respect to the modified result.

## K Algorithms

---

### Algorithm 1 OPEMA

---

**Input:** Logging data  $\mathcal{D} = \{\{s_t^{(i)}, a_t^{(i)}, r_t^{(i)}\}_{t=1}^H\}_{i=1}^n$  from the behavior policy  $\mu$ . A target policy  $\pi$  which we want to evaluate its cumulative reward.

- 1: Calculate the on-policy estimation of initial distribution  $d_1(\cdot)$  by  $\hat{d}_1(s) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}(s_1^{(i)} = s)$ , and set  $\hat{d}_1^\mu(\cdot) := \hat{d}_1(\cdot)$ ,  $\hat{d}_1^\pi(s) := \hat{d}_1(\cdot)$ .
- 2: **for**  $t = 2, 3, \dots, H$  **do**
- 3:   Choose all transition data at time step  $t$ ,  $\{s_t^{(i)}, a_t^{(i)}, r_t^{(i)}\}_{i=1}^n$ .
- 4:   Calculate the on-policy estimation of  $d_t^\mu(\cdot)$  by  $\hat{d}_t^\mu(s) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}(s_t^{(i)} = s)$ .
- 5:   Set the off-policy estimation of  $\hat{P}_t(s_t|s_{t-1}, a_{t-1})$ :

$$\hat{P}_t(s_t|s_{t-1}, a_{t-1}) := \frac{\sum_{i=1}^n \mathbf{1}[(s_t^{(i)}, a_{t-1}^{(i)}, s_{t-1}^{(i)}) = (s_t, s_{t-1}, a_{t-1})]}{n_{s_{t-1}, a_{t-1}}}$$

when  $n_{s_{t-1}, a_{t-1}} > 0$ . Otherwise set it to be zero.

- 6:   Estimate the reward function

$$\hat{r}_t(s_t, a_t) := \frac{\sum_{i=1}^n r_t^{(i)} \mathbf{1}(s_t^{(i)} = s_t, a_t^{(i)} = a_t)}{\sum_{i=1}^n \mathbf{1}(s_t^{(i)} = s_t, a_t^{(i)} = a_t)}.$$

when  $n_{s_t, a_t} > 0$ . Otherwise set it to be zero.

- 7:   Set  $\hat{d}_t^\pi(\cdot, \cdot)$  according to  $\hat{d}_t^\pi = \hat{P}_t^\pi \hat{d}_{t-1}^\pi$ , where  $\hat{d}_t^\pi(\cdot, \cdot)$  is the estimated state-action distribution.

- 8: **end for**

- 9: Substitute the all estimated values above into  $\hat{v}^\pi = \sum_{t=1}^H \langle \hat{d}_t^\pi, \hat{r}_t \rangle$  to obtain  $\hat{v}^\pi$ , the estimated value of  $\pi$ .
- 

**Remark K.1.** In short, we can see Algorithm 2 requires the splitting data size  $M$  which is undecided by Yin and Wang (2020) and that makes the hyper-parameter requiring additional concrete specifications to make the data

---

**Algorithm 2** Data Splitting TMIS in [Yin and Wang \(2020\)](#)

---

**Input:** Logging data  $\mathcal{D} = \{\{s_t^{(i)}, a_t^{(i)}, r_t^{(i)}\}_{t=1}^H\}_{i=1}^n$  from the behavior policy  $\mu$ . A target policy  $\pi$  which we want to evaluate its cumulative reward. Requiring splitting data size  $M$ .

- 1: Randomly splitting the data  $\mathcal{D}$  evenly into  $N$  folds, with each fold  $|\mathcal{D}^{(i)}| = M$ , *i.e.*  $n = M \cdot N$ .
  - 2: **for**  $i = 1, 2, \dots, N$  **do**
  - 3:   Use Algorithm 1 to estimate  $\hat{v}_{(i)}^\pi$  with data  $\mathcal{D}^{(i)}$ .
  - 4: **end for**
  - 5: Use the mean of  $\hat{v}_{(1)}^\pi, \hat{v}_{(2)}^\pi, \dots, \hat{v}_{(N)}^\pi$  as the final estimation of  $v^\pi$ .
- 

*splitting estimator sample efficient. In contrast, OPEMA in Algorithm 1 is defined without ambiguity and can be implemented without extra work.*

*Their results require number of episodes in each splitted data  $M$  to satisfy  $\tilde{O}(\sqrt{nSA}) > M > O(HSA)$ . To achieve data efficiency, they need  $n \approx \Theta(H^2SA/\epsilon^2)$  and by that condition  $M$  has to satisfy  $M \approx C \cdot HSA$ . In this case, data-splitting version needs to create  $N = n/M$  empirical transition dynamics and each dynamics use  $H^3/N \approx C \cdot H^2SA/\epsilon^2$  episodes which is less than the lower bound ( $O(H^3)$ ) required for learning. Most critically, due to data-splitting it has  $N$  empirical transitions hence it is not clear which transition to plan over. Therefore in this sense their result does not enables efficient offline learning. Our Analysis for unsplitted version (OPEMA) addresses all these issues.*