
Federated UCBVI: Communication-Efficient Federated Regret Minimization with Heterogeneous Agents

Safwan Labbi¹ Daniil Tiapkin^{1,2} Lorenzo Mancini¹ Paul Mangold¹ Eric Moulines^{1,3}

¹ CMAP, CNRS, École Polytechnique, Institut Polytechnique de Paris, 91120 Palaiseau, France

² Université Paris-Saclay, CNRS, LMO, 91405, Orsay, France

³ Mohamed bin Zayed University of Artificial Intelligence, UAE

Abstract

In this paper, we present the Federated Upper Confidence Bound Value Iteration algorithm (**Fed-UCBVI**), a novel extension of the UCBVI algorithm (Azar et al., 2017) tailored for the federated learning framework. We prove that the regret of **Fed-UCBVI** scales as $\tilde{O}(\sqrt{H^3|S||A|T/M})$, with a small additional term due to heterogeneity, where $|S|$ is the number of states, $|A|$ is the number of actions, H is the episode length, M is the number of agents, and T is the number of episodes. Notably, in the single-agent setting, this upper bound matches the minimax lower bound up to polylogarithmic factors, while in the multi-agent scenario, **Fed-UCBVI** has linear speed-up. To conduct our analysis, we introduce a new measure of heterogeneity, which may hold independent theoretical interest. Furthermore, we show that, unlike existing federated reinforcement learning approaches, **Fed-UCBVI**'s communication complexity only marginally increases with the number of agents.

1 INTRODUCTION

Federated reinforcement learning (FRL, Zhuo et al., 2019; Qi et al., 2021) adapts the principles of federated learning (FL, McMahan et al., 2017) to the domain of reinforcement learning (RL, Sutton and Barto, 2018). It enables multiple agents, evolving in independent environments, to learn a policy collaboratively without directly exchanging their states/actions. To learn together, agents communicate under the supervision of a central server (CS), aiming to maximize the

expected rewards averaged across all agents. Consequently, agents participating in FRL may learn better policies with fewer interactions with the environment. FRL appears to be a promising solution for reducing the cost of training. However, the efficient implementation of FRL faces significant challenges. Similarly to FL, agents typically evolve in different environments and often have limited computational power and communication bandwidth. Furthermore, the traditional challenges of RL, such as balancing exploration and exploitation, remain. Thus, there is a growing demand for methods tailored for FRL, aiming to reduce communication complexity (i.e., the number of communications) while maintaining efficient exploration and learning.

FRL has attracted considerable attention in recent years, with a strong focus put on federated versions of Q-Learning. This research often relies on one of two following assumptions: either (1) all agents operate in *identical environments* (Chen et al., 2023; Zheng et al., 2024, 2025), or (2) a generative model is available, allowing access to sampling from any state-action pair without exploration (Jin et al., 2022; Wang et al., 2024). Another notable category of methods, called distributed reinforcement learning (Bai et al., 2019; Zhang et al., 2020), enables agents to address RL problems collaboratively. However, these methods require centralizing observational data on a single server, which may not be feasible in real applications.

Unfortunately, the aforementioned approaches do not address the exploration-exploitation trade-off in heterogeneous environments. Furthermore, their high communication complexity poses a major challenge for their use, even with homogeneous agents.

In this paper, we introduce the algorithm **Fed-UCBVI** for tabular episodic FRL and we analyze its *federated regret*, i.e., the regret averaged across all agents, in the presence of environmental heterogeneity. The tabular FRL problem involves M agents, each inter-

Table 1: Comparison with related algorithms in the online setting

Type	Algorithm	Heterogeneity	Communication complexity	Regret
Model-based	Concurrent UCBVI (Azar et al. 2017)	\times	$\mathcal{O}(T)$	$\tilde{\mathcal{O}}(\sqrt{H^3 S A T/M})$
	Byzan-UCBVI (Chen et al. 2023)	\times	$\mathcal{O}(M \cdot S A H \cdot \log(T))$	$\tilde{\mathcal{O}}(\sqrt{H^4 S ^2 A T/M})$
	Fed-UCBVI (our work)	\checkmark	$\mathcal{O}(S A H \cdot \log(T))$	$\tilde{\mathcal{O}}(\sqrt{H^3 S A T/M})$
Model-free	Concurrent UCB-Advantage (Zhang et al. 2020)	\times	$\mathcal{O}(T)$	$\tilde{\mathcal{O}}(\sqrt{H^3 S A T/M})$
	FedQ-Bernstein (Zheng et al. 2024)	\times	$\mathcal{O}(M \cdot S A H^3 \cdot \log(T))$	$\tilde{\mathcal{O}}(\sqrt{H^4 S A T/M})$
	FedQ-Advantage (Zheng et al. 2025)	\times	$\mathcal{O}(M \cdot S A H^2 \log(H) \cdot \log(T))$	$\tilde{\mathcal{O}}(\sqrt{H^3 S A T/M})$
Lower Bound (Jin et al. 2018; Domingues et al. 2021b)		\times	?	$\tilde{\mathcal{O}}(\sqrt{H^3 S A T/M})$

*The results are derived in a homogeneous setting. For all bounds, only the leading term with respect to the dependence on T is shown. H : number of steps per episode; T : total episodes collected per agent; $|S|$: number of states; $|A|$: number of actions; M : number of agents.

acting with its own environment, modeled as a finite-horizon Markov Decision Process (MDP). For an agent $i \in [M]$, a finite-horizon MDP is defined by a tuple $\mathcal{M}^i := (\mathcal{S}, \mathcal{A}, H, \{P_h^i\}_{h \in [H]}, \{r_h^i\}_{h \in [H]})$, where \mathcal{S} is the finite state space, \mathcal{A} is the finite action space, H is the number of steps in one episode (also referred to as a planning horizon), $P_h^i(s'|s, a)$ denotes the probability of transitioning from a state $s \in \mathcal{S}$ to the next state $s' \in \mathcal{S}$ after taking action $a \in \mathcal{A}$ at step h for agent i , and $r_h^i(s, a)$ is a bounded deterministic reward function that satisfies $r_h^i(s, a) \in [0, 1]$ for all $(s, a, h, i) \in \mathcal{S} \times \mathcal{A} \times [H] \times [M]$. Note that both the transition probabilities (kernel) and the reward function can vary depending on the decision-making step $h \in [H]$. The learning process is divided into T episodes, each of length H . Both the transition kernel and the reward function are assumed to be *unknown* to all agents and the central server (CS).

Fed-UCBVI is a model-based approach where each agent independently estimates its local state-action transition kernel. These local estimates are then used to compute state-action value functions, which are aggregated by a CS using an adaptive scheme that accounts for each agent’s level of uncertainty. Communication complexity is managed through an adaptive communication strategy triggered by the optimization process’s progress and ensures efficient coordination. Overall, our contributions are:

- We propose **Fed-UCBVI**, an FRL algorithm designed to aggregate the local estimators of each agent. We prove that the federated regret of **Fed-UCBVI** scales as $\mathcal{O}(\sqrt{H^3|S||A|T/M})$, up to a heterogeneity term which scales proportionally to our heterogeneity measure. This shows that **Fed-UCBVI** achieves a linear speedup and effectively accelerates training compared to single-agent RL. To our knowledge, **Fed-UCBVI** is the first provably efficient algorithm for regret minimization in heterogeneous environments.
- To analyze **Fed-UCBVI**, we introduce a new measure of heterogeneity that quantifies the divergence of each agent’s state-transition kernel from a baseline

kernel, which may be of independent interest.

- We develop a novel method for reducing the communication cost. We prove that the communication complexity of **Fed-UCBVI** is $\mathcal{O}(M \log \log T + \log T)$. This is a significant improvement over existing methods (e.g., Zheng et al., 2024), that require $\mathcal{O}(M \log T)$ communication rounds.
- We validate our theoretical results through numerical experiments on FRL problems, demonstrating that our algorithm outperforms existing FRL baselines with theoretical guarantees. In particular, our simulations show a significant improvement in regret compared to Fed-Q-learning (Zheng et al., 2024) for different degrees of heterogeneity.

The paper is organized as follows: we review the related work in Section 2, and introduce the necessary mathematical background in Section 3. In Section 4, we introduce and analyze the **Fed-UCBVI** algorithm. Then, we present numerical experiments in Section 5.

2 RELATED WORK

Reinforcement Learning. Two main approaches have been proposed for regret minimization in the single-agent, finite-horizon tabular setting: (i) model-based algorithms (Azar et al., 2017; Dann et al., 2017; Zanette and Brunskill, 2019; Zhang et al., 2024b), and (ii) model-free algorithms (Jin et al., 2018; Zhang et al., 2020; Li et al., 2021).

Both approaches offer algorithms that achieve the minimax optimal lower bound up to poly-logarithmic factors, specifically $\Omega(\sqrt{H^3|S||A|T})$ (Jin et al., 2018; Domingues et al., 2021b). Among these, UCBVI (Azar et al., 2017), which is based on the principle of optimism in the face of uncertainty, was the first algorithm to achieve the minimax bound.

Federated Reinforcement Learning. The FRL method most closely related to ours is the Byzantine robust distributed UCBVI algorithm (Chen et al., 2023), which assumes homogeneous agents. This algorithm achieves a regression bound of $\tilde{\mathcal{O}}(\sqrt{H^4|S|^2|A|T/M})$

and a communication complexity that scales logarithmically with the number of episodes T .

In contrast, our method achieves a regret of $\tilde{O}(\sqrt{H^3|\mathcal{S}||\mathcal{A}|T/M})$, which is optimal in single-agent environments. Moreover, we also provide guarantees in heterogeneous environments.

Other FRL approaches are based on model-free methods. Zhang et al. (2020) proposed a federated variant of Q-learning, with regret $\tilde{O}(\sqrt{H^3|\mathcal{S}||\mathcal{A}|T/M})$, and communication complexity linear in T . Zheng et al. (2024) later reduced the communication cost to $O(M \log T)$, but introduced an additional factor of H in the regret bound. More recently, Zheng et al. (2025) improved both regret and communication cost. However, their method still requires homogeneous agents, and the communication complexity remains $O(M \log T)$. FRL’s communication complexity has been studied in several related settings: finding the optimal Q-function Salgia and Chi (2024), policy optimization Lan et al. (2024); Chen et al. (2021) and policy evaluation Mangold et al. (2025).

3 SETTING

3.1 Federated Reinforcement Learning

Policy and Value Functions. A deterministic policy π is a set of functions $\pi_h : \mathcal{S} \rightarrow \mathcal{A}$ where $\pi_h(s) \in \mathcal{A}$, $h \in [H]$. The value function $\mathcal{V}_h^{i,\pi}$, is defined as:

$$\mathcal{V}_h^{i,\pi}(s) = \mathbb{E}_\pi \left[\sum_{h'=h}^H r_{h'}^i(s_{h'}^i, a_{h'}^i) \middle| s_h^i = s \right], \quad (1)$$

where for all $h \leq h' \leq H$, $a_{h'}^i \sim \pi_{h'}^i(\cdot | s_{h'}^i)$ and for all $h \leq h' \leq H-1$, $s_{h'+1}^i \sim P_{h'}^i(\cdot | s_{h'}^i, a_{h'}^i)$. Similarly, the Q-function of a policy π for agent i at step h is

$$\mathcal{Q}_h^{i,\pi}(s, a) := \mathbb{E}_\pi \left[\sum_{h'=h}^H r_{h'}^i(s_{h'}^i, a_{h'}^i) \middle| s_h^i = s, a_h^i = a \right],$$

and satisfies the Bellman equations

$$\begin{aligned} \mathcal{Q}_h^{i,\pi}(s, a) &= r_h^i(s, a) + P_h^i \mathcal{V}_{h+1}^{i,\pi}(s, a), \\ \mathcal{V}_{h+1}^{i,\pi}(s) &= \mathcal{Q}_h^{i,\pi}(s, \pi_h(s)). \end{aligned} \quad (2)$$

Additionally, the optimal Q-value satisfies the optimal Bellman equations

$$\begin{aligned} \mathcal{Q}_h^{i,*}(s, a) &= r_h^i(s, a) + P_h^i \mathcal{V}_{h+1}^{i,*}(s, a), \\ \mathcal{V}_h^{i,*}(s) &= \max_{a \in \mathcal{A}} \mathcal{Q}_h^{i,*}(s, a). \end{aligned} \quad (3)$$

Learning Protocol. At the beginning of each episode $t \in [T]$, all agents select a common policy π_t , which is computed based on the information *exchanged* prior to episode t . Subsequently, each agent generates an independent trajectory of length H . At each step h , an agent observes its state $s_{t,h}^i \in \mathcal{S}$ and

takes an action $a_{t,h}^i = \pi_{t,h}^i(s_{t,h}^i) \in \mathcal{A}$. The agent then observes the next state $s_{t,h+1}^i$ according to the transition probabilities $P_h^i(\cdot | s_{t,h}^i, a_{t,h}^i)$ and receives a deterministic reward $r_{t,h}^i = r_h^i(s_{t,h}^i, a_{t,h}^i)$. After generating these trajectories, agents *may* exchange information through the central server.

Federated Regret. The performance of the learning algorithm is evaluated using the *federated regret*, defined as

$$\mathfrak{R}(T) := \max_{\pi} \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \mathcal{V}_1^{i,\pi}(s_{t,1}^i) - \mathcal{V}_1^{i,\pi_t}(s_{t,1}^i). \quad (4)$$

This regret measures the cumulative difference, in expectation, between the average value of the optimal collaborative policy and the policies used throughout the training procedure.

Communication Complexity and Cost. The *communication complexity*, denoted by $\mathfrak{C}(T)$, is defined as the number of episodes where communication between the CS and the agents occurs. The *communication cost* refer to the total number of bits exchanged between the central server and the agents during the learning process. The objective of the FRL algorithm is to simultaneously minimize both the regret $\mathfrak{R}(T)$ and the communication complexity $\mathfrak{C}(T)$.

3.2 Environmental Heterogeneity

The environments in which agents evolve may differ from one to another. However, since agents aim to learn a shared policy, environmental heterogeneity must be small. To measure this, we introduce a new notion of heterogeneity, decomposing each agent’s state-action transition kernel into a common part, shared by all agents, and an individual part that reflects unique environmental characteristics. Formally, this is captured by the following assumption.

A-1. *There exists a non-homogeneous transition kernel $\{P_h^c\}_{h \in [H]}$, M individual non-homogeneous transition kernels $\{P_h^{\text{ind},i}\}_{h \in [H]}$ for any $i \in [M]$, and a constant $\varepsilon_p \in [0, 1)$, such that for any $i \in [M]$ and $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]$,*

$$P_h^i(s' | s, a) = (1 - \varepsilon_p) P_h^c(s' | s, a) + \varepsilon_p P_h^{\text{ind},i}(s' | s, a).$$

This assumption can be interpreted as follows: the environment of all agents reacts similarly to a given action, but occasionally, due to local specificities, an individual environment reacts differently. This assumption is related to the ε -contamination model proposed by Huber (1964). Likewise, we assume that agents receive comparable rewards for a given state-action pair.

A-2. *There exists a constant $\varepsilon_r \in [0, 1)$ such that for all $(i, j) \in [M]$, and for all $h \in [H]$ it holds that*

$$\|r_h^i - r_h^j\|_\infty \leq \varepsilon_r.$$

Note that **A-1** implies the following bound on the difference between the common transition kernel and each agent’s transition kernel, measured in L_1 -norm,

$$\max_{(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]} \|\mathbf{P}_h^c(\cdot|s,a) - \mathbf{P}_h^i(\cdot|s,a)\|_1 \leq \varepsilon_p. \quad (5)$$

We prove this inequality in Appendix F. Thus, (1) is slightly stronger than (5), which is the typical assumption in other FRL settings, such as FedSARSA (Zhang et al., 2024a) or policy optimization with access to a simulator (Jin et al., 2022; Wang et al., 2024). The motivation for using **A-1** over (5) lies in the need to control how *samples from \mathbf{P}_h^i* relate to *samples from \mathbf{P}_h^c* . This is crucial in RL without a generative model, as the data generation process is not independent and identically distributed, forcing agents to exploit all the samples they have. In Section 4, we discuss in detail the necessity of this assumption for our analysis.

4 FED-UCBVI ALGORITHM

In this section, we present the **Fed-UCBVI** algorithm, which extends the UCBVI algorithm proposed by Azar et al. (2017) to the federated learning framework. The process involves multiple communication rounds with a CS. The number of episodes in each communication round (or epoch) r is random, and each epoch is decomposed into three phases:

- (i) *Data collection:* During this phase, each agent interacts with its environment using the policy $\pi(r)$ provided by the CS, gathering trajectory data.
- (ii) *Synchronization:* Once any agent meets the synchronization conditions, it sends a synchronization signal to the central server, which then broadcasts this information to all other agents.
- (iii) *Policy update:* In this phase, all agents engage in H sequential communications with the CS. At each step $h = H$ to 1, agents send their local estimates of the Q -values and other related information related to step h to the CS. In return, they receive a global estimate of the V -values, along with an updated policy and related information for that step.

The following sections provide a detailed overview of each of these stages.

Data Collection. At the beginning of round r , each agent $i \in [M]$ follows the policy $\pi(r)$ to collect new trajectories. For $\ell \in \mathbb{N}$, denote by $n_{(r,\ell),h}^i(s,a)$ and $n_{(r,\ell),h}^i(s,a,s')$ the number of visits to a state-action

pair (s,a) and the number of transitions from (s,a) to s' at step h after ℓ episodes in the round r .

Synchronization. At the start of epoch r , all agents receive the current global counters

$$N_{(r),h}(s,a) := \sum_{i=1}^M n_{(r),h}^i(s,a), \quad (6)$$

where $n_{(r),h}^i(s,a) := n_{(r,0),h}^i(s,a)$ is the number of visits of a state-action pair by agent i prior to round r .

During epoch r , after ℓ episodes, agent i sends a synchronization signal if a newly visited state-action-step triplet (s,a,h) is identified and one of two synchronization conditions is met. These conditions depend on whether the total number of visits $N_{(r),h}(s,a)$ exceeds a threshold $\nu(\delta, T) = \tilde{\mathcal{O}}(\varepsilon_p THM + M)$ (see Equation (57) in Appendix E for the full expression).

- 1) *Local Doubling Condition.* If $N_{(r),h}(s,a) \leq \nu(\delta, T)$, an agent i sends the synchronization signal if

$$n_{(r,\ell),h}^i(s,a) > 2n_{(r),h}^i(s,a). \quad (7)$$

- 2) *Globally Estimated Doubling Condition.* If $N_{(r),h}(s,a) > \nu(\delta, T)$, agent i sends the synchronization signal if

$$\hat{N}_{(r,\ell),h}^i(s,a) > 2N_{(r),h}(s,a), \quad (8)$$

where $\hat{N}_{(r,\ell),h}^i(s,a)$ is an estimate of $\sum_{i=1}^M n_{(r,\ell),h}^i(s,a)$ based on the information available to agent i .

Policy Update. Upon receiving the synchronization signal, each agent computes its local estimates of transition probabilities as

$$\hat{\mathbf{P}}_{(r+1),h}^i(s'|s,a) := \frac{n_{(r+1),h}^i(s,a,s')}{n_{(r+1),h}^i(s,a)} \quad (9)$$

if $n_{(r+1),h}^i(s,a) > 0$, otherwise $\hat{\mathbf{P}}_{(r+1),h}^i(s'|s,a) := 1/|\mathcal{S}|$. Additionally, each agent builds a local estimate of the reward function, denoted as $\hat{\mathbf{r}}^i$. This estimate is initially set to a null $H \times \mathcal{S} \times \mathcal{A}$ tensor. The entry $\hat{\mathbf{r}}_h^i(s,a)$ is updated with the received reward whenever agent i selects action a in state s at step h .

Next, the agents and the central server exchange their Q - and V -value estimates. For $h = H, \dots, 1$, each agent computes the local Q -value estimate

$$\hat{Q}_{(r+1),h}^i(s,a) := \left[\hat{\mathbf{r}}_h^i + \hat{\mathbf{P}}_{(r+1),h}^i \hat{\mathbf{V}}_{(r+1),h+1} \right](s,a), \quad (10)$$

using the global value estimate $\hat{\mathbf{V}}_{(r+1),h+1}$ previously received from the CS; note that for $h = H$, this value is set to zero and does not require communication.

Then, the CS collects the local Q -value estimates from all agents, along with additional information

Algorithm 1: Fed-UCBVI

Initialization: $t = 1$; $r = 1$; $\hat{r}_h^i(s, a) = 0$; $N_{(1),h}(s, a) = 0$; $n_{(1,0),h}^i(s, a) = 0$; $\hat{Q}_{(1),h}(s, a) = \hat{V}_{(1),h}(s) = H$ for all $(s, a, h, i) \in \mathcal{S} \times \mathcal{A} \times [H] \times [M]$; $\pi_{(1)} = \{\pi_{(1),h}\}_h$ for some policy $\pi_{(1)}$; and $\nu(\delta, T)$ set as in (57).
while $t \leq T$ **do**

for each agent $i = 1$ **to** M **in parallel do**

Set $l = 1$; $n_{(r),h}^i(s, a) = n_{(r,0),h}^i(s, a)$; and $\hat{N}_{(r,0),h}^i(s, a) = N_{(r),h}(s, a)$

while no synchronization signal do

Collect $(s_{t,h}^i, a_{t,h}^i, r_{t,h}^i, s_{t,h+1}^i)_{1 \leq h \leq H}$ using $\pi_{(r)}$

for $h = H$ **to** 1 **do**

Set $n_{(r,\ell),h}^i(s, a) = n_{(r,\ell-1),h}^i(s, a) + 1_{(s,a)}(s_{t,h}^i, a_{t,h}^i)$ for $(s, a) \in \mathcal{S} \times \mathcal{A}$ and
 $n_{(r,\ell),h}^i(s, a, s') = n_{(r,\ell-1),h}^i(s, a, s') + 1_{(s,a,s')}(s_{t,h}^i, a_{t,h}^i, s_{t,h+1}^i)$ for $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$

Set $\hat{N}_{(r,\ell),h}^i(s, a) = \hat{N}_{(r,\ell-1),h}^i(s, a) + M 1_{(s,a)}(s_{t,h}^i, a_{t,h}^i)$ for $(s, a) \in \mathcal{S} \times \mathcal{A}$

Set $\hat{r}_h^i(s_{t,h}^i, a_{t,h}^i) = r_{t,h}^i$ and $\ell = \ell + 1$

if $(N_{(r),h}(s_{t,h}^i, a_{t,h}^i) \leq \nu(\delta, T) \text{ and } n_{(r,\ell),h}^i(s_{t,h}^i, a_{t,h}^i) > 2n_{(r),h}^i(s_{t,h}^i, a_{t,h}^i))$
or $(N_{(r),h}(s_{t,h}^i, a_{t,h}^i) > \nu(\delta, T) \text{ and } \hat{N}_{(r,\ell),h}^i(s_{t,h}^i, a_{t,h}^i) > N_{(r),h}(s_{t,h}^i, a_{t,h}^i))$ **then**

Send synchronization signal

Set $t = t + \ell$; $n_{(r+1),h}^i(s, a) = n_{(r,\ell),h}^i(s, a)$ and update the transition kernels using (9)

Set $\hat{V}_{(r+1),H+1}(s) = 0$ for all $s \in \mathcal{S}$ and broadcast it to all the clients

for $h = H$ **to** 1 **do**

for agent $i = 1$ **to** M **in parallel do**

Compute $\hat{Q}_{(r+1),h}^i$ using (10)

Send $n_{(r+1),h}^i, \hat{P}_{(r+1),h}^i, \hat{V}_{(r+1),h+1}, \hat{P}_{(r+1),h}^i \hat{V}_{(r+1),h+1}^2$, and $\hat{Q}_{(r+1),h}^i$ to the central server

Compute $N_{(r+1),h}, \hat{Q}_{(r+1),h}, \hat{V}_{(r+1),h}(s)$, and $\pi_{(r+1),h}$ using (6), (11), (13), and (14) and broadcast them to all the clients

Set $r = r + 1$

necessary to compute a Bernstein-like bonus function $b_{(r+1),h}(s, a)$ (see (39) in Appendix for an exact expression). The aggregated Q -value is computed as

$$\hat{Q}_{(r+1),h}(s, a) := \min([\mathcal{T}_{(r+1),h}^\omega + b_{(r+1),h}](s, a), H), \quad (11)$$

with $\mathcal{T}_{(r+1),h}^\omega(s, a) = \sum_{i=1}^M \omega_{(r+1),h}^i(s, a) \hat{Q}_{(r+1),h}^i(s, a)$, and

$$\omega_{(r+1),h}^i(s, a) := \frac{n_{(r+1),h}^i(s, a)}{N_{(r+1),h}(s, a)}. \quad (12)$$

Finally, the central server updates the value function and policy according to the equations

$$\hat{V}_{(r+1),h}(s) := \max_{a \in \mathcal{A}} \hat{Q}_{(r+1),h}(s, a), \quad (13)$$

$$\pi_{(r+1),h}(s) := \arg \max_{a \in \mathcal{A}} \hat{Q}_{(r+1),h}(s, a). \quad (14)$$

These updated values are distributed to all agents, and the process continues for all $h = H, \dots, 1$. Once $h = 1$ is reached, the new epoch $r + 1$ begins.

Communication Complexity. Our algorithmic design shares similarities with previous work on reinforcement learning with low switching cost (Bai et al.,

2019; Zhang et al., 2020; Qiao et al., 2022). In particular, the number of times the local data collection policy changes—known as the switching cost—directly corresponds to the number of communication rounds in our framework, which we define as the communication complexity. In its simplest form, the doubling condition in this context can be expressed as:

$$\exists(s, a, h) : N_{(r,\ell),h}(s, a) > 2N_{(r),h}(s, a), \quad (15)$$

where $N_{(r,\ell),h}(s, a) := \sum_{i=1}^M n_{(r,\ell),h}^i(s, a)$ represents the cumulative count across agents.

However, this condition cannot be directly verified in a federated learning setting, as the value of $N_{(r,\ell),h}(s, a)$ is not accessible to any individual agent. One potential solution is to use a weaker local doubling condition, as defined in (7). However, this approach results in communication complexity scaling linearly with the number of agents M , which is impractical for large-scale federated learning environments. Instead, we propose to construct an estimate of the global counter $\hat{N}_{(r,\ell),h}^i(s, a)$ to serve as a plug-in estimate on the left-hand side of (15). This is the core idea behind the

condition in (8). While such estimates may be inaccurate during the initial stages of training, they become reliable once the number of visits exceeds a threshold $\nu(\delta, T) = \tilde{\mathcal{O}}(\varepsilon_p THM + M)$, defined in (57). At that point, $\hat{N}_{(r,\ell),h}^i(s, a)$ can be effectively used as a plug-in estimate. Using this approach, we establish a bound on the communication complexity of **Fed-UCBVI**.

Lemma 4.1 (Communication Complexity). *With probability at least $1 - \delta$, the number of communication rounds of **Fed-UCBVI** is bounded by*

$$\mathfrak{C}(T) \leq \mathcal{O}(|\mathcal{S}||\mathcal{A}|H \log T + M|\mathcal{S}||\mathcal{A}|H \log \log T + M|\mathcal{S}||\mathcal{A}|H \log(1 + \varepsilon_p T)) ,$$

where logarithmic dependence in $|\mathcal{S}|, |\mathcal{A}|, H, 1/\delta$ and M is ignored.

Sketch of the proof: To prove the result, we consider a fixed triplet (s, a, h) and count how many synchronizations this triplet can trigger. Let $k_{s,a,h}^{\min}$ represent the index of the last round where $N_{(r),h}(s, a) \leq \nu(\delta, T)$. To bound the number of synchronizations that occur between the first round and round $k_{s,a,h}^{\min}$, note that agents send an abort signal only when their local visit count of (s, a) at time h has doubled. This can happen at most $\log_2(\nu(\delta, T))$ times for an individual agent, and for all agents combined, the total is upper bounded by $M \log_2(\nu(\delta, T))$.

Next, we bound the number of synchronizations between round $k_{s,a,h}^{\min}$ and the final round. By applying a Bernstein-type concentration inequality, we can show that the synchronization rule (8) implies the equivalent of (15), although with a coefficient of $8/7$ instead of 2 on the right-hand side. Using a similar argument as above, we obtain $\mathcal{O}(\log(M))$ synchronizations triggered by a single state-action-step triplet. We complete the proof by summing these bounds over all (s, a, h) and using the expression of $\nu(\delta, T)$. \square

A complete proof of Lemma 4.1 is provided in Appendix E. Importantly, we observe that, in the homogeneous setting, the linear dependence on M vanishes. Moreover, we can estimate the communication cost, i.e., the number of bits exchanged, by noting that in each communication round, each agent transmits objects of size at most $|\mathcal{S}||\mathcal{A}|H$.

Computational and Space Complexity. First, we remark that, at all times, agents store objects of size $\mathcal{O}(|\mathcal{S}|^2|\mathcal{A}|H)$. At every episode, agents perform $\mathcal{O}(1)$ operations, while they perform $\mathcal{O}(|\mathcal{S}|^2|\mathcal{A}|H)$ operations at communication times. By Lemma 4.1, we deduce that the computation complexity of this algorithm is $\mathcal{O}(T + |\mathcal{S}|^3|\mathcal{A}|^2H^2 \log T + M|\mathcal{S}||\mathcal{A}|H \log \log T + M|\mathcal{S}||\mathcal{A}|H \log(1 + \varepsilon_p T))$ for all T episodes.

Regret Bound. We now state our main result, which bounds the federated regret of **Fed-UCBVI**.

Theorem 4.1. *With probability at least $1 - \delta$, the following bound on the regret of **Fed-UCBVI** holds*

$$\mathfrak{R}(T) = \tilde{\mathcal{O}} \left(\sqrt{H^3|\mathcal{S}||\mathcal{A}|T/M} + H^3|\mathcal{S}|^2|\mathcal{A}| \right) + \tilde{\mathcal{O}}(TH(H\varepsilon_p + \varepsilon_r)) .$$

We give a sketch of the proof below, and postpone the detailed proof to Appendix D.

In the homogeneous setting, where $\varepsilon_p = \varepsilon_r = 0$, we recover the expected linear speedup in number of agents and achieve a minimax optimal regret bound up to logarithmic factors (see Table 1 for comparisons). In contrast, in the heterogeneous setting, an additional term, that scales linearly with the degree of heterogeneity, emerges. We show in Lemma F.9 in Appendix that this is expected, and comes from the fact that, in some cases, a policy optimal for one agent is sub-optimal by at least $\varepsilon_p H^2$ for another agent. This illustrates the trade-off involved in cooperation between heterogeneous agents: if the degree of heterogeneity is too large, cooperation can become counterproductive.

Sketch of the proof: As a first step of the proof, we reduce the problem of minimizing the federated regret (4) to the problem of minimizing a *common* regret. We introduce the common MDP \mathcal{M}^c as follows

$$\mathcal{M}^c := (\mathcal{S}, \mathcal{A}, H, \{r_h^c := \frac{1}{M} \sum_{i=1}^M r_h^i\}_h, \{\mathbf{P}_h^c\}_h) , \quad (16)$$

where $\{\mathbf{P}_h^c\}_h$ is defined in A-1. We set $\mathcal{V}_h^{c,\pi}$ and $\mathcal{V}_h^{c,*}$ the value-function of a policy π and optimal value-function in \mathcal{M}^c . The common regret is defined as

$$\mathfrak{R}^c(T) := \frac{1}{M} \sum_{t=1}^T \sum_{i=1}^M \mathcal{V}_1^{c,*}(s_{t,1}^i) - \mathcal{V}_1^{c,\pi_t}(s_{t,1}^i) . \quad (17)$$

Adapting the performance-difference lemma of Russo (2019) under A-1, it may be shown that

$$\begin{aligned} \mathfrak{R}(T) &= \max_{\pi} \frac{1}{M} \sum_{t=1}^T \sum_{i=1}^M \mathcal{V}_1^{i,\pi}(s_{t,1}^i) - \mathcal{V}_1^{i,\pi_t}(s_{t,1}^i) \\ &\leq \mathfrak{R}^c(T) + 2T\varepsilon_p H^2 + 2T\varepsilon_r H . \end{aligned}$$

As shown in Lemma F.9, the scaling $\mathcal{O}(T(\varepsilon_p H^2 + \varepsilon_r H))$ with H^2 is unavoidable.

The remainder of the proof involves three key steps, outlined below. The first step focuses on estimating the common transition kernel and introduces the primary technical innovations of this work. It also provides justification of A-1.

Step 1: Estimation of the common transition kernel. First, we prove that the weighted average kernel,

$$\hat{P}_{(r),h}(s'|s,a) := \sum_{i=1}^M \omega_{(r),h}^i(s,a) \cdot \hat{P}_{(r),h}^i(s'|s,a),$$

where weights are defined in (12), forms a well-defined (biased) estimator of the common transition kernel P_h^c using data from all agents. Importantly, neither the agents nor the CS have direct access to this quantity.

The analysis of $\hat{P}_{(r),h}$, under **A-1** poses significant challenges compared to both the generative model setting and the case involving homogeneous agents. To illustrate, the kernel can be reformulated as follows, incorporating all samples from the agents:

$$\hat{P}_{(r),h}(s'|s,a) = \frac{1}{N_{(r),h}(s,a)} \sum_{i=1}^M n_{(r),h}^i(s,a,s').$$

In the homogeneous scenario, where $\varepsilon_p = 0$, as explored in prior work (Zheng et al., 2024, 2025), the estimate is derived from an i.i.d. sequence of categorical random variable samples from $P_h^c(\cdot|s,a)$, simplifying the analysis. Moreover, within the generative model framework, such as in (Jin et al., 2022; Wang et al., 2024), we can ensure an equal sample count from each agent's transition kernel P_h^i , resulting in $\hat{P}_{(r),h}$ as a simple mean of independent biased estimates of the common kernel.

However, in our setting, the estimator $\hat{P}_{(r),h}$ incorporates a random and non-stationary number of samples from each agent, making standard techniques of conditioning on a total sample size $N_{(r),h}(s,a)$ inapplicable. Using union-bound arguments to account for the variability in sample sizes across agents results in an exponential number of configurations with respect to M , constraining any possibility of linear speed-up.

Using **A-1**, every kernel P_h^i is a mixture of P_h^c and $P_h^{\text{ind},i}$. The samples obtained by agent i as a mixture of samples coming from the two latter kernels: sample $s_{t,h}^i$ is with probability $1 - \varepsilon_p$ generated from P_h^c , and with probability ε_p from $P_h^{\text{ind},i}$. We define a *virtual estimate of the common kernel*, $\hat{P}_{(r),h}^c$, for each communication round r , representing the estimate we would have obtained if all samples were drawn solely from P_h^c . This estimate is subject to a bias resulting from the heterogeneity.

$$\left\| (\hat{P}_{(r),h} - \hat{P}_{(r),h}^c)(\cdot|s,a) \right\|_1 = \tilde{O}(\varepsilon_p + \frac{1}{N}), \quad (18)$$

where $N = N_{(r),h}(s,a)$, that holds for any $(r,s,a,h) \in [\mathfrak{C}(T)] \times \mathcal{S} \times \mathcal{A} \times \mathcal{S}$.

Step 2: Optimism. In our setting, our estimates are not optimistic due to the presence of heterogeneity; however, we can show the analog of the required

properties $\hat{\mathcal{V}}_{(r),h}(s) \geq \mathcal{V}_h^{c,*}(s) - (2\varepsilon_r + 3\varepsilon_p H)(H + 1 - h)$, for any r and $(s,h) \in \mathcal{S} \times [H]$. The key ingredients are concentration inequalities, an inequality (18) and Lemma 14 of (Zhang et al., 2021); see also Lemma D.1 in Appendix. The proof is carried out by induction on h . Applying the update rule (20), combined with a simple rearranging of the terms, yields

$$\begin{aligned} \hat{Q}_{(r),h}(s,a) &\geq \underbrace{Q_h^{c,*}(s,a) + (\hat{P}_{(r),h}^c - P_h^c) \mathcal{V}_{h+1}^{c,*}(s,a)}_{\text{(IV): concentration error}} \\ &\quad + \underbrace{\sum_{i=1}^M \omega_{(r),h}^i(s,a) \hat{r}_h^i(s,a) - \frac{1}{M} \sum_{i=1}^M r_h^i(s,a)}_{\text{(I): reward heterogeneity error}} \\ &\quad + \underbrace{\hat{P}_{(r),h}(\hat{\mathcal{V}}_{(r),h+1}(s,a) - \mathcal{V}_{h+1}^{c,*}(s,a))}_{\text{(II): correction error}} \\ &\quad + \underbrace{(\hat{P}_{(r),h} - \hat{P}_{(r),h}^c) \mathcal{V}_{h+1}^{c,*}(s,a)}_{\text{(III): transition heterogeneity error}} + b_{(r),h}(s,a). \end{aligned}$$

Terms (II) and (IV), which represent the correction and the concentration errors, are standard and are controlled using respectively induction hypothesis, Lemma D.1 and standard deviation inequalities. We control (I) by applying **A-2** and noticing that the convex combination of $\hat{r}_h^i(s,a)$ is also a convex combination of the true rewards $r_h^i(s,a)$. Finally, to control (III) we combine Holder's inequality and inequality (18). An appropriate choice of the exploration bonus concludes the statement.

Step 3: Bounding the regret. For each quantity indexed by the number of communication rounds r (e.g. $\hat{\mathcal{V}}_{(r),h}$), we introduce a corresponding quantity indexed by the episode number t (e.g. $\hat{\mathcal{V}}_{t,h}$), defined as the value of the former at the last communication round before t (see (27) in Appendix for formal definitions). Next, following the approach of Azar et al. (2017), we define $\delta_{t,h}^i = \hat{\mathcal{V}}_{t,h}(s_{t,h}^i) - \mathcal{V}_1^{c,\pi_t}(s_{t,h}^i)$ and analyze this term independently

$$\begin{aligned} \delta_{t,h}^i &\leq \delta_{t,h+1}^i + \underbrace{[\hat{P}_{t,h} - \hat{P}_{t,h}^c] \hat{\mathcal{V}}_{t,h+1}(s_{t,h}^i, a_{t,h}^i)}_{\text{(A): heterogeneity error}} + \zeta_{t,h}^i \\ &\quad + \underbrace{[\hat{P}_{t,h}^c - P_h^c] [\hat{\mathcal{V}}_{t,h+1} - \mathcal{V}_{h+1}^{c,*}](s_{t,h}^i, a_{t,h}^i)}_{\text{(B): correction error}} + 2\varepsilon_r \\ &\quad + \underbrace{[\hat{P}_{t,h}^c - P_h^c] \mathcal{V}_{h+1}^{c,*}(s_{t,h}^i, a_{t,h}^i)}_{\text{(C): concentration error}} + b_{t,h}(s_{t,h}^i, a_{t,h}^i) \\ &\quad + \underbrace{[P_h^c - P_h^i] [\hat{\mathcal{V}}_{t,h+1} - \mathcal{V}_{h+1}^{c,\pi_t}](s_{t,h}^i, a_{t,h}^i)}_{\text{(D): heterogeneity error}}, \end{aligned}$$

where $\zeta_{t,h}^i$ is a martingale term defined in (54). The analysis of (C) and $\zeta_{t,h}^i$ is standard in the literature.

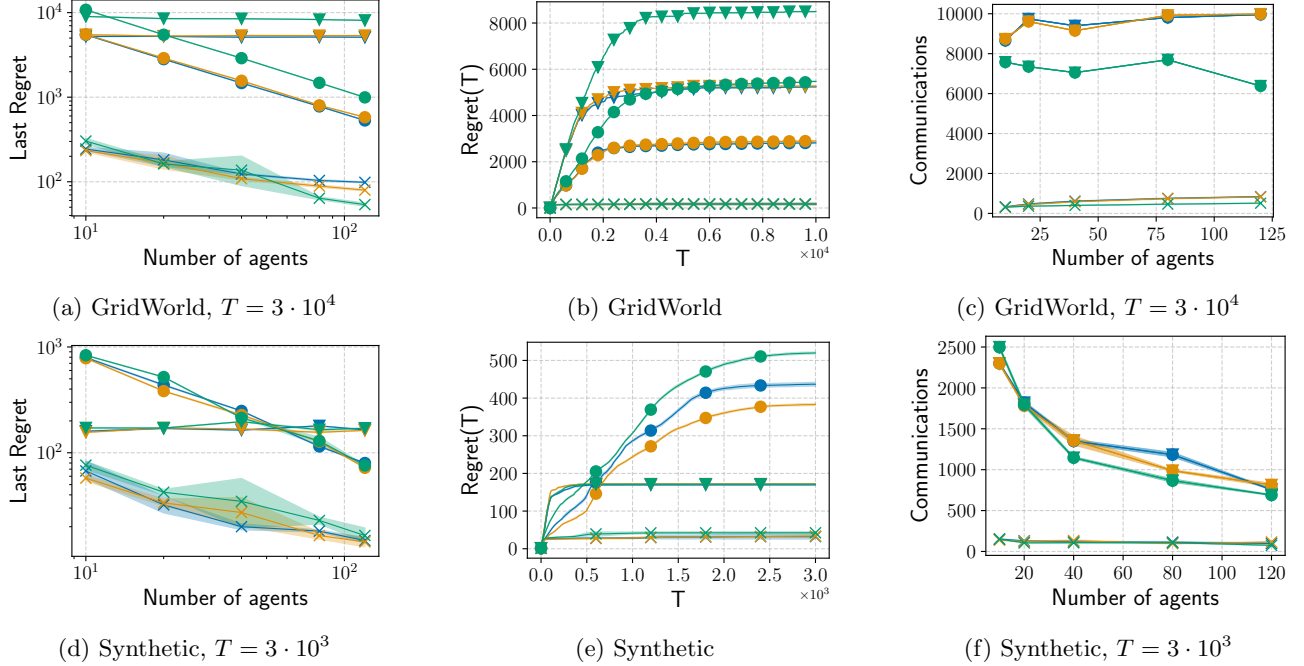


Figure 1: Comparison of different properties of **Fed-UCBVI** (represented by crosses), **FedQ-Bernstein** (represented by circles) and **FedQ-Advantage** (represented by triangles) in two environments: **synthetic** (below) and **GridWorld** (above). From left to right: plot of the common regret (*lower is better*), as a function of M for different ε_p in a log-log scale; plot of the common regret (*lower is better*) for $M = 20$ agents as a function of T for different ε_p ; plot of the number of communication (*lower is better*) as a function of M for different ε_p .

To bound (A) we employ a combination of (18) and Holder’s inequality. The bound on (D) also combines Holder’s inequality and Lemma F.1. The standard recursion argument concludes the proof. \square

5 EXPERIMENTS

In this section, we study the empirical performance of **Fed-UCBVI**, and compare it with the **FedQ-Bernstein** (Zheng et al., 2024) and **FedQ-Advantage** (Zheng et al., 2025) algorithms on two environments.¹

Environments. We consider two environments specifically designed to satisfy A-1 and A-2. In both environments, transitions are defined using two distinct kernels: with probability $1 - \varepsilon_p$, the agent follow the global kernel, and with probability ε_p , it follows an individualized kernel. The first environment is based on **GridWorld** (Domingues et al., 2021a), where the agent navigates a grid to reach a target. Upon reaching the target, the agent receives a reward of +1; otherwise, the reward is 0. At each step, the agent selects one of four possible directions (up, down, left, or right). Under the global transition kernel, the agent moves to the intended square with a probability of 0.8, and to a random neighboring square with the remaining probability. In the individual transition kernels, the agent’s

movement to neighboring squares follows a probability distribution unique to each agent. We use a 3×3 grid with a wall located at coordinate (1, 1), resulting in $|\mathcal{S}| = 8$ possible states. The planning horizon is set to $H = 10$, with the agent starting at coordinate (0, 0) and aiming to reach the target at (2, 2).

The second environment is a **synthetic** setting, modeled after Zheng et al. (2024), with $|\mathcal{S}| = 5$, $|\mathcal{A}| = 5$, and $H = 5$. All agents share the same reward function $r_h(s, a)$, with rewards drawn uniformly from $[0, 1]$ for each $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$. For each s, a, h , the common and individual transition kernels are drawn uniformly at random from the $|\mathcal{S}|$ -dimensional simplex.

In all results, we report the common regret instead of the federated regret to simplify computations. Experiments were conducted on a computer with an Intel Xeon 6534 and 196GB RAM. We report the average over 5 runs and the standard deviation in all the plots. The code is provided in the supplementary material.

Fed-UCBVI has linear speed-up. In Figures 1a and 1d, we evaluate the regret after T iterations of training with varying numbers of agents M across different levels of heterogeneity. As shown in Theorem 4.1, the regret decreases as M increases. Notably, this trend persists even in high-heterogeneity settings, highlighting the robust empirical performance

¹Our code is available online on GitHub: <https://github.com/Labbi-Safwan/Fed-UCBVI>

of our approach. This phenomenon also holds for FedQ-Bernstein. However, FedQ-Advantage does not have linear speed-up.

Impact of Heterogeneity. In Figures 1b and 1e, we present the regret of Fed-UCBVI for various values of ε_p . Fed-UCBVI’s regret is significantly lower than that of FedQ-Bernstein and FedQ-Advantage, reflecting similar performance gaps as observed in the single-agent setting. Moreover, as predicted by our theoretical analysis, increasing ε_p only incurs a slight increase in Fed-UCBVI’s regret, due to the additional term scaling linearly with T .

Fed-UCBVI’s communication complexity is small. In Figures 1c and 1f, we observe that the communication complexity of Fed-UCBVI is significantly lower than the number of iterations T and increases only marginally with the number of agents M . This aligns with the results of Lemma 4.1. In contrast, FedQ-Bernstein and FedQ-Advantage exhibit very similar and consistently high communication complexity. The reduced communication in Fed-UCBVI results from our novel method for triggering communication rounds based on local estimates of global counters, validating the effectiveness of this approach.

6 CONCLUSION

In this paper, we presented Fed-UCBVI, a federated reinforcement learning method based on a new aggregation strategy that reduces communication cost and handles heterogeneous agents. We introduced a novel measure of heterogeneity, under which we provide a formal analysis of Fed-UCBVI’s regret, showing that it nearly matches minimax optimal regret bounds. To our knowledge, this is the first federated regret analysis with guarantees in heterogeneous environments. Furthermore, our method provably removes the linear dependence of the communication complexity on $M \log T$. A promising direction for future work is to reduce the communication cost further, by developing new methods that correct for heterogeneity.

ACKNOWLEDGEMENTS

The work of S. Labbi, L. Mancini and P. Mangold has been supported by Technology Innovation Institute (TII), project Fed2Learn. The work of D. Tiapkin has been supported by the Paris Île-de-France Région in the framework of DIM AI4IDF. The work of E. Moulines has been partly funded by the European Union (ERC-2022-SYG-OCEAN-101071601). Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Research Council Executive Agency. Neither the European Union

nor the granting authority can be held responsible for them.

References

- Azar, M. G., Osband, I., and Munos, R. (2017). Minimax regret bounds for reinforcement learning. In *International conference on machine learning*, pages 263–272. PMLR.
- Bai, Y., Xie, T., Jiang, N., and Wang, Y.-X. (2019). Provably efficient q-learning with low switching cost. *Advances in Neural Information Processing Systems*, 32.
- Chen, T., Zhang, K., Giannakis, G. B., and Başar, T. (2021). Communication-efficient policy gradient methods for distributed reinforcement learning. *IEEE Transactions on Control of Network Systems*, 9(2):917–929.
- Chen, Y., Zhang, X., Zhang, K., Wang, M., and Zhu, X. (2023). Byzantine-robust online and offline distributed reinforcement learning. In *International Conference on Artificial Intelligence and Statistics*, pages 3230–3269. PMLR.
- Dann, C., Lattimore, T., and Brunskill, E. (2017). Unifying pac and regret: Uniform pac bounds for episodic reinforcement learning. *Advances in Neural Information Processing Systems*, 30.
- Domingues, O. D., Flet-Berliac, Y., Leurent, E., Ménard, P., Shang, X., and Valko, M. (2021a). rlberry - A Reinforcement Learning Library for Research and Education.
- Domingues, O. D., Ménard, P., Kaufmann, E., and Valko, M. (2021b). Episodic reinforcement learning in finite mdps: Minimax lower bounds revisited. In *Algorithmic Learning Theory*, pages 578–598. PMLR.
- Domingues, O. D., Ménard, P., Pirota, M., Kaufmann, E., and Valko, M. (2021c). Kernel-based reinforcement learning: A finite-time analysis. In *International Conference on Machine Learning*, pages 2783–2792. PMLR.
- Doob, J. (1953). *Stochastic Processes*. Probability and Statistics Series. Wiley.
- Huber, P. J. (1964). Robust Estimation of a Location Parameter. *The Annals of Mathematical Statistics*, 35(1):73 – 101.
- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. (2018). Is q-learning provably efficient? *Advances in neural information processing systems*, 31.
- Jin, H., Peng, Y., Yang, W., Wang, S., and Zhang, Z. (2022). Federated reinforcement learning with environment heterogeneity. In *International Conference*

on *Artificial Intelligence and Statistics*, pages 18–37. PMLR.

Jonsson, A., Kaufmann, E., Ménard, P., Darwiche Domingues, O., Leurent, E., and Valko, M. (2020). Planning in markov decision processes with gap-dependent sample complexity. *Advances in Neural Information Processing Systems*, 33:1253–1263.

Lan, G., Han, D.-J., Hashemi, A., Aggarwal, V., and Brinton, C. G. (2024). Asynchronous federated reinforcement learning with policy gradient updates: Algorithm design and convergence analysis. *arXiv preprint arXiv:2404.08003*.

Li, G., Shi, L., Chen, Y., Gu, Y., and Chi, Y. (2021). Breaking the sample complexity barrier to regret-optimal model-free reinforcement learning. *Advances in Neural Information Processing Systems*, 34:17762–17776.

Mangold, P., Samsonov, S., Labbi, S., Levin, I., Alami, R., Naumov, A., and Moulines, E. (2025). Scaffisa: Taming heterogeneity in federated linear stochastic approximation and td learning. *Advances in Neural Information Processing Systems*, 37:13927–13981.

Maurer, A. and Pontil, M. (2009). Empirical bernstein bounds and sample variance penalization. *arXiv preprint arXiv:0907.3740*.

McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR.

Qi, J., Zhou, Q., Lei, L., and Zheng, K. (2021). Federated reinforcement learning: techniques, applications, and open challenges. *Intelligence & Robotics*, 1(1).

Qiao, D., Yin, M., Min, M., and Wang, Y.-X. (2022). Sample-efficient reinforcement learning with loglog (t) switching cost. In *International Conference on Machine Learning*, pages 18031–18061. PMLR.

Ross, S. and Bagnell, D. (2010). Efficient reductions for imitation learning. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 661–668. JMLR Workshop and Conference Proceedings.

Russo, D. (2019). Worst-case regret bounds for exploration via randomized value functions. *Advances in Neural Information Processing Systems*, 32.

Salgia, S. and Chi, Y. (2024). The sample-communication complexity trade-off in federated q-learning. In Globerson, A., Mackey, L., Belgrave, D.,

Fan, A., Paquet, U., Tomczak, J., and Zhang, C., editors, *Advances in Neural Information Processing Systems*, volume 37, pages 39694–39747. Curran Associates, Inc.

Sobel, M. J. (1982). The variance of discounted markov decision processes. *Journal of Applied Probability*, 19(4):794–802.

Sutton, R. S. and Barto, A. G. (2018). *Reinforcement Learning: An Introduction*. A Bradford Book, Cambridge, MA, USA.

Talebi, M. S. and Maillard, O.-A. (2018). Variance-aware regret bounds for undiscounted reinforcement learning in mdps. In *Algorithmic Learning Theory*, pages 770–805. PMLR.

Tiapkin, D., Belomestny, D., Calandriello, D., Moulines, E., Munos, R., Naumov, A., Perrault, P., Tang, Y., Valko, M., and Menard, P. (2023). Fast rates for maximum entropy exploration. In *International Conference on Machine Learning*, pages 34161–34221. PMLR.

Wang, M., Yang, P., and Su, L. (2024). On the convergence rates of federated q-learning across heterogeneous environments. In *International Workshop on Federated Foundation Models in Conjunction with NeurIPS 2024*.

Zanette, A. and Brunskill, E. (2019). Tighter problem-dependent regret bounds in reinforcement learning without domain knowledge using value function bounds. In *International Conference on Machine Learning*, pages 7304–7312. PMLR.

Zhang, C., Wang, H., Mitra, A., and Anderson, J. (2024a). Finite-time analysis of on-policy heterogeneous federated reinforcement learning. In *The Twelfth International Conference on Learning Representations*.

Zhang, Z., Chen, Y., Lee, J. D., and Du, S. S. (2024b). Settling the sample complexity of online reinforcement learning. In *The Thirty Seventh Annual Conference on Learning Theory*, pages 5213–5219. PMLR.

Zhang, Z., Ji, X., and Du, S. (2021). Is reinforcement learning more difficult than bandits? a near-optimal algorithm escaping the curse of horizon. In *Conference on Learning Theory*, pages 4528–4531. PMLR.

Zhang, Z., Zhou, Y., and Ji, X. (2020). Almost optimal model-free reinforcement learning via reference-advantage decomposition. *Advances in Neural Information Processing Systems*, 33:15198–15207.

Zheng, Z., Gao, F., Xue, L., and Yang, J. (2024). Federated q-learning: Linear regret speedup with low

communication cost. In *The Twelfth International Conference on Learning Representations*.

Zheng, Z., Zhang, H., and Xue, L. (2025). Federated q-learning with reference-advantage decomposition: Almost optimal regret and logarithmic communication cost. In *The Thirteenth International Conference on Learning Representations*.

Zhuo, H. H., Feng, W., Lin, Y., Xu, Q., and Yang, Q. (2019). Federated deep reinforcement learning. *arXiv preprint arXiv:1901.08277*.

Checklist

1. For all models and algorithms presented, check if you include:

- (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model.

Answer: Yes

Justification: The mathematical setting, the assumptions, and the algorithm are extensively described in Section 3 and Section 4.

- (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm.

Answer: Yes

Justification: All properties associated with the algorithm are detailed in Section 4, along with detailed proof.

- (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries.

Answer: Yes

2. For any theoretical claim, check if you include:

- (a) Statements of the full set of assumptions of all theoretical results.

Answer: Yes

Justification: All the assumptions required to conduct the proof are described in Section 3.

- (b) Complete proofs of all theoretical results.

Answer: Yes

Justification: A proof sketch for all the theorems is provided in Section 4, with the full proofs included in the appendix.

- (c) Clear explanations of any assumptions.

Answer: Yes

Justification: The assumptions are clearly stated and explained in Section 3.2.

3. For all figures and tables that present empirical results, check if you include:

- (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL).

Answer: Yes, our code is available online on GitHub: <https://github.com/Labbi-Safwan/Fed-UCBVI>

- (b) All the training details (e.g., data splits, hyperparameters, how they were chosen).

Answer: Yes

- (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times).

Answer: Yes

- (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider).

Answer: Yes

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:

- (a) Citations of the creator If your work uses existing assets.

Answer: Yes

- (b) The license information of the assets, if applicable.

Answer: Yes

- (c) New assets either in the supplemental material or as a URL, if applicable.

Answer: Yes

- (d) Information about consent from data providers/curators.

Answer: Not Applicable

- (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content.

Answer: Not Applicable

5. If you used crowdsourcing or conducted research with human subjects, check if you include:

- (a) The full text of instructions given to participants and screenshots.

Answer: Not Applicable

- (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable.

Answer: Not Applicable

- (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation.

Answer: Not Applicable

Supplementary Materials

A NOTATION

For clarity, we summarize here the notations that we use

Symbols	Meaning	Definition
$\mathfrak{C}(T)$	Number of communication rounds performed in average	Section 3
$\mathfrak{R}(T)$	Federated regret of the algorithm	Equation (4)
\mathcal{S}	State space	Section 3
\mathcal{A}	Action space	Section 3
M	Number of agents	Section 3
T	Total number of collected episodes per agent	Section 3
H	Length of an episode	Section 3
R_{\max}	Maximum number of communication rounds	Equation (58)
\mathbf{P}_h^i	Transition kernel at step h of agent i	Section 3
\mathbf{P}_h^c	Common transition kernel at step h	A-1
$\mathbf{P}_h^{\text{ind},i}$	Individual transition kernel at step h	A-1
ε_p	Degree of heterogeneity on the transition kernels	A-1
ε_r	Degree of heterogeneity on the rewards	A-2
\mathbf{r}_h^i	Reward at step h of agent i	Section 3
\mathbf{r}_h^c	Reward function of the common MDP	Equation (16)
$Q_h^{i,\pi}$	Q-function of a policy π at step h of agent i	Equation (2)
$\mathcal{V}_h^{i,\pi}$	Value function of a policy π at step h in the i -th MDP	Equation (2)
$Q_h^{i,*}$	Optimal Q-function at step h of agent i in the i -th environment	Equation (3)
$\mathcal{V}_h^{i,*}$	Optimal value function at step h of agent i in the i -th environment	Equation (3)
$\nu(\delta, T)$	Threshold for defining the condition on initiating the aggregation signal	Equation (57)
$\hat{\mathbf{P}}_{(r),h}^i$	Estimated transition kernel during the round r by agent i at step h	Equation (9)
$\hat{\mathbf{P}}_{(r),h}^c$	Virtual estimate of common transition kernel by agent i at step h	Equation (24)
$\hat{\mathbf{r}}_h^i$	Estimated reward at step h of agent i	Fed-UCBVI
$n_{(r,\ell),h}^i$	Local counter of the cumulative number of visits at the level of agent i	Fed-UCBVI
$N_{(r),h}$	Global counter of the cumulative number of visits over all the agents	Equation (6)
$\hat{N}_{(r,\ell),h}^i$	Local estimator of agent i of the <i>true</i> cumulative number of visits	Fed-UCBVI
$b_{(r),h}$	Bonus function used in round r and step H	Equation (39)
$\hat{Q}_{(r),h}^i(s, a)$	Estimator of the Q-function at the level of agent i	Equation (10)
$\hat{\mathbf{P}}_{(r),h}$	Weighted average of $\{\hat{\mathbf{P}}_{(r),h}^i\}_i$ during the round r at step h	Equation (24)
$\hat{Q}_{(r),h}(s, a)$	Global estimator of the Q-function	Equation (11)
$\hat{\mathcal{V}}_{(r),h}(s)$	Global estimator of the value function	Equation (13)
$\text{Var}_{\hat{\mathbf{P}}_{(r),h}}(f)(s, a)$	Variance of a function f with respect to $\hat{\mathbf{P}}_{(r),h}(\cdot s, a)$	Equation (19)
$\text{Var}_{\hat{\mathbf{P}}_{(r),h}^c}(f)(s, a)$	Variance of a function f with respect to $\hat{\mathbf{P}}_{(r),h}^c(\cdot s, a)$	Equation (19)

Table 2: Summary of the notations.

Let (X, \mathcal{X}) be a measurable space. For any probability measures P and Q on (X, \mathcal{X}) , and for any $f : X \rightarrow \mathbb{R}$ we define

$$Pf := \mathbb{E}_{s \sim P}[f(s)], \quad \text{Var}_P f := \mathbb{E}_{s \sim P}[(f(s) - Pf)^2]. \quad (19)$$

For any probability measures P and Q on (X, \mathcal{X}) , the Kullback-Leibler divergence $\text{KL}(P\|Q)$ is given by

$$\text{KL}(P\|Q) := \begin{cases} \mathbb{E}_P \left[\log \frac{dP}{dQ} \right], & P \ll Q, \\ +\infty, & \text{otherwise} \end{cases}.$$

Let A be an element of the σ -algebra \mathcal{X} . We define the indicator function of A as

$$\begin{aligned} 1_A(\cdot) : X &\longrightarrow \{0, 1\} \\ x &\longmapsto \begin{cases} 1 & \text{if } x \in A, \\ 0, & \text{otherwise} \end{cases} \end{aligned}$$

We define the indicator function of an element $x \in X$ as

$$\begin{aligned} 1_x(\cdot) : X &\longrightarrow \{0, 1\} \\ y &\longmapsto \begin{cases} 1 & \text{if } x = y, \\ 0, & \text{otherwise} \end{cases} \end{aligned}$$

We write $f(|\mathcal{S}|, |\mathcal{A}|, H, T, M) = \mathcal{O}(g(|\mathcal{S}|, |\mathcal{A}|, H, T, M, \delta))$ if there exists $S_0, A_0, H_0, T_0, \delta_0$ and a constant C such that for any $|\mathcal{S}| \geq S_0, |\mathcal{A}| \geq A_0, H \geq H_0, T \geq T_0$, and $\delta \leq \delta_0$, we have $f(|\mathcal{S}|, |\mathcal{A}|, H, T, M) \leq C \cdot g(|\mathcal{S}|, |\mathcal{A}|, H, T, M, \delta)$. We say that $f(|\mathcal{S}|, |\mathcal{A}|, H, T, M) = \tilde{\mathcal{O}}(g(|\mathcal{S}|, |\mathcal{A}|, H, T, M, \delta))$ if in the previous bound C is a poly-logarithmic function with respect to the variables $|\mathcal{S}|, |\mathcal{A}|, H, T, M, \delta$.

For $a \in \mathbb{N}$, define $[a]$ as the set of all natural numbers from 1 to a :

$$[a] := \{k \in \mathbb{N} \mid 1 \leq k \leq a\}.$$

Additionally, for $(a, b) \in \mathbb{N} \times \bar{\mathbb{N}}$, where $\bar{\mathbb{N}} = \mathbb{N} \cup \{+\infty\}$, such that $a \leq b$, define the set $\llbracket a, b \rrbracket$ as the set of all natural numbers between a and b , inclusive:

$$\llbracket a, b \rrbracket := \{k \in \mathbb{N} \mid a \leq k \leq b\}.$$

B PSEUDO CODE

For clarity of exposition, we provide the complete pseudo-code of the server-side and client-side algorithms in Algorithm 2 and Algorithm 3.

C CONCENTRATION EVENTS

Before we proceed, let us define several essential quantities.

Change of epoch notation We notice that the set of all regular episodes $t \in [T]$ is separated into a sequence of different *random* epochs E_1, E_2, \dots . To define them properly, let us define the epoch-changing timestamps as follows

$$T_1 := 0, \quad T_{r+1} := \min\{t > T_r \mid \text{Sync}_r(t) = \text{True}\}. \quad (20)$$

where the epoch-switching predicate is defined as

$$\text{Sync}_r(t) = \begin{cases} \exists i \in [M] : n_{(r,\ell),h}^i(s_{t,h}^i, a_{t,h}^i) \geq 2n_{(r),h}^i(s_{t,h}^i, a_{t,h}^i) & \text{if } N_{(r),h}(s_{t,h}^i, a_{t,h}^i) < \nu(\delta, T) \\ \exists i \in [M] : \tilde{N}_{(r,\ell),h}^i(s_{t,h}^i, a_{t,h}^i) \geq 2N_{(r),h}(s_{t,h}^i, a_{t,h}^i) & \text{if } N_{(r),h}(s_{t,h}^i, a_{t,h}^i) \geq \nu(\delta, T) \end{cases}, \quad (21)$$

for $\ell = t - T_r - 1$ and $\nu(\delta, T)$ is defined in (57). In particular, this condition exactly corresponds to the synchronization condition used by Fed-UCBVI. Then, the epoch E_r is defined as $E_r := \llbracket T_r + 1; T_{r+1} \rrbracket$. In particular, for any $t \in [T]$, we define r_t as a unique index r such that $t \in E_r$:

$$r_t = \min\{r \geq 1 \mid t > T_r\}. \quad (22)$$

Algorithm 2: Fed-UCBVI (Central Server)

Initialize: $t = 1$, $r = 1$, $N_{(1),h}(s, a) = 0$, $\hat{Q}_{(1),h}(s, a) = \hat{V}_{(1),h}(s) = H$ for all $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$,
 $\pi_{(1)} = \{\pi_{(1),h}\}_{1 \leq h \leq H}$ an arbitrary deterministic policy
while $t \leq T$ **do**
 Broadcast $\pi_{(r)} = \{\pi_{(r),h}\}_{1 \leq h \leq H}$, $\{N_{(r),h}\}_{1 \leq h \leq H}$, r , and t to all clients;
 Wait until receiving the synchronization signal and an updated episode number t and forward the
 abortion signal to all clients;
 Set $\hat{V}_{(r+1),H+1}(s) = 0$ for all $s \in \mathcal{S}$ and send it to all clients;
 for $h = H$ **to** 1 **do**
 Receive $\{\hat{Q}_{(r+1),h}^i\}_i$, $\{n_{(r+1),h}^i\}_i$, $\{\hat{P}_{(r+1),h}^i \hat{V}_{(r+1),h+1}(s, a)\}_{s,a,i}$, and $\{\hat{P}_{(r+1),h}^i \hat{V}_{(r+1),h+1}^2(s, a)\}_{s,a,i}$
 from the different clients;
 for $(s, a) \in \mathcal{S} \times \mathcal{A}$ **do**
 Compute $N = N_{(r+1),h}(s, a) = \sum_{i=1}^M n_{(r+1),h}^i(s, a)$;
 Set $n^i = n_{(r+1),h}^i(s, a)$ for $i \in [M]$;
 Set $V := \text{Var}_{\hat{P}_{(r),h}}(\hat{V}_{(r),h+1})(s, a) = \frac{1}{N} \sum_{i=1}^M n^i \hat{P}_{(r+1),h}^i \hat{V}_{(r+1),h+1}^2(s, a)$

$$- \left(\frac{1}{N} \sum_{i=1}^M n^i \hat{P}_{(r+1),h}^i \hat{V}_{(r+1),h+1}(s, a) \right)^2$$

 Compute $b_{(r),h}(s, a) = \begin{cases} \frac{28\beta^*(\delta)H+11\beta^c(\delta,N)}{N} + \sqrt{\frac{8\beta^*(\delta)}{N} \cdot V}, & N \geq 2, \\ H, & N \leq 1; \end{cases}$
 Set $\hat{Q}_{(r+1),h}(s, a) = \begin{cases} \min(\sum_{i=1}^M \frac{n^i}{N} \hat{Q}_{(r+1),h}^i(s, a) + b_{(r+1),h}(s, a), H) & \text{if } N > 0, \\ H & \text{otherwise;} \end{cases}$
 for $s \in \mathcal{S}$ **do**
 Compute $\hat{V}_{(r+1),h}(s) = \max_{a \in \mathcal{A}} \hat{Q}_{(r+1),h}(s, a)$;
 Compute $\pi_{(r+1),h}(s) = \arg \max_{a \in \mathcal{A}} \hat{Q}_{(r+1),h}(s, a)$;
 Broadcast $\hat{V}_{(r+1),h}$ to all clients;
 Set $r = r + 1$;
 Send a signal to inform the clients of the end of training.

Definitions First of all, let us recall that by [A-1](#) the transition kernel P_h^i for the agent i is a mixture of common kernel P_h^c and individual kernel $P_h^{\text{ind},i}$, thus any sample $s_{t,h+1}^i \sim P_h^i(s_{t,h}^i, a_{t,h}^i)$ for $(t, h, i) \in [T] \times [H] \times [M]$ can be represented via the following experiment

$$s_{t,h+1}^i = \begin{cases} s_{t,h+1}^{c,i} \sim P_h^c(s_{t,h}^i, a_{t,h}^i), & \xi_{t,h}^i = 0, \\ s_{t,h+1}^{\text{ind},i} \sim P_h^{\text{ind},i}(s_{t,h}^i, a_{t,h}^i), & \xi_{t,h}^i = 1, \end{cases} \quad (23)$$

where $\xi_{t,h}^i \sim \text{Ber}(\varepsilon_p)$ is a choice of component of the mixture. Using this representation, we can define a *virtual estimate of the common kernel* for a step t as follows

$$\hat{P}_{(r),h}^c(s'|s, a) := \frac{1}{N_{(r),h}(s, a)} \sum_{i=1}^M \sum_{t=1}^{T_r} \mathbf{1}_{(s,a,s')}(s_{t,h}^i, a_{t,h}^i, s_{t,h+1}^{c,i}), \quad (24)$$

where T_r is defined in [\(20\)](#). We emphasize that $\hat{P}_{(r),h}^c$ is never computed explicitly by the algorithm since the values of $\xi_{t,h}^i$ are never observed, however we are very interested in the analysis of it.

Algorithm 3: Fed-UCBVI (i-th Client Side)

Initialize: $n_{(1,0),h}^i(s, a) = 0$, $\hat{r}_h^i(s, a) = 0$, $n_{(1,0),h}^i(s, a, s') = 0$ for all $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]$;
 Compute $\nu(\delta, T) = 14\varepsilon_p THM + 182M\beta^c(\delta, T)$;
while *signal of end of training not received* **do**
 Receive $\{\pi_{(r),h}\}_{1 \leq h \leq H}$, $(N_{(r),h})_{1 \leq h \leq H}$, r , and t from the central server;
 Set $\ell = 1$;
 Set $\hat{N}_{(r,0),h}^i(s, a) = N_{(r),h}(s, a)$ for all $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$;
 Set $n_{(r,0),h}^i(s, a) = n_{(r),h}^i(s, a)$ for all $(s, a, s', h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times [H]$;
 while *no synchronization signal from central server and* $t \leq T$ **do**
 synchronize = **False**;
 while *synchronize = False* **do**
 Collect a new trajectory $(s_{t,h}^i, a_{t,h}^i, r_{t,h}^i)_{1 \leq h \leq H}$ using the policy $\pi_{(r)}$;
 for $h = 1$ **to** H **do**
 Set $\hat{r}_h^i(s_{t,h}^i, a_{t,h}^i) = r_{t,h}^i$;
 Set $n_{(r,\ell),h}^i(s, a) = n_{(r,\ell-1),h}^i(s, a) + \mathbf{1}_{(s,a)}(s_{t,h}^i, a_{t,h}^i)$ for $(s, a) \in \mathcal{S} \times \mathcal{A}$;
 $n_{(r,\ell),h}^i(s, a, s') = n_{(r,\ell-1),h}^i(s, a, s') + \mathbf{1}_{(s,a,s')}(s_{t,h}^i, a_{t,h}^i, s_{t,h+1}^i)$ for $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$;
 Set $\hat{N}_{(r,\ell),h}^i(s, a) = \hat{N}_{(r,\ell-1),h}^i(s, a) + M\mathbf{1}_{(s,a)}(s_{t,h}^i, a_{t,h}^i)$ for $(s, a) \in \mathcal{S} \times \mathcal{A}$;
 if $N_{(r),h}(s_{t,h}^i, a_{t,h}^i) < \nu(\delta, T)$ **and** $n_{(r,l),h}^i(s_{t,h}^i, a_{t,h}^i) \geq 2n_{(r),h}^i(s_{t,h}^i, a_{t,h}^i)$ **then**
 synchronize = **True**;
 else if $N_{(r),h}(s_{t,h}^i, a_{t,h}^i) \geq \nu(\delta, T)$ **and** $\hat{N}_{(r,l),h}^i(s_{t,h}^i, a_{t,h}^i) \geq 2N_{(r),h}(s_{t,h}^i, a_{t,h}^i)$ **then**
 synchronize = **True**;
 Set $\ell = \ell + 1$ and $t = t + 1$;
 ;
 Send an abortion signal and an episode number t to the central server;
 Set $n_{(r+1),h}^i(s, a) = n_{(r,\ell),h}^i(s, a)$ and $n_{(r+1),h}^i(s, a, s') = n_{(r,\ell),h}^i(s, a, s')$ for all $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$;
 Set $\hat{\mathbf{P}}_{(r+1),h}^i(s'|s, a) = \begin{cases} \frac{n_{(r+1),h}^i(s, a, s')}{n_{(r+1),h}^i(s, a)} & \text{if } n_{(r+1),h}^i(s, a) > 0 \\ \frac{1}{|\mathcal{S}|} & \text{else;} \end{cases}$
 for $h = H$ **to** 1 **do**
 Receive $\hat{\mathcal{V}}_{(r+1),h+1}$ from the central server;
 Compute $\hat{\mathcal{Q}}_{(r+1),h}^i(s, a) = \hat{r}_h^i(s, a) + \hat{\mathbf{P}}_{(r+1),h}^i \hat{\mathcal{V}}_{(r+1),h+1}(s, a)$ for all $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$;
 Send $\hat{\mathcal{Q}}_{(r+1),h}^i$, $n_{(r+1),h}^i$, $\{\hat{\mathbf{P}}_{(r+1),h}^i \hat{\mathcal{V}}_{(r+1),h+1}(s, a)\}_{s,a}$, and $\{\hat{\mathbf{P}}_{(r+1),h}^i \hat{\mathcal{V}}_{(r+1),h+1}^2(s, a)\}_{s,a}$ to the central server.

Additionally, let us define the weighted average kernel

$$\hat{\mathbf{P}}_{(r),h}(s'|s, a) := \sum_{i=1}^N \frac{n_{(r),h}^i(s, a)}{N_{(r),h}(s, a)} \hat{\mathbf{P}}_{(r),h}^i(s'|s, a) = \frac{N_{(r),h}(s, a, s')}{N_{(r),h}(s, a)}, \quad (25)$$

where $N_{(r),h}(s, a) = \sum_{i=1}^M n_{(r),h}^i(s, a)$ was defined in (6), and $\hat{\mathbf{P}}_{(r),h}^i$ was defined in (9) as

$$\hat{\mathbf{P}}_{(r),h}^i(s'|s, a) = \begin{cases} \frac{n_{(r),h}^i(s, a, s')}{n_{(r),h}^i(s, a)} & \text{if } n_{(r),h}^i(s, a) > 0 \\ \frac{1}{|\mathcal{S}|} & \text{else} \end{cases}. \quad (26)$$

Notably, the kernel $\hat{\mathbf{P}}_{(r),h}$ is never revealed to any agent or to a central server, but it is very useful in the analysis. Also, for any time t we define r_t as an index of the previous epoch. For convenience and ease of reading, we introduce the transition kernels and counters in the *regular timescale*

$$\hat{\mathbf{P}}_{t,h}^i := \hat{\mathbf{P}}_{(r_t),h}^i, \quad \hat{\mathbf{P}}_{t,h} := \hat{\mathbf{P}}_{(r_t),h}, \quad n_{t,h}^i = n_{(r_t),h}^i, \quad \text{and } N_{t,h}^i = N_{(r_t),h}^i, \quad (27)$$

where r_t is defined in (22).

Let $\beta^{\text{KL}}, \beta^c, \beta^{\text{Var}}: (0, 1) \times \mathbb{N} \rightarrow \mathbb{R}_+$ and $\beta^*, \beta, \beta^{\text{max}}: (0, 1) \rightarrow \mathbb{R}_+$ be some functions defined later on in Lemma C.1, and R_{max} be the maximal number of communications defined (58). We define the following favorable events

$$\begin{aligned}
 \mathcal{E}^{\text{KL}}(\delta) &:= \left\{ \forall r \in \mathbb{N}, \forall h \in [H], \forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad \text{KL} \left(\widehat{\mathbf{P}}_{(r),h}^c(s, a) \parallel \mathbf{P}_h^c(s, a) \right) \leq \frac{|\mathcal{S}| \beta^{\text{KL}}(\delta, N_{(r),h}(s, a))}{N_{(r),h}(s, a)} \right\}, \\
 \mathcal{E}^c(\delta) &:= \left\{ \forall r \in [R_{\text{max}}], \forall h \in [H], \forall (s, a) \in \mathcal{S} \times \mathcal{A} : \left\| \widehat{\mathbf{P}}_{(r),h}^c(s, a) - \widehat{\mathbf{P}}_{(r),h}^c(s, a) \right\|_1 \leq \frac{9}{8} \varepsilon_p + \frac{11 \beta^c(\delta, N_{(r),h}(s, a))}{N_{(r),h}(s, a)} \right\}, \\
 \mathcal{E}^*(\delta) &:= \left\{ \forall r \in [R_{\text{max}}], \forall h \in [H], \forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad \left| \widehat{\mathbf{P}}_{(r),h}^c - \mathbf{P}_h^c \mathcal{V}_{h+1}^{c,*}(s, a) \right| \right. \\
 &\quad \left. \leq 1_{\llbracket 2; +\infty \rrbracket}(N_{(r),h}(s, a)) \left(\sqrt{\frac{2 \text{Var}_{\widehat{\mathbf{P}}_{(r),h}^c}(\mathcal{V}_{h+1}^{c,*})(s, a) \beta^*(\delta)}{N_{(r),h}(s, a) - 1}} + \frac{7 \beta^*(\delta)}{N_{(r),h}(s, a) - 1} \right) + H 1_{\llbracket 0; 2 \rrbracket}(N_{(r),h}(s, a)) \right\}, \\
 \mathcal{E}^{\text{Var}}(\delta) &:= \left\{ \forall t \in [T] : \quad \sum_{(t' \geq 1, h \geq 1, i \geq 1)}^{(t, H, M)} \text{Var}_{\mathbf{P}_h^i}(\mathcal{V}_{h+1}^{i, \pi_{t'}})(s_{t',h}^i, a_{t',h}^i) \leq \sqrt{2H^5 M t \beta^{\text{Var}}(\delta, t)} + 3H^3 \beta^{\text{Var}}(\delta, t) + H^2 M t \right\}, \\
 \mathcal{E}^{\text{count}}(\delta) &:= \left\{ \forall t \in [T], \forall h \in [H], \forall (s, a) \in \mathcal{S} \times \mathcal{A}, \forall i \in [M] : \quad |\widetilde{N}_{t,h}^M(s, a) - \hat{N}_{t,h}^i(s, a)| \right. \\
 &\quad \left. \leq \frac{2}{7} \hat{N}_{t,h}^i(s, a) + 2\varepsilon_p T H M + 26M \beta^c(\delta, T) \right\}, \\
 \mathcal{E}(\delta) &:= \left\{ \forall h \in [H] : \quad \sum_{(t \geq 1, h' \geq h, i \geq 1)}^{(T, H, M)} \gamma_{h'-1} \left(\mathbf{P}_{h'}^i \left[\hat{\mathcal{V}}_{t,h'+1} - \mathcal{V}_{h'+1}^{c, \pi_t} \right] (s_{t,h'}^i, a_{t,h'}^i) - \left[\hat{\mathcal{V}}_{t,h'+1} - \mathcal{V}_{h'+1}^{c, \pi_t} \right] (s_{t,h'+1}^i) \right) \right. \\
 &\quad \left. \leq \sqrt{8e^2 H^2 \cdot T H M \cdot \beta(\delta)}, \quad \gamma_h := \left(1 + \frac{1}{H} \right)^{H-h}, \text{ and} \right. \\
 &\quad \left. \sum_{(t \geq 1, h' \geq h, i \geq 1)}^{(T, H, M)} \left(\mathbf{P}_{h'}^i \left[\hat{\mathcal{V}}_{t,h'+1} - \mathcal{V}_{h'+1}^{c, \pi_t} \right] (s_{t,h'}^i, a_{t,h'}^i) - \left[\hat{\mathcal{V}}_{t,h'+1} - \mathcal{V}_{h'+1}^{c, \pi_t} \right] (s_{t,h'+1}^i) \right) \leq \sqrt{8H^2 \cdot T H M \cdot \beta(\delta)} \right\}.
 \end{aligned}$$

We also introduce the intersection of these events, $\mathcal{G}(\delta) := \mathcal{E}^{\text{KL}}(\delta) \cap \mathcal{E}^c(\delta) \cap \mathcal{E}^*(\delta) \cap \mathcal{E}^{\text{Var}}(\delta) \cap \mathcal{E}^{\text{count}}(\delta) \cap \mathcal{E}(\delta)$. We prove that for the right choice of the functions $\beta^{\text{KL}}, \beta^c, \beta^*, \beta^{\text{Var}}$, and β the above events hold with high probability.

Lemma C.1. *For any $\delta \in (0, 1)$ and for the following choices of functions β ,*

$$\begin{aligned}
 \beta^{\text{KL}}(\delta, n) &:= \log(6|\mathcal{S}||\mathcal{A}|H/\delta) + \log(e(1+n)), & \beta^c(\delta, n) &:= \log(6|\mathcal{S}||\mathcal{A}|H/\delta) + \log(6e(2n+1)), \\
 \beta^*(\delta) &:= \log(12|\mathcal{S}||\mathcal{A}|H/\delta), & \beta^{\text{Var}}(\delta, t) &:= \log(24e(2Mt+1)/\delta), \\
 \beta(\delta) &:= \log(48H/\delta).
 \end{aligned}$$

it holds that

$$\begin{aligned}
 \mathbb{P}[\mathcal{E}^{\text{KL}}(\delta)] &\geq 1 - \delta/6, & \mathbb{P}[\mathcal{E}^c(\delta)] &\geq 1 - \delta/6, & \mathbb{P}[\mathcal{E}^*(\delta)] &\geq 1 - \delta/6, & \mathbb{P}[\mathcal{E}^{\text{Var}}(\delta)] &\geq 1 - \delta/6, \\
 \mathbb{P}[\mathcal{E}^{\text{count}}(\delta)] &\geq 1 - \delta/6, & \mathbb{P}[\mathcal{E}(\delta)] &\geq 1 - \delta/6.
 \end{aligned}$$

In particular, $\mathbb{P}[\mathcal{G}(\delta)] \geq 1 - \delta$.

Proof. First, let us define an appropriate filtration for martingale and optional skipping-based arguments. A natural federated online filtration is defined as

$$\mathcal{F}_{t,h}^i = \sigma \left(\{s_{t',h'}^{i'}, a_{t',h'}^{i'}\}_{(t', h', i') \preceq (t, h, i)} \right), \quad (28)$$

where the order over triplets (t', h', i') is lexicographic. With respect to this filtration, for any fixed state-action-step triplet $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ we define the partial global counters that form a sequence of excursion times on an extended time $(t, i) \in \mathbb{N} \times [M]$ and a first extended timestamp to reach a particular partial counter value $j \in [T \cdot M]$

$$\tilde{N}_{t,h}^i(s, a) = \sum_{(t', i') \preceq (t, i)} \mathbf{1}_{(s, a)}(s_{t,h}^i, a_{t,h}^i), \quad (t_{s,a,h,j}, i_{s,a,h,j}) := \min\{(t, i) \in \mathbb{N} \times [M] \mid \tilde{N}_{t,h}^i(s, a) = j\}. \quad (29)$$

For a given time t , we also define $\psi_t := T_{r_t}$ representing the number of episodes visited before r_t . In particular, we have $N_{(r),h}(s, a) = \tilde{N}_{T_{r,h}}^M(s, a)$.

Event $\mathcal{E}^{\text{KL}}(\delta)$ To analyze it, we need first to represent the virtual estimate of the common transition kernel as follows

$$\begin{aligned} \hat{\mathbf{P}}_{(r),h}^c(s'|s, a) &= \frac{1}{N_{(r),h}(s, a)} \sum_{t=1}^{T_r} \sum_{i=1}^M \mathbf{1}_{(s,a)}(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{s'}(s_{t,h+1}^{c,i}) \\ &= \frac{1}{N_{(r),h}(s, a)} \sum_{j=1}^{N_{(r),h}(s,a)} \mathbf{1}_{s'}(s_{t_{s,a,h,j},h+1}^{c,i_{s,a,h,j}}). \end{aligned} \quad (30)$$

By the optional skipping argument (see, e.g., [Doob, 1953](#), Chapter III, p. 145), the sampled states $\{\tilde{s}_{s,a,h,j}^c\}_{j \in [TM]} := \{s_{t_{s,a,h,j},h+1}^{c,i_{s,a,h,j}}\}_{j \in [TM]}$, conditioned on the value of $N_{(r),h}(s, a)$, form an i.i.d. sequence of categorical random variables from the distribution $\mathbf{P}_h^c(s, a)$. Thus, we have for any fixed $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ by Lemma [F.4](#)

$$\mathbb{P} \left[\exists r \geq 1 : \text{KL} \left(\hat{\mathbf{P}}_{(r),h}^c(s, a) \parallel \mathbf{P}_h^c(s, a) \right) \geq \frac{\log(6|\mathcal{S}||\mathcal{A}|H/\delta) + |\mathcal{S}| \log(e(1+n))}{N_{(r),h}(s, a)} \right] \leq \frac{\delta}{6|\mathcal{S}||\mathcal{A}|H}.$$

By a union bound argument and noticing that $\log(6|\mathcal{S}||\mathcal{A}|H/\delta) + |\mathcal{S}| \log(e(1+n)) \leq |\mathcal{S}| \beta^{KL}(\delta, n)$, we conclude the first statement.

Event $\mathcal{E}^c(\delta)$ By a union bound argument, it is enough to show that each of the following events

$$\bar{\mathcal{E}}^c(\delta, s, a, h) = \left\{ \exists r \in [R_{\max}] : \left\| \hat{\mathbf{P}}_{(r),h}(\cdot|s, a) - \hat{\mathbf{P}}_{(r),h}^c(\cdot|s, a) \right\|_1 \geq \frac{9}{8} \varepsilon_p + \frac{11\beta^c(\delta, N_{(r),h}(s, a))}{N_{(r),h}(s, a)} \right\}$$

holds with probability less or equal than $\delta' := \delta/(6|\mathcal{S}||\mathcal{A}|H)$ for any $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$. To do it, let us analyze the difference between kernels. By the definitions [\(24\)](#)-([25](#)):

$$[\hat{\mathbf{P}}_{(r),h} - \hat{\mathbf{P}}_{(r),h}^c](s'|s, a) = \frac{1}{N_{(r),h}(s, a)} \sum_{i=1}^M \sum_{t=1}^{T_r} \left(\mathbf{1}_{(s,a,s')}(s_{t,h}^i, a_{t,h}^i, s_{t,h}^{i,i}) - \mathbf{1}_{(s,a,s')}(s_{t,h}^i, a_{t,h}^i, s_{t,h}^{c,i,i}) \right).$$

Next, we notice that, using representation [\(23\)](#), we can rewrite the first indicator as follows

$$\mathbf{1}_{(s,a,s')}(s_{t,h}^i, a_{t,h}^i, s_{t,h}^{i,i}) = (1 - \xi_{t,h}^i) \cdot \mathbf{1}_{(s,a,s')}(s_{t,h}^i, a_{t,h}^i, s_{t,h}^{c,i,i}) + \xi_{t,h}^i \mathbf{1}_{(s,a,s')}(s_{t,h}^i, a_{t,h}^i, s_{t,h}^{\text{ind},i,i}),$$

thus, we have the following expression for the difference between kernels

$$[\hat{\mathbf{P}}_{(r),h} - \hat{\mathbf{P}}_{(r),h}^c](s'|s, a) = \frac{1}{N_{(r),h}(s, a)} \sum_{i=1}^M \sum_{t=1}^{T_r} \xi_{t,h}^i \mathbf{1}_{(s,a)}(s_{t,h}^i, a_{t,h}^i) \cdot \left(\mathbf{1}_{s'}(s_{t,h}^{\text{ind},i,i}) - \mathbf{1}_{s'}(s_{t,h}^{c,i,i}) \right),$$

and thus, using $|\mathbf{1}_{s'}(s_{t,h}^{\text{ind},i,i}) - \mathbf{1}_{s'}(s_{t,h}^{c,i,i})| \leq 1$ and Definition [\(29\)](#), we obtain

$$\left\| [\hat{\mathbf{P}}_{(r),h} - \hat{\mathbf{P}}_{(r),h}^c](s, a) \right\|_1 \leq \frac{1}{N_{(r),h}(s, a)} \sum_{t=1}^{T_r} \sum_{i=1}^M \xi_{t,h}^i \mathbf{1}_{(s,a)}(s_{t,h}^i, a_{t,h}^i) = \frac{1}{N_{(r),h}(s, a)} \sum_{j=1}^{N_{(r),h}(s,a)} \xi_{t_{s,a,h,j}}^{i_{s,a,h,j}}. \quad (31)$$

Again, by the optional skipping argument, conditioned on the event $N_{(r),h}(s, a) = N$ the sequence $\{\tilde{\xi}_{s,a,h,j}\}_{j \in [TM]} := \{\xi_{t_{s,a,h,j}}^{i_{s,a,h,j}}\}_{j \in [TM]}$ is i.i.d., thus Corollary [F.1](#) implies

$$\mathbb{P}[\bar{\mathcal{E}}^c(\delta, s, a, h)] \leq \mathbb{P} \left[\exists N \geq 1 : \sum_{j=1}^N \tilde{\xi}_{s,a,h,j} > \frac{9}{8} N \varepsilon_p + 11 \log \left(\frac{24|\mathcal{S}||\mathcal{A}|H e(2N+1)}{\delta} \right) \right] \leq \frac{\delta}{6|\mathcal{S}||\mathcal{A}|H}.$$

Event $\mathcal{E}^*(\delta)$ To analyze this event, we use the representation of the kernel (30) and optional skipping argument conditioned on $N_{(r),h}(s, a) = N$ where $N \geq 2$

$$[\widehat{\mathbf{P}}_{(r),h}^c - \mathbf{P}_h^c] \mathcal{V}_{h+1}^{c,*}(s, a) = \frac{1}{N} \sum_{j=1}^N \mathcal{V}_{h+1}^{c,*}(s_{s,a,h,j}^c) - \mathbf{P}_h^c \mathcal{V}_{h+1}^{c,*}(s, a).$$

Thus, we have a sum of centered i.i.d. random variables, and thus we can apply Lemma F.7

$$\mathbb{P} \left[\left| [\widehat{\mathbf{P}}_{(r),h}^c - \mathbf{P}_h^c] \mathcal{V}_{h+1}^{c,*}(s, a) \right| \geq \sqrt{\frac{2 \text{Var}_{\widehat{\mathbf{P}}_{(r),h}^c}(\mathcal{V}_{h+1}^{c,*})(s, a) \beta^*(\delta)}{N-1}} + \frac{7\beta^*(\delta)}{N-1} \middle| N_{(r),h}(s, a) = N \right] \leq \frac{\delta}{6|\mathcal{S}||\mathcal{A}|HT},$$

where $\beta^*(\delta) = \log(12|\mathcal{S}||\mathcal{A}|H/\delta)$. We then conclude by a union bound over $(s, a, h, N) \in \mathcal{S} \times \mathcal{A} \times [H] \times \{2, \dots, MT\}$. If $N_{(r),h}(s, a) \leq 1$, we have the trivial bound $|\widehat{\mathbf{P}}_{(r),h}^c - \mathbf{P}_h^c| \mathcal{V}_{h+1}^{c,*}(s, a) \leq H$.

Event $\mathcal{E}^{\text{Var}}(\delta)$ For any $t' \in [T]$, define

$$X_{t'}^i = \sum_{h=1}^H \text{Var}_{\mathbf{P}_h^i}(\mathcal{V}_{h+1}^{i,\pi_{t'}})(s_{t',h}^i, a_{t',h}^i) - \sigma \mathcal{V}_1^{i,\pi_{t'}}(s_{t',1}^i),$$

where $\sigma \mathcal{V}_1^{i,\pi_{t'}}$ is defined in (59). This sequence forms a martingale-difference sequence with respect to the following filtration

$$\mathcal{F}_t^i = \sigma \left(\{s_{t',h}^{i'}, a_{t',h}^{i'}\}_{(t',i') \preceq (t,i), h \in [H]} \right),$$

where the order over the pairs (t', i') is lexicographic. Applying Theorem F.1 yields

$$\mathbb{P} \left[\exists t \geq 1, \sum_{(t' \geq 1, i \geq 1)}^{(t,M)} X_{t'}^i \leq \sqrt{2 \sum_{(t' \geq 1, i \geq 1)}^{(t,M)} \mathbb{E}_{\pi}[(X_{t'}^i)^2 | \mathcal{F}_{t_{\text{prev}}^{i_{\text{prev}}}}^{i_{\text{prev}}}] \log(24e(2Mt+1)/\delta)} + 3H^3 \log(24e(2Mt+1)/\delta) \right] \leq \frac{\delta}{6},$$

as we have $|X_{t'}^i| \leq H^3$ and where $(t'_{\text{prev}}, i_{\text{prev}})$ is a previous element in a lexicographic order with respect to (t', i) . Now, we bound the conditional second-order moment of $X_{t'}^i$ as follows

$$\mathbb{E}_{\pi_{t'}}[(X_{t'}^i)^2 | \mathcal{F}_{t_{\text{prev}}^{i_{\text{prev}}}}^{i_{\text{prev}}}] \leq \mathbb{E}_{\pi_{t'}} \left[\left(\sum_{h=1}^H \text{Var}_{\mathbf{P}_h^i}(\mathcal{V}_{h+1}^{i,\pi_{t'}})(s_{t',h}^i, a_{t',h}^i) \right)^2 \middle| \mathcal{F}_t^{i-1} \right] \leq H^3 \mathbb{E}_{\pi_{t'}} \left[\sum_{h=1}^H \text{Var}_{\mathbf{P}_h^i}(\mathcal{V}_{h+1}^{i,\pi_{t'}})(s_{t',h}^i, a_{t',h}^i) \right].$$

By Lemma F.3, we have

$$\mathbb{E}_{\pi_{t'}} \left[\sum_{h=1}^H \text{Var}_{\mathbf{P}_h^i}(\mathcal{V}_{h+1}^{i,\pi_{t'}})(s_{t',h}^i, a_{t',h}^i) \right] = \mathbb{E}_{\pi_{t'}} \left[\left(\sum_{h=1}^H r_h^i(s_h^i, a_h^i) - \mathcal{V}_1^{i,\pi_{t'}}(s_1^i) \right)^2 \right] \leq \mathbb{E}_{\pi_{t'}} \left[\left(\sum_{h=1}^H r_h^i(s_h^i, a_h^i) \right)^2 \right] \leq H^2.$$

By combining the previous inequalities, we obtain

$$\sum_{(t' \geq 1, i \geq 1)}^{(t,M)} X_{t'}^i \leq \sqrt{2H^5 M t \log(24e(2Mt+1)/\delta)} + 3H^3 \log(24e(2Mt+1)/\delta)$$

Now using Lemma F.3 again we get

$$\begin{aligned} \sum_{(t' \geq 1, i \geq 1)}^{(t,M)} \sum_{h=1}^H \text{Var}_{\mathbf{P}_h^i}(\mathcal{V}_{h+1}^{i,\pi_{t'}})(s_{t',h}^i, a_{t',h}^i) &= \sum_{(t' \geq 1, i \geq 1)}^{(t,M)} X_{t'}^i + \sigma \mathcal{V}_1^{i,\pi_{t'}}(s_{t',1}^i) \\ &\leq \sqrt{2H^5 M t \log(24e(2Mt+1)/\delta)} + 3H^3 \log(24e(2Mt+1)/\delta) + H^2 M t. \end{aligned}$$

Event $\mathcal{E}^{\text{count}}(\delta)$ For any fixed $(s, a, h, i, t_1) \in \mathcal{S} \times \mathcal{A} \times [H] \times [M] \times [T]$, we have by Corollary F.1

$$\begin{aligned} \mathbb{P}[\exists t_2 \in \mathbb{N} : \left| \sum_{t'=t_1}^{t_2} \mathbf{1}_{(s,a)}(s_{t',h}^i, a_{t',h}^i) - d_h^{i,\pi_{t'}}(s, a) \right| \\ \geq \frac{1}{8} \sum_{t'=t_1}^{t_2} d_h^{i,\pi_{t'}}(s, a) + 11\beta^c(\delta, t_2 - t_1 + 1)] \leq \frac{\delta}{6|\mathcal{S}||\mathcal{A}|MTH}, \end{aligned}$$

holds with probability less or equal than $\delta' := \delta/(6|\mathcal{S}||\mathcal{A}|MTH)$. Thus, by a union bound argument, the following event

$$\begin{aligned} \bar{\mathcal{E}}^{\text{dev}}(\delta) := \left\{ \forall (t_1, t_2) \in [T]^2, \forall h \in [H], \forall (s, a) \in \mathcal{S} \times \mathcal{A}, \forall i \in [M] : \left| \sum_{t'=t_1}^{t_2} \mathbf{1}_{(s,a)}(s_{t',h}^i, a_{t',h}^i) - d_h^{i,\pi_{t'}}(s, a) \right| \right. \\ \left. \geq \frac{1}{8} \sum_{t'=t_1}^{t_2} d_h^{i,\pi_{t'}}(s, a) + 11\beta^c(\delta, t_2 - t_1 + 1) \right\}, \end{aligned}$$

holds with probability less or equal to $\delta/6$. Now, to conclude the proof, it is enough to show $\mathcal{E}^{\text{count}}(\delta) \subset \bar{\mathcal{E}}^{\text{dev}}(\delta)$. Let's recall the definition of the estimated counter by agent i

$$\hat{N}_{t,h}^i(s, a) = N_{(r_t),h}(s, a) + M \sum_{t'=\psi_t}^t \mathbf{1}_{(s,a)}(s_{t',h}^i, a_{t',h}^i).$$

Using (29), the definition of $\hat{N}_{t,h}^i(s, a)$, and the triangular inequality, we have for any fixed $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned} |\tilde{N}_{t,h}^M(s, a) - \hat{N}_{t,h}^i(s, a)| &= \left| \sum_{j \neq i} \sum_{t'=\psi_t}^t \mathbf{1}_{(s,a)}(s_{t',h}^j, a_{t',h}^j) - \mathbf{1}_{(s,a)}(s_{t',h}^i, a_{t',h}^i) \right| \\ &\leq \underbrace{\left| \sum_{j \neq i} \sum_{t'=\psi_t}^t \mathbf{1}_{(s,a)}(s_{t',h}^j, a_{t',h}^j) - d_h^{j,\pi_{t'}}(s, a) \right|}_{(1)} + \underbrace{\left| \sum_{j \neq i} \sum_{t'=\psi_t}^t d_h^{j,\pi_{t'}}(s, a) - d_h^{i,\pi_{t'}}(s, a) \right|}_{(2)} \\ &\quad + \underbrace{(M-1) \left| \sum_{t'=\psi_t}^t d_h^{i,\pi_{t'}}(s, a) - \mathbf{1}_{(s,a)}(s_{t',h}^i, a_{t',h}^i) \right|}_{(3)}. \end{aligned}$$

Term (2): Heterogeneity error Using Lemma F.2 combined with the triangular inequality, it holds that

$$(2) = \left| \sum_{j \neq i} \sum_{t'=\psi_t}^t d_h^{j,\pi_{t'}}(s, a) - d_h^{i,\pi_{t'}}(s, a) \right| \leq \varepsilon_p HMT.$$

Terms (1) and (3): concentration error On the event $\bar{\mathcal{E}}^{\text{dev}}(\delta)$, we can bound (1) as follows

$$\begin{aligned} (1) &\leq \sum_{j \neq i} \sum_{t'=\psi_t}^t \left| \mathbf{1}_{(s,a)}(s_{t',h}^j, a_{t',h}^j) - d_h^{j,\pi_{t'}}(s, a) \right| \\ &\leq \sum_{j \neq i} \frac{1}{8} \sum_{t'=\psi_t}^t d_h^{j,\pi_{t'}}(s, a) + 11M\beta^c(\delta, T) \\ &\leq \underbrace{\frac{1}{8} \sum_{j \neq i} \left| \sum_{t'=\psi_t}^t d_h^{j,\pi_{t'}}(s, a) - d_h^{i,\pi_{t'}}(s, a) \right|}_{(2)} + 11M\beta^c(\delta, T) + \frac{M}{8} \sum_{t'=\psi_t}^t d_h^{i,\pi_{t'}}(s, a). \end{aligned}$$

Now using the latter bound on (2) combined with the inequality $\frac{7}{8} \sum_{t'=\psi_t}^t d_h^{i,\pi_{t'}}(s, a) - 11\beta^c(\delta, T) \leq \sum_{t'=\psi_t}^t \mathbf{1}_{(s,a)}(s_{t',h}^i, a_{t',h}^i)$ that follows from $\mathcal{E}^{\text{dev}}(\delta)$, we get

$$(1) \leq \frac{1}{8} \varepsilon_p HMT + \frac{88M}{7} \beta^c(\delta, T) + \frac{1}{7} \hat{N}_{t,h}^i(s, a).$$

We proceed similarly to bound (3)

$$\begin{aligned} (3) &= (M-1) \left| \sum_{t'=\psi_t}^t d_h^{i,\pi_{t'}}(s, a) - \mathbf{1}_{(s,a)}(s_{t',h}^i, a_{t',h}^i) \right| \\ &\leq \frac{M}{8} \sum_{t'=\psi_t}^t d_h^{j,\pi_{t'}}(s, a) + 11M\beta^c(\delta, T) \leq \frac{1}{7} \hat{N}_{t,h}^i(s, a) + \frac{88M}{7} \beta^c(\delta, T). \end{aligned}$$

Finally combining the bounds on (1), (2) and (3) yields the desired result.

Event $\mathcal{E}(\delta)$ Notice that the two following sequences

$$\begin{aligned} X_{t,h}^i &:= \left(1 + \frac{1}{H}\right)^{H-h'-1} \left(\mathbf{P}_h^i \left[\hat{\mathcal{V}}_{t,h+1} - \mathcal{V}_{h+1}^{c,\pi_t} \right] (s_{t,h}^i, a_{t,h}^i) - \left[\hat{\mathcal{V}}_{t,h+1} - \mathcal{V}_{h+1}^{c,\pi_t} \right] (s_{t,h+1}^i) \right), \\ Y_{t,h}^i &:= \mathbf{P}_h^i \left[\hat{\mathcal{V}}_{t,h+1} - \mathcal{V}_{h+1}^{c,\pi_t} \right] (s_{t,h}^i, a_{t,h}^i) - \left[\hat{\mathcal{V}}_{t,h+1} - \mathcal{V}_{h+1}^{c,\pi_t} \right] (s_{t,h+1}^i), \end{aligned}$$

forms a martingale-difference sequence with respect to filtration $\mathcal{F}_{t,h}^i$ defined in (28). Thus, applying Azuma-Hoeffding inequality with a union bound over h and over the two events allows us to conclude the statement. \square

Lemma C.2. *Conditioned on $\mathcal{E}^{\text{KL}}(\delta)$, for any function $f : \mathcal{S} \mapsto [0, H]$, $h \in [H]$, $(s, a) \in \mathcal{S} \times \mathcal{A}$, and any $r \in [R_{\max}]$, we have*

$$\begin{aligned} (\hat{\mathbf{P}}_{(r),h}^c - \mathbf{P}_h^c) f(s, a) &\leq \frac{1}{H} \mathbf{P}_h^c f(s, a) + \frac{2H^2 |\mathcal{S}| \beta^{\text{KL}}(\delta, N_{(r),h}(s, a))}{N_{(r),h}(s, a)}, \\ \|\hat{\mathbf{P}}_{(r),h}^c - \mathbf{P}_h^c\|_1 &\leq \sqrt{\frac{2|\mathcal{S}| \beta^{\text{KL}}(\delta, N_{(r),h}(s, a))}{N_{(r),h}(s, a)}}. \end{aligned}$$

Proof. Using Lemma F.5 with $\mathbf{P} = \hat{\mathbf{P}}_{(r),h}^c(\cdot | s, a)$ and $\mathbf{Q} = \mathbf{P}_h^c(\cdot | s, a)$ it holds that

$$(\hat{\mathbf{P}}_{(r),h}^c - \mathbf{P}_h^c) f(\cdot | s, a) \leq \sqrt{2 \text{Var}_{\mathbf{P}_h^c(\cdot | s, a)}(f) \text{KL}(\hat{\mathbf{P}}_{(r),h}^c(\cdot | s, a) \| \mathbf{P}_h^c(\cdot | s, a))} + \frac{2}{3} H \text{KL}(\hat{\mathbf{P}}_{(r),h}^c(\cdot | s, a) \| \mathbf{P}_h^c(\cdot | s, a)). \quad (32)$$

Now, since f 's values are in $[0, H]$, we can write

$$\text{Var}_{\mathbf{P}_h^c(\cdot | s, a)}(f) \leq \mathbf{P}_h^c(f^2)(s, a) \leq H \mathbf{P}_h^c(f)(s, a). \quad (33)$$

Combining the latter inequality with the fact that for all $a, b \geq 0$, $\sqrt{2ab} \leq a + b$, we obtain

$$\begin{aligned} \sqrt{2 \text{Var}_{\mathbf{P}_h^c(\cdot | s, a)}(f) \text{KL}(\hat{\mathbf{P}}_{(r),h}^c(\cdot | s, a) \| \mathbf{P}_h^c(\cdot | s, a))} &= \sqrt{\frac{2}{H} \mathbf{P}_h^c(f)(s, a) \cdot H^2 \text{KL}(\hat{\mathbf{P}}_{(r),h}^c(\cdot | s, a) \| \mathbf{P}_h^c(\cdot | s, a))} \\ &\leq \frac{1}{H} \mathbf{P}_h^c(f)(s, a) + H^2 \text{KL}(\hat{\mathbf{P}}_{(r),h}^c(\cdot | s, a) \| \mathbf{P}_h^c(\cdot | s, a)). \end{aligned} \quad (34)$$

Furthermore, since $\mathcal{E}^{\text{KL}}(\delta)$ holds, we have the inequality $\text{KL}(\hat{\mathbf{P}}_{(r),h}^c(s, a) \| \mathbf{P}_h^c(s, a)) \leq \frac{|S| \beta^{\text{KL}}(\delta, N_{(r),h}(s, a))}{N_{(r),h}(s, a)}$. Plugging this bound in (34), we can upper bound (32) as

$$(\hat{\mathbf{P}}_{(r),h}^c - \mathbf{P}_h^c) f(\cdot | s, a) \leq \frac{1}{H} \mathbf{P}_h^c(f)(s, a) + H^2 \frac{|S| \beta^{\text{KL}}(\delta, N_{(r),h}(s, a))}{N_{(r),h}(s, a)} + \frac{2H}{3} \frac{|S| \beta^{\text{KL}}(\delta, N_{(r),h}(s, a))}{N_{(r),h}(s, a)},$$

which gives the result. The second inequality follows from the combination of Pinsker inequality and the definition of $\mathcal{E}^{\text{KL}}(\delta)$. \square

D REGRET ANALYSIS

We define the common MDP \mathcal{M}^c as

$$\mathcal{M}^c := (\mathcal{S}, \mathcal{A}, H, \{r_h^c := \frac{1}{M} \sum_{i=1}^M r_h^i\}_h, \{P_h^c\}_h). \quad (35)$$

We denote by $\mathcal{V}_h^{c,*}$ and $\mathcal{Q}_h^{c,*}$ the value function and the Q-function at step h in the common environment \mathcal{M}^c . In particular, these functions satisfy Bellman's equations and Bellman's optimality equations (Sutton and Barto, 2018)

$$\mathcal{Q}_h^{c,\pi}(s, a) = r_h^c(s, a) + P_h^c \mathcal{V}_{h+1}^{c,\pi}(s, a), \quad \mathcal{V}_h^{c,\pi}(s) = \mathcal{Q}_h^{c,\pi}(s, \pi_h(s)) \quad (36)$$

$$\mathcal{Q}_h^{c,*}(s, a) = r_h^c(s, a) + P_h^c \mathcal{V}_{h+1}^{c,*}(s, a), \quad \mathcal{V}_h^{c,*}(s) = \max_{a \in \mathcal{A}} \mathcal{Q}_h^{c,*}(s, a), \quad (37)$$

D.1 Optimism

Let us define the following event

$$\mathcal{E}^{\text{optimism}} = \left\{ \forall r \in [R_{\max}], \forall (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H] : \hat{\mathcal{V}}_{(r),h}(s) \geq \mathcal{V}_h^{c,*}(s) - (2\varepsilon_r + 3\varepsilon_p H)(H + 1 - h), \right. \\ \left. \hat{\mathcal{Q}}_{(r),h}(s, a) \geq \mathcal{Q}_h^{c,*}(s, a) - (2\varepsilon_r + 3\varepsilon_p H)(H + 1 - h) \right\}.$$

Then, we will show that this event holds on event $\mathcal{G}(\delta)$. To prove the optimism of our estimates, we use the same monotonicity arguments as in Zhang et al. (2021), see also Zhang et al. (2024b). Define

$$g(P, f, \alpha) = Pf + \max(\sqrt{\alpha \text{Var}_P(f)}, \alpha H), \quad (38)$$

where P be is a probability measure on \mathcal{S} , $f \in \mathbb{R}^{|\mathcal{S}|}$ is a non negative vector satisfying $\|f\|_\infty \leq H$, and α is a positive real number.

Lemma D.1 (Lemma 14 by Zhang et al. 2021). *The function g is non-decreasing in each entry of f .*

For completeness, we provide the proof below.

Proof. To justify this claim, consider any $s \in \mathcal{S}$, and let us fix P, α and all but the s -th entries of f . It then suffices to observe that (i) g is a differentiable almost everywhere function, and (ii) except for at most two possible choices of $f(s)$ that obey $\sqrt{\alpha \text{Var}_P(f)} = \alpha H$, one can use the properties of P and f to calculate

$$\begin{aligned} \frac{\partial g(P, f, \alpha)}{\partial f(s)} &= P(s) + \sqrt{\alpha} \mathbf{1} \left\{ \sqrt{\alpha \text{Var}_P(f)} \geq \alpha H \right\} \frac{P(s)(f(s) - \mathbb{E}_{s' \sim P}[f(s')])}{\sqrt{\text{Var}_P(f)}} \\ &= P(s) + \mathbf{1} \left\{ \sqrt{\alpha \text{Var}_P(f)} \geq \alpha H \right\} \frac{\alpha H}{\sqrt{\alpha \text{Var}_P(f)}} \cdot \frac{P(s)(f(s) - \mathbb{E}_{s' \sim P}[f(s')])}{H} \\ &\geq \min \left\{ P(s) + P(s) \frac{(f(s) - \mathbb{E}_{s' \sim P}[f(s')])}{H}, P(s) \right\} \\ &\geq P(s) \min \left\{ \frac{H + f(s) - \mathbb{E}_{s' \sim P}[f(s')]}{H}, 1 \right\} \geq 0, \end{aligned}$$

where in the end we used the fact that $\|f\|_\infty \leq H$. \square

We define the bonus function as

$$b_{(r),h}(s, a) := \begin{cases} \frac{28\beta^*(\delta)H + 11\beta^c(\delta, N)}{N} + \sqrt{\frac{8\beta^*(\delta)}{N} \cdot \text{Var}_{\hat{P}_{(r)}}(\hat{\mathcal{V}}_{(r),h+1}(s, a))}, & N \geq 2 \\ H, & N \leq 1 \end{cases} \quad (39)$$

for $N = N_{(r),h}(s, a)$ and where β^* and β^c are defined in Lemma C.1.

Lemma D.2. *Under conditions of Lemma C.1, it holds $\mathcal{E}^{\text{optimism}} \subseteq \mathcal{G}(\delta)$ for any $\delta \in (0, 1)$.*

Proof. We process the proof by backward induction over h .

Base case For $h = H + 1$ and for all $(s, a, r) \in \mathcal{S} \times \mathcal{A} \times [R_{\max}]$, we have

$$\hat{\mathcal{V}}_{(r),h}(s) = 0 \geq \mathcal{V}_h^{c,\star}(s) - 0 = 0 \quad \text{and} \quad \hat{\mathcal{Q}}_{(r),h}(s, a) = 0 \geq \mathcal{Q}_h^{c,\star}(s, a) - 0 = 0 ,$$

which gives the desired result.

Induction Let $h \in [H]$ such that for all $(s, a, r) \in \mathcal{S} \times \mathcal{A} \times [R_{\max}]$ and $h' \geq h$

$$\hat{\mathcal{V}}_{(r),h'}(s) \geq \mathcal{V}_{h'}^{c,\star}(s) - (2\varepsilon_r + 3\varepsilon_p H)(H + 1 - h) , \quad \text{and} \quad (40)$$

$$\hat{\mathcal{Q}}_{(r),h'}(s, a) \geq \mathcal{Q}_{h'}^{c,\star}(s, a) - (2\varepsilon_r + 3\varepsilon_p H)(H + 1 - h) . \quad (41)$$

First, let us consider a trivial case $\hat{\mathcal{Q}}_{(r),h}(s, a) = H$. The result is trivial since $H \geq \mathcal{Q}_h^{c,\star}(s, a)$.

Next, we assume that $\hat{\mathcal{Q}}_{(r),h}(s, a) < H$. In particular, by the definition of bonuses, it automatically follows that $N_{(r),h}(s, a) \geq 2$. In this case, according to the update rule (11), we have

$$\begin{aligned} \hat{\mathcal{Q}}_{(r),h}(s, a) &\geq \sum_{i=1}^M \frac{n_{(r),h}^i(s, a)}{N_{(r),h}(s, a)} \hat{\mathcal{Q}}_{(r),h}^i(s, a) + b_{(r),h}(s, a) \\ &= \sum_{i=1}^M \frac{n_{(r),h}^i(s, a)}{N_{(r),h}(s, a)} \hat{r}_h^i(s, a) + \hat{\mathbf{P}}_{(r),h} \hat{\mathcal{V}}_{(r),h+1}(s, a) + b_{(r),h}(s, a) \\ &= \frac{1}{M} \sum_{i=1}^M \hat{r}_h^i(s, a) + \mathbf{P}_h^c \mathcal{V}_{h+1}^{c,\star}(s, a) + b_{(r),h}(s, a) + \underbrace{\sum_{i=1}^M \frac{n_{(r),h}^i(s, a)}{N_{(r),h}(s, a)} \hat{r}_h^i(s, a) - \frac{1}{M} \sum_{i=1}^M \hat{r}_h^i(s, a)}_{\text{(I)}} \\ &\quad + \underbrace{\hat{\mathbf{P}}_{(r),h} (\hat{\mathcal{V}}_{(r),h+1}(s, a) - \mathcal{V}_{h+1}^{c,\star}(s, a))}_{\text{(II)}} + \underbrace{(\hat{\mathbf{P}}_{(r),h} - \hat{\mathbf{P}}_{(r),h}^c) \mathcal{V}_{h+1}^{c,\star}(s, a)}_{\text{(III)}} + \underbrace{(\hat{\mathbf{P}}_{(r),h} - \mathbf{P}_h^c) \mathcal{V}_{h+1}^{c,\star}(s, a)}_{\text{(IV)}} . \end{aligned} \quad (42)$$

Terms (I) and (III): heterogeneity errors First, let us handle the terms that come from the presence of heterogeneity between agents. To analyse (I), recall that since for all $(s, a, i, h, r) \in \mathcal{S} \times \mathcal{A} \times [M] \times [H] \times [R_{\max}]$, either 1) $n_{(r),h}^i(s, a) = 0$ and the value of $\hat{r}_h^i(s, a)$ is ignored in the weighted sum, or 2) $n_{(r),h}^i(s, a) > 0$ and $\hat{r}_h^i(s, a) = r_h^i(s, a)$. Thus, $\sum_{i=1}^M \frac{n_{(r),h}^i(s, a)}{N_{(r),h}(s, a)} \hat{r}_h^i(s, a)$ is a convex combination of the true rewards over i , which ensures that

$$\text{(I)} = \sum_{i=1}^M \frac{n_{(r),h}^i(s, a)}{N_{(r),h}(s, a)} \hat{r}_h^i(s, a) - \frac{1}{M} \sum_{i=1}^M \hat{r}_h^i(s, a) \geq -2\varepsilon_r . \quad (43)$$

Conditioned on $\mathcal{E}^c(\delta)$, Hölder's inequality yields the following bound on (III)

$$\text{(III)} = (\hat{\mathbf{P}}_{(r),h} - \hat{\mathbf{P}}_{(r),h}^c) \mathcal{V}_{h+1}^{c,\star}(s, a) \geq -\|\hat{\mathbf{P}}_{(r),h} - \hat{\mathbf{P}}_{(r),h}^c\|_1 \cdot \|\mathcal{V}_{h+1}^{c,\star}\|_\infty \geq -2\varepsilon_p H - \frac{11\beta^c(\delta, N_{(r),h}(s, a))}{N_{(r),h}(s, a)} H . \quad (44)$$

Term (II): correction error To control this term, we aim to apply Lemma D.1. We first define the *shifted estimator* $\bar{\mathcal{V}}_{(r),h+1}$ as

$$\bar{\mathcal{V}}_{(r),h+1}(s) := \hat{\mathcal{V}}_{(r),h+1}(s) + (2\varepsilon_r + 3\varepsilon_p H)(H - h) . \quad (45)$$

By the induction hypothesis (40), we know that $\bar{\mathcal{V}}_{(r),h+1}(s) \geq \mathcal{V}_{h+1}^{c,\star}(s, a)$. We decompose further (II) as

$$\begin{aligned} \text{(II)} &= \hat{\mathbf{P}}_{(r),h} \hat{\mathcal{V}}_{(r),h+1}(s, a) - \hat{\mathbf{P}}_{(r),h} \mathcal{V}_{h+1}^{c,\star}(s, a) \\ &\geq \hat{\mathbf{P}}_{(r),h} \bar{\mathcal{V}}_{(r),h+1}(s, a) + \max \left(\sqrt{\frac{4\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\bar{\mathcal{V}}_{(r),h+1}(s))(s, a)}{N_{(r),h}(s, a)}}, \frac{4\beta^*(\delta)H}{N_{(r),h}(s, a)} \right) - \hat{\mathbf{P}}_{(r),h} \mathcal{V}_{h+1}^{c,\star}(s, a) \end{aligned}$$

$$- \sqrt{\frac{4\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\bar{\mathcal{V}}_{(r),h+1}(s,a))}{N_{(r),h}(s,a)}} - \frac{4\beta^*(\delta)H}{N_{(r),h}(s,a)} - (2\varepsilon_r + 3\varepsilon_p H)(H-h),$$

where we used in the last inequality that for any $a, b \in \mathbb{R}_+$, $\max(a, b) \leq a+b$ and the fact that $\hat{\mathbf{P}}_{(r),h} \hat{\mathcal{V}}_{(r),h+1}(s,a) - \hat{\mathbf{P}}_{(r),h} \bar{\mathcal{V}}_{(r),h+1}(s,a) = -(2\varepsilon_r + 3\varepsilon_p H)(H-h)$. Now by applying Lemma D.1, we get

$$\begin{aligned} \text{(II)} &\geq \max \left(\sqrt{\frac{4\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\mathcal{V}_{(r),h+1}^{c,*})(s,a)}{N_{(r),h}(s,a)}}, \frac{4\beta^*(\delta)H}{N_{(r),h}(s,a)} \right) - \sqrt{\frac{4\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\bar{\mathcal{V}}_{(r),h+1}(s,a))}{N_{(r),h}(s,a)}} - \frac{4\beta^*(\delta)H}{N_{(r),h}(s,a)} \\ &\geq \underbrace{\sqrt{\frac{4\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\mathcal{V}_{(r),h+1}^{c,*})(s,a)}{N_{(r),h}(s,a)}}}_{(1)} - \underbrace{\sqrt{\frac{4\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\bar{\mathcal{V}}_{(r),h+1}(s,a))}{N_{(r),h}(s,a)}}}_{(2)} - \frac{4\beta^*(\delta)H}{N_{(r),h}(s,a)}. \end{aligned} \quad (46)$$

We want now to control the variance terms that appear in (1) and (2). Using inequalities (62) and (61) of Lemma F.6, we have

$$\begin{aligned} \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\mathcal{V}_{h+1}^{c,*})(s,a) &\geq \text{Var}_{\hat{\mathbf{P}}_{(r),h}^c}(\mathcal{V}_{h+1}^{c,*})(s,a) - 3H^2\varepsilon_p, \\ \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\bar{\mathcal{V}}_{(r),h+1})(s,a) &\leq 2 \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\hat{\mathcal{V}}_{(r),h+1})(s,a) + 2\hat{\mathbf{P}}_{(r),h}^c[\bar{\mathcal{V}}_{(r),h+1} - \hat{\mathcal{V}}_{(r),h+1}] \\ &\leq 2 \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\hat{\mathcal{V}}_{(r),h+1})(s,a) + 2(2\varepsilon_r + 3\varepsilon_p H)(H-h), \end{aligned}$$

where in the last inequality we used the induction hypothesis. Besides, as for any $a, b, c \in \mathbb{R}_+$, we have $a \geq b - c \implies \sqrt{a} \geq \sqrt{b} - \sqrt{c}$, and also for any $d, f \in \mathbb{R}_+$ we have $\sqrt{d+f} \leq \sqrt{d} + \sqrt{f}$, we get

$$\begin{aligned} (1) &:= \sqrt{\frac{4\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\mathcal{V}_{(r),h+1}^{c,*})(s,a)}{N_{(r),h}(s,a)}} \geq \sqrt{\frac{4\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{(r),h}^c}(\mathcal{V}_{(r),h+1}^{c,*})(s,a)}{N_{(r),h}(s,a)}} - \sqrt{\frac{12\varepsilon_p H^2 \beta^*(\delta)}{N_{(r),h}(s,a)}}, \text{ and} \\ (2) &:= \sqrt{\frac{4\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\bar{\mathcal{V}}_{(r),h+1})(s,a)}{N_{(r),h}(s,a)}} \leq \sqrt{\frac{8\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\hat{\mathcal{V}}_{(r),h+1})(s,a)}{N_{(r),h}(s,a)}} + \sqrt{\frac{8\beta^*(\delta)(3\varepsilon_p H + 2\varepsilon_r)(H+1-h)}{N_{(r),h}(s,a)}} \end{aligned} \quad (47)$$

Plugging the inequalities (47) and (48) in (46), we obtain

$$\begin{aligned} \text{(II)} &\geq \sqrt{\frac{4\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{(r),h}^c}(\mathcal{V}_{(r),h+1}^{c,*})(s,a)}{N_{(r),h}(s,a)}} - \sqrt{\frac{8\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\hat{\mathcal{V}}_{(r),h+1})(s,a)}{N_{(r),h}(s,a)}} - \sqrt{\frac{12\varepsilon_p H^2 \beta^*(\delta)}{N_{(r),h}(s,a)}} - \frac{4\beta^*(\delta)H}{N_{(r),h}(s,a)} \\ &\quad - \sqrt{\frac{8\beta^*(\delta)(3\varepsilon_p H + 2\varepsilon_r)(H+1-h)}{N_{(r),h}(s,a)}}. \end{aligned} \quad (48)$$

Finally, as for any $a, b \in \mathbb{R}_+$ we have $\sqrt{2ab} \leq a+b$ then

$$\begin{aligned} \text{(II)} &\geq \sqrt{\frac{4\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{(r),h}^c}(\mathcal{V}_{(r),h+1}^{c,*})(s,a)}{N_{(r),h}(s,a)}} - \sqrt{\frac{8\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{(r)}}(\hat{\mathcal{V}}_{(r),h+1})(s,a)}{N_{(r),h}(s,a)}} - \varepsilon_p H - \frac{10\beta^*(\delta)H}{N_{(r),h}(s,a)} \\ &\quad - (3\varepsilon_p H + 2\varepsilon_r)(H-h) - \frac{4\beta^*(\delta)H}{N_{(r),h}(s,a)} \end{aligned} \quad (49)$$

Term (IV): concentration error Conditioned on $\mathcal{E}^*(\delta)$, we have

$$\text{(IV)} = (\hat{\mathbf{P}}_{(r),h}^c - \mathbf{P}_h^c) \mathcal{V}_{h+1}^{c,*}(s,a) \geq - \left| [\hat{\mathbf{P}}_{(r),h}^c - \mathbf{P}_h^c] \mathcal{V}_{h+1}^{c,*}(s,a) \right|$$

$$\geq -\sqrt{\frac{2 \operatorname{Var}_{\widehat{\mathbf{P}}_{(r)}^c}(\mathcal{V}_{h+1}^{c,*})(s, a) \beta^*(\delta)}{N_{(r),h}(s, a) - 1}} - \frac{7\beta^*(\delta)}{N_{(r),h}(s, a) - 1} \geq -\sqrt{\frac{4 \operatorname{Var}_{\widehat{\mathbf{P}}_{(r)}^c}(\mathcal{V}_{h+1}^{c,*})(s, a) \beta^*(\delta)}{N_{(r),h}(s, a)}} - \frac{14\beta^*(\delta)}{N_{(r),h}(s, a)}, \quad (51)$$

as for $n \geq 2$, we have $\frac{2}{n} \geq \frac{1}{n-1}$.

Combine everything together By plugging in the bounds on (I), (II), (III), and (IV) in (42), we get

$$\begin{aligned} \hat{\mathcal{Q}}_{(r),h}(s, a) &\geq \frac{1}{M} \sum_{i=1}^M r_h^i(s, a) + \mathbf{P}_h^c \mathcal{V}_{h+1}^{c,*}(s, a) + b_{(r),h}(s, a) - \frac{11\beta^c(\delta, N_{(r),h}(s, a))}{N_{(r),h}(s, a)} H - \frac{28\beta^*(\delta)H}{N_{(r),h}(s, a)} \\ &\quad - \sqrt{\frac{8\beta^*(\delta) \operatorname{Var}_{\widehat{\mathbf{P}}_{(r)}^c}(\hat{\mathcal{V}}_{(r),h+1})(s, a)}{N_{(r),h}(s, a)}} - (3\varepsilon_p H + 2\varepsilon_r)(H + 1 - h) = \mathcal{Q}_h^{c,*}(s, a) - (2\varepsilon_r + 3\varepsilon_p H)(H + 1 - h), \end{aligned}$$

where the last inequality is a consequence of the definition of the bonus (39) and optimal Bellman equations (37). \square

D.2 Regret decomposition

We will start by writing down a regret decomposition. Let us define the essential technical quantities, such as common regret and partial common upper regret

$$\mathfrak{R}^c(T) := \frac{1}{M} \sum_{t=1}^T \sum_{i=1}^M \mathcal{V}_1^{c,*}(s_{t,1}^i) - \mathcal{V}_1^{c,\pi_t}(s_{t,1}^i), \quad \overline{\mathfrak{R}}_h^c(T) := \frac{1}{M} \sum_{t=1}^T \sum_{i=1}^M \hat{\mathcal{V}}_{t,h}(s_{t,h}^i) - \mathcal{V}_h^{c,\pi_t}(s_{t,h}^i)$$

Lemma D.3. Assume conditions of Lemma C.1. Then, on the event $\mathcal{G}(\delta)$, the following inequality for any partial upper common regret holds

$$\overline{\mathfrak{R}}_h^c(T) \leq U_h^T := A_h^T + B_h^T + C_h^T + 7eTH^2\varepsilon_p + 2eTH\varepsilon_r + \sqrt{8H^2 \cdot TH \cdot \beta(\delta)/M} + \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h'=h}^H 2eH \mathbf{1}_{[0;1]}(\bar{N}_{t,h'}^i),$$

where

$$\begin{aligned} \bar{N}_{t,h}^i &:= N_{(r_t),h}(s_{t,h}^i, a_{t,h}^i), \\ A_h^T &:= \frac{e}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h'=h}^H \sqrt{\frac{4\beta^*(\delta) \operatorname{Var}_{\widehat{\mathbf{P}}_{t,h'}^c}(\mathcal{V}_{h+1}^{c,*})(s_{t,h'}^i, a_{t,h'}^i)}{\bar{N}_{t,h'}^i}} \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h'}^i), \\ B_h^T &:= \frac{e}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h'=h}^H \sqrt{\frac{8\beta^*(\delta) \operatorname{Var}_{\widehat{\mathbf{P}}_{t,h}^c}(\hat{\mathcal{V}}_{t,h'+1})(s_{t,h'}^i, a_{t,h'}^i)}{\bar{N}_{t,h'}^i}} \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h'}^i), \\ C_h^T &:= \frac{e}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h'=h}^H \frac{22\beta^c(\delta, \bar{N}_{t,h}^i) + 46H\beta^*(\delta) + 2H^2|\mathcal{S}|\beta^{\text{KL}}(\delta, \bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i} \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h'}^i). \end{aligned}$$

Proof. Let us define $\delta_{t,h}^i = \hat{\mathcal{V}}_{t,h}(s_{t,h}^i) - \mathcal{V}_1^{c,\pi_t}(s_{t,h}^i)$ and let us study this term separately. Since the policy is deterministic, i.e., $a_{t,h}^i = \pi_{t,h}(s_{t,h}^i)$, and satisfies $\hat{\mathcal{V}}_{t,h}(s_{t,h}^i) = \hat{\mathcal{Q}}_{t,h}(s_{t,h}^i, a_{t,h}^i)$, we have

$$\delta_{t,h}^i = \hat{\mathcal{Q}}_{t,h}(s_{t,h}^i, a_{t,h}^i) - \mathcal{Q}_h^{c,*}(s_{t,h}^i, a_{t,h}^i) + \mathcal{Q}_h^{c,*}(s_{t,h}^i, a_{t,h}^i) - \mathcal{Q}_h^{c,\pi_t}(s_{t,h}^i, a_{t,h}^i).$$

Next, for empirical Q-values, we have the following bound due to the clipping mechanism and A-1

$$\hat{\mathcal{Q}}_{t,h}(s, a) \leq \sum_{i=1}^N \frac{n_{t,h}^i(s, a)}{\bar{N}_{t,h}^i} \hat{r}_h^i(s, a) + \hat{\mathbf{P}}_{t,h} \hat{\mathcal{V}}_{t,h+1}(s, a) + b_{t,h}(s, a)$$

$$\leq r_h^c(s, a) + \hat{P}_{t,h} \hat{V}_{t,h+1}(s, a) + b_{t,h}(s, a) + 2\varepsilon_r,$$

thus, applying Bellman equations (36) and optimal Bellman equations (37), we have after a simple rearranging

$$\begin{aligned} \delta_{t,h}^i &\leq \hat{P}_{t,h} \hat{V}_{t,h+1}(s_{t,h}^i, a_{t,h}^i) - P_h^c \mathcal{V}_{h+1}^{c,*}(s_{t,h}^i, a_{t,h}^i) + b_{t,h}(s_{t,h}^i, a_{t,h}^i) + P_h^c [\mathcal{V}_{h+1}^{c,*} - \mathcal{V}_{h+1}^{c,\pi_t}](s_{t,h}^i, a_{t,h}^i) + 2\varepsilon_r \\ &= [\hat{P}_{t,h} - P_h^c] \hat{V}_{t,h+1}(s_{t,h}^i, a_{t,h}^i) + P_h^c [\hat{V}_{t,h+1} - \mathcal{V}_{h+1}^{c,\pi_t}](s_{t,h}^i, a_{t,h}^i) + b_{t,h}(s_{t,h}^i, a_{t,h}^i) + 2\varepsilon_r. \end{aligned}$$

In the decomposition above, we further rearrange it, using a virtual estimate of $\hat{P}_{t,h}^c$ defined in (24) and re-introducing again the kernel for i -th agent P_h^i

$$\delta_{t,h}^i \leq \underbrace{[\hat{P}_{t,h} - \hat{P}_{t,h}^c] \hat{V}_{t,h+1}(s_{t,h}^i, a_{t,h}^i)}_{\text{(A)}} + \underbrace{[\hat{P}_{t,h}^c - P_h^c] [\hat{V}_{t,h+1} - \mathcal{V}_{h+1}^{c,*}](s_{t,h}^i, a_{t,h}^i)}_{\text{(B)}} + \underbrace{[\hat{P}_{t,h}^c - P_h^c] \mathcal{V}_{h+1}^{c,*}(s_{t,h}^i, a_{t,h}^i)}_{\text{(C)}} \quad (52)$$

$$+ \underbrace{[P_h^c - P_h^i] [\hat{V}_{t,h+1} - \mathcal{V}_{h+1}^{c,\pi_t}](s_{t,h}^i, a_{t,h}^i)}_{\text{(D)}} + \underbrace{P_h^i [\hat{V}_{t,h+1} - \mathcal{V}_{h+1}^{c,\pi_t}](s_{t,h}^i, a_{t,h}^i) - [\hat{V}_{t,h+1} - \mathcal{V}_{h+1}^{c,\pi_t}](s_{t,h}^i, a_{t,h}^i)}_{=:\zeta_{t,h}^i} \quad (53)$$

$$+ \underbrace{[\hat{V}_{t,h+1} - \mathcal{V}_{h+1}^{c,\pi_t}](s_{t,h}^i, a_{t,h}^i)}_{\delta_{t,h+1}^i} + b_{t,h}(s_{t,h}^i, a_{t,h}^i) + 2\varepsilon_r. \quad (54)$$

Next, we analyze each term separately. With a slight abuse of notation, let us define $\bar{N}_{t,h}^i = N_{(r_t),h}(s_{t,h}^i, a_{t,h}^i)$. In the sequel, we analyze only such $(t, i, h) \in [T] \times [M] \times [H]$ such that $\bar{N}_{t,h}^i \geq 2$. In the case where $\bar{N}_{t,h}^i \leq 1$, we have the trivial bound $\delta_{t,h}^i \leq H$.

Terms (A) and (D): heterogeneity errors First, let us handle the terms that come from the presence of heterogeneity between agents. To analyze (A), let us apply the definition of the event $\mathcal{E}^c(\delta) \subseteq \mathcal{G}(\delta)$ combined with Holder's inequality

$$\text{(A)} \leq H \|\hat{P}_{t,h}(s_{t,h}^i, a_{t,h}^i) - \hat{P}_{t,h}^c(s_{t,h}^i, a_{t,h}^i)\|_1 \leq 2H\varepsilon_p + \frac{11H\beta^c(\delta, \bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i}.$$

For (D) we apply Holder's inequality, A-1 and Corollary F.2

$$\text{(D)} \leq 2H \|P_h^c(s_{t,h}^i, a_{t,h}^i) - P_h^i(s_{t,h}^i, a_{t,h}^i)\|_1 \leq 2H\varepsilon_p.$$

Term (B): correction error To analyze this term, we apply Lemma C.2 with $f(s') := [\hat{V}_{t,h+1} - \mathcal{V}_{h+1}^{c,*}](s')$ and get

$$\begin{aligned} \text{(B)} &\leq \frac{1}{H} P_h^c [\hat{V}_{t,h+1} - \mathcal{V}_{h+1}^{c,*}](s_{t,h}^i, a_{t,h}^i) + \frac{2H^2 |\mathcal{S}| \beta^{\text{KL}}(\delta, \bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i} \\ &\stackrel{(1)}{\leq} \frac{1}{H} P_h^c [\hat{V}_{t,h+1} - \mathcal{V}_{h+1}^{c,\pi_t}](s_{t,h}^i, a_{t,h}^i) + \frac{2H^2 |\mathcal{S}| \beta^{\text{KL}}(\delta, \bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i} \\ &\stackrel{(2)}{\leq} \frac{1}{H} \text{(D)} + \frac{1}{H} \delta_{t,h+1}^i + \frac{1}{H} \zeta_{t,h+1}^i + \frac{2H^2 |\mathcal{S}| \beta^{\text{KL}}(\delta, \bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i}. \end{aligned}$$

where (1) follows from the definition of optimal policy, and (2) follows from a simple rearranging of terms, similar to the decomposition of $\delta_{t,h}^i$. Additional term (D) appeared compared to a standard decomposition.

Term (C): concentration error From the definition of the event $\mathcal{E}^*(\delta) \subseteq \mathcal{G}(\delta)$ defined in Lemma C.1, and from the analysis of the case $\bar{N}_{t,h}^i \geq 2$ it follows that

$$\text{(C)} \leq \sqrt{\frac{2 \text{Var}_{\hat{P}_{t,h}^c}(\mathcal{V}_{h+1}^{c,*})(s_{t,h}^i, a_{t,h}^i) \beta^*(\delta)}{\bar{N}_{t,h}^i - 1}} + \frac{7\beta^*(\delta)}{\bar{N}_{t,h}^i - 1} \leq \sqrt{\frac{4 \text{Var}_{\hat{P}_{t,h}^c}(\mathcal{V}_{h+1}^{c,*})(s_{t,h}^i, a_{t,h}^i) \beta^*(\delta)}{\bar{N}_{t,h}^i}} + \frac{14\beta^*(\delta)}{\bar{N}_{t,h}^i}.$$

Bounding the bonus From the definition of the bonus (39), we have for all $(t, i, h) \in [T] \times [M] \times [H]$ such that $\bar{N}_{t,h}^i \geq 2$

$$b_{t,h}(s_{t,h}^i, a_{t,h}^i) = \frac{28\beta^*(\delta)H + 11\beta^c(\delta, \bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i} + \sqrt{\frac{8\beta^*(\delta)}{\bar{N}_{t,h}^i} \cdot \text{Var}_{\hat{\mathbf{P}}_{(r)}(\hat{\mathcal{V}}_{(r),h+1})(s, a)}.$$

Using the inequality (62) of Lemma F.6, we have $\text{Var}_{\hat{\mathbf{P}}_{t,h}}(\hat{\mathcal{V}}_{t,h+1})(s_{t,h}^i, a_{t,h}^i) \leq \text{Var}_{\hat{\mathbf{P}}_{t,h}^c}(\hat{\mathcal{V}}_{t,h+1})(s_{t,h}^i, a_{t,h}^i) + 3H^2\varepsilon_p$. Besides, as for any $a, b \in \mathbb{R}_+$, we have $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$, we get

$$\sqrt{\frac{8\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{t,h}}(\hat{\mathcal{V}}_{t,h+1})(s_{t,h}^i, a_{t,h}^i)}{\bar{N}_{t,h}^i}} \leq \sqrt{\frac{8\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{t,h}^c}(\hat{\mathcal{V}}_{t,h+1})(s_{t,h}^i, a_{t,h}^i)}{\bar{N}_{t,h}^i}} + \sqrt{\frac{24\varepsilon_p H^2 \beta^*(\delta)}{\bar{N}_{t,h}^i}}.$$

Now as for any $a, b \in \mathbb{R}_+$, we have $\sqrt{2ab} \leq a + b$, we get

$$\sqrt{\frac{8\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{t,h}}(\hat{\mathcal{V}}_{t,h+1})(s_{t,h}^i, a_{t,h}^i)}{\bar{N}_{t,h}^i}} \leq \sqrt{\frac{8\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{t,h}^c}(\hat{\mathcal{V}}_{t,h+1})(s_{t,h}^i, a_{t,h}^i)}{\bar{N}_{t,h}^i}} + 3\varepsilon_p H + \frac{4H\beta^*(\delta)}{\bar{N}_{t,h}^i}.$$

Combine everything together After combining all the terms, we have for all $(t, i, h) \in [T] \times [M] \times [H]$ such that $\bar{N}_{t,h}^i \geq 2$

$$\begin{aligned} \delta_{t,h}^i &\leq 7H^2\varepsilon_p + 2\varepsilon_r + \left(1 + \frac{1}{H}\right) \delta_{t,h+1}^i + \left(1 + \frac{1}{H}\right) \zeta_{t,h+1}^i + \frac{22\beta^c(\delta, \bar{N}_{t,h}^i) + 46H\beta^*(\delta) + 2H^2|\mathcal{S}|\beta^{\text{KL}}(\delta, \bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i} \\ &\quad + \sqrt{\frac{4 \text{Var}_{\hat{\mathbf{P}}_{t,h}^c}(\mathcal{V}_{h+1}^{c,*})(s_{t,h}^i, a_{t,h}^i)\beta^*(\delta)}{\bar{N}_{t,h}^i}} + \sqrt{\frac{8\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{t,h}^c}(\hat{\mathcal{V}}_{t,h+1})(s_{t,h}^i, a_{t,h}^i)}{\bar{N}_{t,h}^i}}. \end{aligned}$$

Let us define $\gamma_h = (1 + 1/H)^{H-h}$. Notice that for any $h \geq 0$ it holds $\gamma_h \leq e$. Then by summing and expanding over $h \in [H]$ we have

$$\begin{aligned} \bar{\mathfrak{R}}_h^c(T) &\leq 7eTH^2\varepsilon_p + 2eTH\varepsilon_r + \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h'=h}^H \gamma_{h'-1} \zeta_{t,h'+1}^i + \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h'=h}^H 2eH \mathbf{1}_{[0;1]}(\bar{N}_{t,h'}^i) \\ &\quad + \frac{e}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h'=h}^H \sqrt{\frac{4\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{t,h'}^c}(\mathcal{V}_{h+1}^{c,*})(s_{t,h'}^i, a_{t,h'}^i)}{\bar{N}_{t,h'}^i}} \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h'}^i) \quad =: A_h^T \\ &\quad + \frac{e}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h'=h}^H \sqrt{\frac{8\beta^*(\delta) \text{Var}_{\hat{\mathbf{P}}_{t,h'}^c}(\hat{\mathcal{V}}_{t,h'+1})(s_{t,h'}^i, a_{t,h'}^i)}{\bar{N}_{t,h'}^i}} \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h'}^i) \quad =: B_h^T \\ &\quad + \frac{e}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h'=h}^H \frac{22\beta^c(\delta, \bar{N}_{t,h}^i) + 46H\beta^*(\delta) + 2H^2|\mathcal{S}|\beta^{\text{KL}}(\delta, \bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i} \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h'}^i). \quad =: C_h^T \end{aligned}$$

To conclude the statement, we apply a definition of the event $\mathcal{E}(\delta)$ to the third term in the decomposition above. \square

Lemma D.4. Define $\bar{N}_{t,h}^i = N_{(r_t),h}(s_{t,h}^i, a_{t,h}^i)$. Assume conditions of Lemma C.1. Then, on the event $\mathcal{G}(\delta)$, the following inequalities holds:

$$\begin{aligned} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \frac{\mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i} &\leq 4|\mathcal{S}||\mathcal{A}|H \log \left(\frac{eMTH}{|\mathcal{S}||\mathcal{A}|} \right), \\ \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \sqrt{\frac{\mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i)}{N_{h,i}^{(r,l)}}} &\leq 8H\sqrt{|\mathcal{S}||\mathcal{A}|MT}, \end{aligned}$$

$$\sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \mathbf{1}_{[0;1]}(\bar{N}_{t,h}^i) \leq 4H|\mathcal{S}||\mathcal{A}|.$$

Proof. The quantity $\tilde{N}_{t,h}^i(s, a)$ represents the exact number of visits of the pair (s, a) at step h until episode t , and after the first i agents executed the h -step. We want to bound $\bar{N}_{t,h}^i$ using $\tilde{N}_{t,h}^i(s_{t,h}^i, a_{t,h}^i)$ so that we can compute the latter sums by applying the pigeon-hole principle. To derive such a bound, we distinguish two cases:

Case 1: $N_{r_t,h}(s, a) < \nu(\delta, T)$ In this case, by the synchronization rule described in Algorithm 2, we have $n_{t,h}^i(s, a) < 2n_{(r_t),h}^i(s, a)$. If we sum the latter inequality over all the agents, we obtain $\tilde{N}_{t,h}^M(s, a) \leq 2N_{(r_t),h}(s, a)$. Now using definition of $\tilde{N}_{t,h}^i(s, a)$ yields

$$N_{(r_t),h}(s, a) \leq \tilde{N}_{t,h}^i(s, a) \leq 2N_{(r_t),h}(s, a).$$

Case 2: $N_{r_t,h}(s, a) \geq \nu(\delta, T)$ In this case, the synchronization rule ensures $\hat{N}_{t,h}^i(s, a) \leq 2N_{r_t,h}(s, a)$. Conditioned on $\mathcal{E}^{\text{count}}(\delta)$, we have $\tilde{N}_{t,h}^M(s, a) \leq \frac{10}{7}\hat{N}_{t,h}^i(s, a)$. Combining the two latter inequalities gives

$$N_{(r_t),h}(s, a) \leq \tilde{N}_{t,h}^i(s, a) \leq 4N_{(r_t),h}(s, a),$$

where the lower bound follows directly from the definition of $\tilde{N}_{t,h}^i$. Using the two previous inequalities, we derive the following bound

$$\bar{N}_{t,h}^i \leq \tilde{N}_{t,h}^i(s_{t,h}^i, a_{t,h}^i) \leq 4\bar{N}_{t,h}^i.$$

Applying the latter inequality in the first sum of the lemma yields

$$\sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \frac{\mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i} \leq \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \frac{4 \cdot \mathbf{1}_{[1;+\infty]}(\tilde{N}_{t,h}^i)}{\tilde{N}_{t,h}^i}.$$

By construction, this counter is thus incremented by at most 1 every time and we can apply the pigeon-hole principle on this counter which yields

$$\begin{aligned} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \frac{\mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i} &\leq \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \frac{4 \cdot \mathbf{1}_{[1;+\infty]}(\tilde{N}_{t,h}^i)}{\tilde{N}_{t,h}^i} \leq 4 \sum_{h=1}^H \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}}^{N_{T,h}(s,a)} \sum_{n=1} \frac{1}{n} \\ &\leq 4 \sum_{h=1}^H \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} (\log(N_{T,h}(s, a)) + 1) \leq 4|\mathcal{S}||\mathcal{A}|H \log \left(\frac{eMTH}{|\mathcal{S}||\mathcal{A}|} \right), \end{aligned}$$

where we used the concavity of the logarithm in the last inequality. Similarly, we have

$$\begin{aligned} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \sqrt{\frac{\mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i}} &\leq \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \sqrt{\frac{4 \cdot \mathbf{1}_{[1;+\infty]}(\tilde{N}_{t,h}^i)}{\tilde{N}_{t,h}^i}} \\ &\leq 2 \sum_{h=1}^H \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{n=1}^{N_{T,h}(s,a)} \sqrt{\frac{1}{n}} \leq 8 \sum_{h=1}^H \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sqrt{N_{T,h}(s, a)} \leq 8H \sqrt{|\mathcal{S}||\mathcal{A}|MT}, \end{aligned}$$

where we used the concavity of the square root in the last inequality. Now as $\tilde{N}_{t,h}^i(s_{t,h}^i, a_{t,h}^i) \leq 4\bar{N}_{t,h}^i$, then we have $\mathbf{1}_{[0;1]}(\bar{N}_{t,h}^i) \leq \mathbf{1}_{[0;4]}(\tilde{N}_{t,h}^i)$. Plugging in the latter inequality in the last sum of the lemma yields

$$\sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \mathbf{1}_{[0;1]}(\bar{N}_{t,h}^i) \leq \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \mathbf{1}_{[0;4]}(\tilde{N}_{t,h}^i) \leq \sum_{h=1}^H \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{n=1}^4 1 \leq 4H|\mathcal{S}||\mathcal{A}|.$$

□

For ease of reading, we define $\beta^{\max}(\delta)$ as

$$\beta^{\max}(\delta) := \max \left(\beta^{\text{KL}}(\delta, MT), \beta^c(\delta, MT), \beta^*(\delta), \beta(\delta), \beta^{\text{Var}}(\delta, T), \log \left(\frac{eMT H}{|\mathcal{S}||\mathcal{A}|} \right) \right). \quad (55)$$

Lemma D.5. *Assume conditions of Lemma C.1. Then, on the event $\mathcal{G}(\delta)$, the following inequality holds*

$$\begin{aligned} \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \text{Var}_{\hat{\mathbf{P}}_{t,h}^c} (\mathcal{V}_{h+1}^{c,*})(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i) &\leq 2H^2T + 2H^2U_1^T + 11H^3T\varepsilon_p + 6HT\varepsilon_r \\ &\quad + 30H^3\beta^{\max}(\delta)|\mathcal{S}||\mathcal{A}|^{1/2}T^{1/2}M^{-1/2}, \end{aligned}$$

and we also have

$$\begin{aligned} \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \text{Var}_{\hat{\mathbf{P}}_{t,h}^c} (\hat{\mathcal{V}}_{t,h+1})(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i) &\leq 2H^2T + 2H^2U_1^T + 17H^3T\varepsilon_p + 10HT\varepsilon_r \\ &\quad + 30H^3\beta^{\max}(\delta)|\mathcal{S}||\mathcal{A}|^{1/2}T^{1/2}M^{-1/2}, \end{aligned}$$

where $\beta^{\max}(\delta)$ is defined in Lemma C.1 as a worst-case concentration logarithmic factor.

Proof. Using inequality (62) of Lemma F.6, we have

$$\begin{aligned} \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \text{Var}_{\hat{\mathbf{P}}_{t,h}^c} (\mathcal{V}_{h+1}^{c,*})(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i) &\leq \underbrace{\frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \text{Var}_{\mathbf{P}_h^i} (\mathcal{V}_{h+1}^{c,*})(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i)}_{(\mathbf{W})} \\ &\quad + \underbrace{3H^2 \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \|(\mathbf{P}_h^i - \hat{\mathbf{P}}_{t,h}^c)(s_{t,h}^i, a_{t,h}^i)\|_1 \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i)}_{(\mathbf{X})}. \end{aligned}$$

Term (X): We have for any $(s, a) \in \mathcal{S} \times \mathcal{A}$

$$\|\mathbf{P}_h^i(s, a) - \hat{\mathbf{P}}_{t,h}^c(s, a)\|_1 \leq \|\mathbf{P}_h^i(s, a) - \mathbf{P}_h^c(s, a)\|_1 + \|\mathbf{P}_h^c(s, a) - \hat{\mathbf{P}}_{t,h}^c(s, a)\|_1 \leq \varepsilon_p + \sqrt{\frac{2|\mathcal{S}|\beta^{\text{KL}}(\delta, \bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i}}, \quad (56)$$

where the bound on $\|\mathbf{P}_h^i(s, a) - \mathbf{P}_h^c(s, a)\|_1$ is provided by Lemma F.1 and the bound on $\|\mathbf{P}_h^c(s, a) - \hat{\mathbf{P}}_{t,h}^c(s, a)\|_1$ is provided by the second inequality of Lemma C.2. Thus we get

$$3H^2 \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \|(\mathbf{P}_h^c - \hat{\mathbf{P}}_{t,h}^c)(s_{t,h}^i, a_{t,h}^i)\|_1 \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i) \leq \frac{3H^2}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \varepsilon_p + \sqrt{\frac{2|\mathcal{S}|\beta^{\text{KL}}(\delta, \bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i}} \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i).$$

Finally applying Lemma D.4 yields

$$(\mathbf{X}) \leq 3H^3T\varepsilon_p + 24\sqrt{\frac{2H^6\beta^{\text{KL}}(\delta, MT)|\mathcal{S}|^2|\mathcal{A}|T}{M}}.$$

Term (W): Using inequality (61) of Lemma F.6, we have

$$\begin{aligned} \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \text{Var}_{\mathbf{P}_h^i} (\mathcal{V}_{h+1}^{c,*})(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i) &\leq \underbrace{\frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H 2 \text{Var}_{\mathbf{P}_h^i} (\mathcal{V}_{h+1}^{i,\pi_t})(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i)}_{(\mathbf{Y})} \\ &\leq \underbrace{\frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H 2H \mathbf{P}_h^i |\mathcal{V}_{h+1}^{c,*} - \mathcal{V}_{h+1}^{i,\pi_t}|(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i)}_{(\mathbf{Z})}, \end{aligned}$$

where we recall that $\mathcal{V}_{h+1}^{i,\pi_t}$ is the value function of the policy π_t in the environment of the i -th agent (1). Conditioned on $\mathcal{E}^{\text{Var}}(\delta)$, we have

$$(\mathbf{Y}) \leq \sqrt{\frac{8H^5 T \beta^{\text{Var}}(\delta, T)}{M}} + \frac{6H^3 \beta^{\text{Var}}(\delta, T)}{M} + 2H^2 T.$$

Now by Corollary F.2, we have conditioned on $\mathcal{E}^{\text{optimism}}$ for all $s \in \mathcal{S}$

$$\begin{aligned} |\mathcal{V}_{h+1}^{\mathbf{c},\star}(s) - \mathcal{V}_{h+1}^{i,\pi_t}(s)| &\leq |\mathcal{V}_{h+1}^{\mathbf{c},\pi_t}(s) - \mathcal{V}_{h+1}^{i,\pi_t}(s)| + \mathcal{V}_{h+1}^{\mathbf{c},\star}(s) - \mathcal{V}_{h+1}^{\mathbf{c},\pi_t}(s) \\ &\leq \varepsilon_p H^2 + \varepsilon_r H + \hat{\mathcal{V}}_{t,h+1}(s) - \mathcal{V}_{h+1}^{\mathbf{c},\pi_t}(s) + (2\varepsilon_r + 3\varepsilon_p H)(H - h) \\ &\leq 4\varepsilon_p H^2 + 3\varepsilon_r H + \hat{\mathcal{V}}_{t,h+1}(s) - \mathcal{V}_{h+1}^{\mathbf{c},\pi_t}(s). \end{aligned}$$

Using the definition of $\delta_{t,h}^i$ and $\zeta_{t,h+1}^i$ introduced in Lemma D.3, we have

$$\begin{aligned} (\mathbf{Z}) &\leq \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H 2H(4\varepsilon_p H^2 + 3\varepsilon_r + \delta_{t,h}^i + \zeta_{t,h+1}^i) \\ &\leq 8H^3 T \varepsilon_p + 6HT \varepsilon_r + \sqrt{\frac{8H^5 T \beta(\delta)}{M}} + 2H \sum_{h=1}^H \bar{\mathfrak{R}}_{h+1}^{\mathbf{c}}(T) \leq 8H^3 T \varepsilon_p + 6HT \varepsilon_r + \sqrt{\frac{8H^5 T \beta(\delta)}{M}} + 2H^2 U_1^T, \end{aligned}$$

where the second inequality holds conditioned on $\mathcal{E}(\delta)$. Combining everything yields

$$\begin{aligned} \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \text{Var}_{\mathbf{P}_h^i}(\mathcal{V}_{h+1}^{\mathbf{c},\star})(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i) &\leq (\mathbf{X}) + (\mathbf{Y}) + (\mathbf{Z}) \\ &\leq 11H^3 T \varepsilon_p + 24 \sqrt{\frac{2H^6 \beta^{\text{KL}}(\delta, MT) |\mathcal{S}|^2 |\mathcal{A}| T}{M}} + \sqrt{\frac{8H^5 T \beta^{\text{Var}}(\delta, T)}{M}} + \frac{6H^3 \beta^{\text{Var}}(\delta, T)}{M} + 2H^2 T \\ &\quad + 6HT \varepsilon_r + \sqrt{\frac{8H^5 T \beta(\delta)}{M}} + 2H^2 U_1^T \end{aligned}$$

Now let's move to the second inequality of this lemma. Again by using inequality (62) of Lemma F.6, we have

$$\begin{aligned} \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \text{Var}_{\hat{\mathbf{P}}_{t,h}^i}(\hat{\mathcal{V}}_{t,h+1})(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i) &\leq \underbrace{\frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \text{Var}_{\mathbf{P}_h^i}(\hat{\mathcal{V}}_{t,h+1})(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i)}_{(\mathbf{W}')} \\ &\quad + \underbrace{3H^2 \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \|(\mathbf{P}_h^i - \hat{\mathbf{P}}_{t,h}^i)(s_{t,h}^i, a_{t,h}^i)\|_1 \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i)}_{(\mathbf{X})}. \end{aligned}$$

Term (\mathbf{W}') : Using inequality (61) of Lemma F.6, we have

$$\begin{aligned} \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \text{Var}_{\mathbf{P}_h^i}(\hat{\mathcal{V}}_{t,h+1})(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i) &\leq \underbrace{\frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H 2 \text{Var}_{\mathbf{P}_h^i}(\mathcal{V}_{h+1}^{i,\pi_t})(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i)}_{(\mathbf{Y})} \\ &\leq \underbrace{\frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H 2H \mathbf{P}_h^i |\hat{\mathcal{V}}_{t,h+1} - \mathcal{V}_{h+1}^{i,\pi_t}|(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i)}_{(\mathbf{Z}')}. \end{aligned}$$

Now again by Corollary F.2, we have conditioned on $\mathcal{E}^{\text{optimism}}$ for all $s \in \mathcal{S}$

$$|\hat{\mathcal{V}}_{t,h+1}(s) - \mathcal{V}_{h+1}^{i,\pi_t}(s)| \leq |\mathcal{V}_{h+1}^{\mathbf{c},\pi_t}(s) - \mathcal{V}_{h+1}^{i,\pi_t}(s)| + \hat{\mathcal{V}}_{t,h+1}(s) + (2\varepsilon_r + 3\varepsilon_p H)(H - h) - \mathcal{V}_{h+1}^{\mathbf{c},\pi_t} + (2\varepsilon_r + 3\varepsilon_p H)(H - h)$$

$$\begin{aligned}
 &\leq \varepsilon_p H^2 + \varepsilon_r H + \hat{\mathcal{V}}_{t,h+1}(s) - \mathcal{V}_{h+1}^{c,\pi_t}(s) + 2(2\varepsilon_r + 3\varepsilon_p H)(H - h) \\
 &\leq 7\varepsilon_p H^2 + 5\varepsilon_r H + \hat{\mathcal{V}}_{t,h+1}(s) - \mathcal{V}_{h+1}^{c,\pi_t}(s).
 \end{aligned}$$

By combining the bounds that we have on (\mathbf{X}) , (\mathbf{Y}) , and (\mathbf{Z}') , we derive the following bound

$$\begin{aligned}
 \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \text{Var}_{\hat{\mathbf{p}}_{t,h}^c}(\hat{\mathcal{V}}_{t,h+1})(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i) &\leq (\mathbf{X}) + (\mathbf{Y}) + (\mathbf{Z}') \\
 &\leq 17H^3 T \varepsilon_p + 24 \sqrt{\frac{2H^6 \beta^{\text{KL}}(\delta, MT) |\mathcal{S}|^2 |\mathcal{A}| T}{M}} + \sqrt{\frac{8H^5 T \beta^{\text{Var}}(\delta, T)}{M}} + \frac{6H^3 \beta^{\text{Var}}(\delta, T)}{M} + 10H^2 T \\
 &\quad + 10HT \varepsilon_r + \sqrt{\frac{8H^5 T \beta(\delta)}{M}} + 2H^2 U_1^T.
 \end{aligned}$$

Finally, as we have

$$\begin{aligned}
 24 \sqrt{\frac{2H^6 \beta^{\text{KL}}(\delta, MT) |\mathcal{S}|^2 |\mathcal{A}| T}{M}} &\leq 48H^3 \beta^{\max}(\delta) |\mathcal{S}| |\mathcal{A}|^{1/2} T^{1/2} M^{-1/2} \\
 \sqrt{\frac{8H^5 T \beta(\delta)}{M}} &\leq 3H^3 \beta^{\max}(\delta) |\mathcal{S}| |\mathcal{A}|^{1/2} T^{1/2} M^{-1/2} \\
 \sqrt{\frac{8H^5 T \beta^{\text{Var}}(\delta, T)}{M}} &\leq 3H^3 \beta^{\max}(\delta) |\mathcal{S}| |\mathcal{A}|^{1/2} T^{1/2} M^{-1/2} \\
 \frac{6H^3 \beta^{\text{Var}}(\delta, T)}{M} &\leq 6H^3 \beta^{\max}(\delta) |\mathcal{S}| |\mathcal{A}|^{1/2} T^{1/2} M^{-1/2},
 \end{aligned}$$

then

$$\begin{aligned}
 \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \text{Var}_{\hat{\mathbf{p}}_{t,h}^c}(\hat{\mathcal{V}}_{t,h+1})(s_{t,h}^i, a_{t,h}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h}^i) &\leq 17H^3 T \varepsilon_p + 10HT \varepsilon_r + 2H^2 T + 2H^2 U_1^T \\
 &\quad + 60H^3 \beta^{\max}(\delta) |\mathcal{S}| |\mathcal{A}|^{1/2} T^{1/2} M^{-1/2}.
 \end{aligned}$$

□

Lemma D.6. Assume conditions of Lemma C.1. Then, on the event $\mathcal{G}(\delta)$, the following inequality holds

$$\begin{aligned}
 A_1^T &\leq 23e \cdot \beta^{\max}(\delta) \cdot \sqrt{H^3 |\mathcal{S}| |\mathcal{A}| T M^{-1}} + 23e \beta^{\max}(\delta) \sqrt{H^3 |\mathcal{S}| |\mathcal{A}| U_1^T M^{-1}} \\
 &\quad + 16e \beta^{\max}(\delta) T (6H^2 \varepsilon_p + 3H \varepsilon_r) + 96e H^3 |\mathcal{S}|^{3/2} |\mathcal{A}| M^{-1/2} (\beta^{\max}(\delta))^2, \\
 B_1^T &\leq 46e \cdot \beta^{\max}(\delta) \cdot \sqrt{H^3 |\mathcal{S}| |\mathcal{A}| T M^{-1}} + 46e \beta^{\max}(\delta) \sqrt{H^3 |\mathcal{S}| |\mathcal{A}| U_1^T M^{-1}} \\
 &\quad + 32e \beta^{\max}(\delta) T (6H^2 \varepsilon_p + 3H \varepsilon_r) + 192e H^3 |\mathcal{S}|^{3/2} |\mathcal{A}| M^{-1/2} (\beta^{\max}(\delta))^2, \\
 C_1^T &\leq 272 |\mathcal{S}|^2 |\mathcal{A}| M^{-1} H^3 (\beta^{\max}(\delta))^2.
 \end{aligned}$$

Proof. **Term A_1^T .** To bound the term A_1^T , we start by applying Cauchy-Schwartz inequality

$$\begin{aligned}
 A_1^T &= \frac{e}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h'=h}^H \sqrt{\frac{4\beta^*(\delta) \text{Var}_{\hat{\mathbf{p}}_{t,h'}^c}(\mathcal{V}_{h+1}^{c,*})(s_{t,h'}^i, a_{t,h'}^i)}{\bar{N}_{t,h'}^i}} \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h'}^i) \\
 &\leq \frac{e}{\sqrt{M}} \sqrt{\frac{1}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h'=h}^H \text{Var}_{\hat{\mathbf{p}}_{t,h'}^c}(\mathcal{V}_{h+1}^{c,*})(s_{t,h'}^i, a_{t,h'}^i) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h'}^i)} \sqrt{\sum_{i=1}^M \sum_{t=1}^T \sum_{h'=h}^H \frac{4\beta^*(\delta) \mathbf{1}_{[2;+\infty]}(\bar{N}_{t,h'}^i)}{\bar{N}_{t,h'}^i}},
 \end{aligned}$$

Now, applying Lemma D.5, Lemma D.4, and the subadditivity of the square root, we obtain

$$A_1^T \leq \frac{16\beta^{\max}(\delta)e}{\sqrt{M}} \left(\sqrt{2H^3 |\mathcal{S}| |\mathcal{A}| T} + \sqrt{2H^3 |\mathcal{S}| |\mathcal{A}| U_1^T} + \sqrt{H^2 |\mathcal{S}| |\mathcal{A}| T \cdot (11H^2 \varepsilon_p + 6H \varepsilon_r)} \right)$$

$$+ \sqrt{60H^4|\mathcal{S}|^2|\mathcal{A}|^{3/2}T^{1/2}M^{-1/2} \cdot \beta^{\max}(\delta)} \Bigg) .$$

Next, we analyze the last two terms in the upper bound above. For the third one, by a standard inequality $\sqrt{2ab} \leq a + b$ it holds that

$$\sqrt{H^2|\mathcal{S}||\mathcal{A}|T \cdot (11H^2\varepsilon_p + 6H\varepsilon_r)} \leq 3HT\varepsilon_r + 6H^2T\varepsilon_p + H^2|\mathcal{S}||\mathcal{A}| .$$

Notably, the first two terms already appeared in the regret decomposition; see Lemma D.3. For the last term, the decomposition is more standard

$$\sqrt{60H^4|\mathcal{S}|^2|\mathcal{A}|^{3/2}T^{1/2}M^{-1/2}\beta^{\max}(\delta)} \leq 6\sqrt{H^3|\mathcal{S}||\mathcal{A}|T} + 5H^{5/2}|\mathcal{S}|^{3/2}|\mathcal{A}|M^{-1/2}\beta^{\max}(\delta) .$$

Thus, by a simple rearranging of the terms and applying inequalities $M \geq 1, H \geq 1$, we have

$$\begin{aligned} A_1^T &\leq 23e \cdot \beta^{\max}(\delta) \cdot \sqrt{H^3|\mathcal{S}||\mathcal{A}|TM^{-1}} + 23e\beta^{\max}(\delta)\sqrt{H^3|\mathcal{S}||\mathcal{A}|U_1^TM^{-1}} \\ &\quad + 16e\beta^{\max}(\delta)T(6H^2\varepsilon_p + 3H\varepsilon_r) + 96eH^3|\mathcal{S}|^{3/2}|\mathcal{A}|M^{-1/2}(\beta^{\max}(\delta))^2 . \end{aligned}$$

Term B_1^T . Similarly, the bound for the term B_1^T is derived using a combination of Cauchy-Schwartz, Lemma D.5, Lemma D.4, and the subadditivity of the square root

$$\begin{aligned} B_1^T &\leq 46e \cdot \beta^{\max}(\delta) \cdot \sqrt{H^3|\mathcal{S}||\mathcal{A}|TM^{-1}} + 46e\beta^{\max}(\delta)\sqrt{H^3|\mathcal{S}||\mathcal{A}|U_1^TM^{-1}} \\ &\quad + 32e\beta^{\max}(\delta)T(6H^2\varepsilon_p + 3H\varepsilon_r) + 192eH^3|\mathcal{S}|^{3/2}|\mathcal{A}|M^{-1/2}(\beta^{\max}(\delta))^2 . \end{aligned}$$

Term C_1^T . Finally to estimate C_1^T , we apply Lemma D.4

$$\begin{aligned} \frac{e}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h'=h}^H \frac{22\beta^c(\delta, \bar{N}_{t,h}^i) + 46H\beta^*(\delta) + 2H^2|\mathcal{S}|\beta^{\text{KL}}(\delta, \bar{N}_{t,h}^i)}{\bar{N}_{t,h}^i} 1_{[2;+\infty]}(\bar{N}_{t,h'}^i) \\ \leq \frac{68H^2|\mathcal{S}|\beta^{\max}(\delta)}{M} \sum_{i=1}^M \sum_{t=1}^T \sum_{h=1}^H \frac{1}{\bar{N}_{t,h}^i} 1_{[2;+\infty]}(\bar{N}_{t,h}^i) \leq 272|\mathcal{S}|^2|\mathcal{A}|M^{-1}H^3(\beta^{\max}(\delta))^2 . \end{aligned}$$

□

D.3 Proof of Theorem 4.1

Hereafter, we establish the following bound on the regret

$$\mathfrak{R}(T) \leq 138e\sqrt{H^3|\mathcal{S}||\mathcal{A}|TM^{-1}(\beta^{\max}(\delta))^2} + 6020e^2H^3|\mathcal{S}|^2|\mathcal{A}|(\beta^{\max}(\delta))^2 + e\beta^{\max}(\delta)TH(595H\varepsilon_p + 148H\varepsilon_r) .$$

Proof. Let us start by moving from our regret to a regret connected to a common kernel, using a combination of Corollary F.2 and A-1

$$\mathfrak{R}(T) = \max_{\pi} \frac{1}{M} \sum_{t=1}^T \sum_{i=1}^M \mathcal{V}_1^{i,\pi}(s_{t,1}^i) - \mathcal{V}_1^{i,\pi_t}(s_{t,1}^i) \leq \underbrace{\frac{1}{M} \sum_{t=1}^T \sum_{i=1}^M \mathcal{V}_1^{c,*}(s_{t,1}^i) - \mathcal{V}_1^{c,\pi_t}(s_{t,1}^i)}_{\mathfrak{R}^c(T)} + 2T\varepsilon_pH^2 + 2T\varepsilon_rH .$$

Next, we assume that the event $\mathcal{G}(\delta)$, defined in Lemma C.1 holds. Then Lemma D.2 implies

$$\mathfrak{R}^c(T) \leq \overline{\mathfrak{R}}_1^c(T) + 3T\varepsilon_pH^2 + 2T\varepsilon_rH = \frac{1}{M} \sum_{t=1}^T \sum_{i=1}^M \hat{\mathcal{V}}_{t,1}(s_{t,1}^i) - \mathcal{V}_1^{c,\pi_t}(s_{t,1}^i) + 3T\varepsilon_pH^2 + 2T\varepsilon_rH .$$

By Lemma D.3 we have

$$\overline{\mathfrak{R}}_1^c(T) \leq U_1^T = A_1^T + B_1^T + C_1^T + 7eTH^2\varepsilon_p + 2eTH\varepsilon_r + \sqrt{8H^3T \cdot \beta(\delta)/M} + 4eH^2|\mathcal{S}||\mathcal{A}| ,$$

and, applying Lemma D.6, we have the following quadratic inequality on U_1^T

$$\begin{aligned} U_1^T &\leq 69e \cdot \beta^{\max}(\delta) \cdot \sqrt{H^3|\mathcal{S}||\mathcal{A}|TM^{-1}} + 69e\beta^{\max}(\delta)\sqrt{H^3|\mathcal{S}||\mathcal{A}|U_1^TM^{-1}} \\ &\quad + 144e\beta^{\max}(\delta)T(2H^2\varepsilon_p + \varepsilon_r H) + 288eH^3|\mathcal{S}|^{3/2}|\mathcal{A}|M^{-1/2}(\beta^{\max}(\delta))^2 \\ &\quad + 272|\mathcal{S}|^2|\mathcal{A}|M^{-1}H^3(\beta^{\max}(\delta))^2 + eT(7H^2\varepsilon_p + 2H\varepsilon_r) + \sqrt{8H^3TM^{-1}\beta^{\max}(\delta)} + 4eH^2|\mathcal{S}||\mathcal{A}|. \end{aligned}$$

After some rearranging of the terms, we have the following simplified version

$$\begin{aligned} U_1^T &\leq 69e \cdot \beta^{\max}(\delta)\sqrt{H^3|\mathcal{S}||\mathcal{A}|U_1^TM^{-1}} + 71e \cdot \beta^{\max}(\delta) \cdot \sqrt{H^3|\mathcal{S}||\mathcal{A}|TM^{-1}} \\ &\quad + e\beta^{\max}(\delta)T(295H^2\varepsilon_p + 146H\varepsilon_r) + 560H^3|\mathcal{S}|^2|\mathcal{A}|(\beta^{\max}(\delta))^2. \end{aligned}$$

Finally, using inequality $2ab \leq a^2 + b^2$, we have

$$69e \cdot \beta^{\max}(\delta)\sqrt{H^3|\mathcal{S}||\mathcal{A}|U_1^TM^{-1}} \leq \frac{1}{2}U_1^T + 2450e^2H^3|\mathcal{S}||\mathcal{A}|M^{-1}(\beta^{\max}(\delta))^2,$$

thus

$$U_1^T \leq 138e\sqrt{H^3|\mathcal{S}||\mathcal{A}|TM^{-1}(\beta^{\max}(\delta))^2} + 6020e^2H^3|\mathcal{S}|^2|\mathcal{A}|(\beta^{\max}(\delta))^2 + e\beta^{\max}(\delta)TH(590H\varepsilon_p + 144H\varepsilon_r).$$

□

E COMMUNICATION COMPLEXITY

In the sequel, we prove the bound on the communication complexity of **Fed-UCBVI** stated in Lemma 4.1.

Lemma 4.1 (Communication Complexity). *With probability at least $1 - \delta$, the number of communication rounds of **Fed-UCBVI** is bounded by*

$$\begin{aligned} \mathfrak{C}(T) &\leq \mathcal{O}(|\mathcal{S}||\mathcal{A}|H \log T + M|\mathcal{S}||\mathcal{A}|H \log \log T \\ &\quad + M|\mathcal{S}||\mathcal{A}|H \log(1 + \varepsilon_p T)), \end{aligned}$$

where logarithmic dependence in $|\mathcal{S}|, |\mathcal{A}|, H, 1/\delta$ and M is ignored.

Proof. Let us fix $(s, a, h) \in |\mathcal{S}| \times \mathcal{A} \times [H]$ and bound the maximum number of abortion signals triggered by this triplet. We define R as the value of the variable r that indicates the current round of the communication, defined in **Fed-UCBVI**, during iteration T . Let us also define $k_{s,a,h,i}$ as the number of times agent i triggered the synchronization rule because of the triplet (s, a, h) and $k_{s,a,h}$ the number of times the synchronization rule was triggered because of the triplet (s, a, h) . Recall that

$$\nu(\delta, T) = 14\varepsilon_p THM + 182M\beta^c(\delta, T). \quad (57)$$

We distinguish two cases:

- (1) $N_{(R),h}(s, a) \leq \nu(\delta, T)$: Thus, it holds that $2^{k_{s,a,h,i}} \leq n_{(R),h}^i \leq \nu(\delta, T)$. Thereby $k_{s,a,h,i} \leq \log(\nu(\delta, T))$. Hence, we have $k_{s,a,h} \leq M \log(\nu(\delta, T))$.
- (2) $N_{(R),h}(s, a) > \nu(\delta, T)$: In this case, we can define

$$k_{s,a,h}^{\min} = \min\{r \in [R] : N_{(r),h}(s, a) \leq \nu(\delta, T) \text{ and } N_{(r+1),h}(s, a) > \nu(\delta, T)\}.$$

By the precedent case, we have $k_{s,a,h}^{\min} \leq M \log(\nu(\delta, T))$. Now let us denote by r_1, \dots, r_p , where $p = k_{s,a,h} - k_{s,a,h}^{\min}$, the indices of the rounds where the synchronization rule was triggered because of the triplet (s, a, h) starting from round $k_{s,a,h}^{\min}$. Thus, for a certain $i \in [M]$, we have

$$\hat{N}_{(r_{t+1}),h}^i(s, a) \geq 2N_{(r_{t+1}-1),h}(s, a).$$

Under the event $\mathcal{E}^{\text{count}}(\delta)$, we have for any $t \in [1; p]$,

$$\tilde{N}_{(r_{t+1}),h}^M(s,a) \geq \frac{3}{7} \hat{N}_{(r_{t+1}),h}^i(s,a) + \frac{1}{7} \nu(\delta, T) \geq \frac{4}{7} \hat{N}_{(r_{t+1}),h}^i(s,a),$$

where $\tilde{N}_{(r_{t+1}),h}^M(s,a)$ is defined in (29). Combining the two previous inequalities, it gives

$$\tilde{N}_{(r_{t+1}),h}^M(s,a) \geq 2 \cdot (4/7) N_{(r_{t+1}-1),h}(s,a) \geq (8/7) N_{(r_t),h}(s,a) = (8/7) \tilde{N}_{(r_t),h}^M(s,a),$$

where the second inequality comes from $r_{t+1} > r_t$ and monotonicity of the counters. Unrolling the last recursion yields

$$TM \geq \tilde{N}_{(r_p),h}^M(s,a) \geq (8/7)^{k_{s,a,h} - k_{s,a,h}^{\min}} \nu(\delta, T).$$

Thus, we obtain

$$k_{s,a,h} \leq k_{s,a,h}^{\min} + \frac{\log(TM/\nu(\delta, T))}{\log(8/7)} \leq M \log(\nu(\delta, T)) + \frac{\log(TM/\nu(\delta, T))}{\log(8/7)},$$

which yields

$$\mathfrak{C}(T) \leq R_{\max} := M|\mathcal{S}||\mathcal{A}|H \log(\nu(\delta, T)) + |\mathcal{S}||\mathcal{A}|H \frac{\log(TM/\nu(\delta, T))}{\log(8/7)}. \quad (58)$$

□

F TECHNICAL LEMMAS

Lemma F.1 (ℓ_1 -norm Bound). *Assume A-1, then*

$$\max_{(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]} \|\mathbf{P}_h^c(\cdot|s,a) - \mathbf{P}_h^i(\cdot|s,a)\|_1 \leq \varepsilon_{\mathbf{p}}.$$

Proof. Let $(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]$. Using A-1, we get

$$\|\mathbf{P}_h^c(\cdot|s,a) - \mathbf{P}_h^i(\cdot|s,a)\|_1 = \sum_{s' \in \mathcal{S}} |\mathbf{P}_h^c(s'|s,a) - \mathbf{P}_h^i(s'|s,a)| = \sum_{s' \in \mathcal{S}} \varepsilon_{\mathbf{p}} |\mathbf{P}_h^c(s'|s,a) - \mathbf{P}_h^{\text{ind},i}(s'|s,a)| \leq \varepsilon_{\mathbf{p}}.$$

□

Lemma F.2. *For any policy π , for any $(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]$, and for any $(i,j) \in [M]^2$, we have*

$$|d_{s_0,h}^{i,\pi}(s,a) - d_{s_0,h}^{j,\pi}(s,a)| \leq \varepsilon_{\mathbf{p}} H.$$

Proof. Let us consider two following MDPs $\mathcal{M}_1 = (\mathcal{S}, \mathcal{A}, H, \{\mathbf{P}_{h',i}\}_{1 \leq h' \leq H}, \{\mathbf{1}_{(s,a)}(\cdot) \mathbf{1}_h(h')\}_{1 \leq h' \leq H})$ and $\mathcal{M}_2 = (\mathcal{S}, \mathcal{A}, H, \{\mathbf{P}_{h',j}\}_{1 \leq h' \leq H}, \{\mathbf{1}_{(s,a)}(\cdot) \mathbf{1}_h(h')\}_{1 \leq h' \leq H})$. Let's denote by $\tilde{V}_h^{i,\pi}$ and $\tilde{V}_h^{j,\pi}$ the values function associated with the policy π in these two respective environments. We have

$$\tilde{V}_h^{i,\pi}(s_0) = \mathbb{E}_{\pi} \left[\sum_{h'=1}^H \mathbf{1}_{(s,a)}(s_h^i, a_h^i) \mathbf{1}_h(h') \right] = \mathbb{E}_{\pi} [\mathbf{1}_{(s,a)}(s_h^i, a_h^i)] = d_{s_0,h}^{i,\pi}(s,a).$$

Similarly, we have $\tilde{V}_h^{j,\pi}(s_0) = d_{s_0,h}^{j,\pi}(s,a)$. Finally applying Lemma F.8 combined with Holder's inequality, and the fact that $\|\tilde{V}_h^{j,\pi}\|_{\infty} \leq 1$ yields

$$|d_{s_0,h}^{i,\pi}(s,a) - d_{s_0,h}^{j,\pi}(s,a)| \leq \varepsilon_{\mathbf{p}} H.$$

□

F.1 Bellman type equations for the variance

For a deterministic policy π and an agent i , we recall the following definitions of the Bellman-type equations for the variances as follows

$$\begin{aligned}\sigma \mathcal{Q}_h^{i,\pi}(s, a) &:= \text{Var}_{\mathbf{P}_h^i}(\mathcal{V}_{h+1}^{i,\pi})(s, a) + \mathbf{P}_h^i \sigma \mathcal{V}_{h+1}^{i,\pi}(s, a) \\ \sigma \mathcal{V}_h^{i,\pi}(s) &:= \sigma \mathcal{Q}_h^{i,\pi}(s, \pi(s)) \\ \sigma \mathcal{V}_{H+1}^{i,\pi}(s) &:= 0,\end{aligned}\tag{59}$$

where $\text{Var}_{\mathbf{P}_h^i}(f)(s, a) := \mathbb{E}_{s' \sim \mathbf{P}_h^i(\cdot|s, a)} \left[(f(s') - \mathbf{P}_h^i f(s, a))^2 \right]$ denotes the variance operator. Unrolling the precedent relation yields

$$\sigma \mathcal{V}_1^{i,\pi}(s) = \sum_{h=1}^H \sum_{s', a'} d_{s,h}^{i,\pi}(s', a') \text{Var}_{\mathbf{P}_h^i}(\mathcal{V}_{h+1}^{i,\pi})(s', a'),$$

where $d_{s,h}^{i,\pi}(s', a')$ is the probability of visiting a pair (s', a') in the i -th environment while following the policy π and starting from a state s . Next, we state the well-known Bellman equation for variances (see, e.g., [Sobel 1982](#); [Azar et al. 2017](#)).

Lemma F.3. *For any deterministic policy π , for all $h \in [H]$, and for all $i \in [M]$,*

$$\mathbb{E}_\pi \left[\left(\sum_{h'=h}^H r_{h'}^i(s_{h'}^i, a_{h'}^i) - \mathcal{Q}_h^{i,\pi}(s_h^i, a_h^i) \right)^2 \middle| (s_h^i, a_h^i) = (s, a) \right] = \mathcal{Q}_h^{i,\pi}(s, a).\tag{60}$$

In particular,

$$\mathbb{E}_\pi \left[\left(\sum_{h=1}^H r_h^i(s_h^i, a_h^i) - \mathcal{V}_1^{i,\pi}(s_1^i) \right)^2 \right] = \sigma \mathcal{V}_1^{i,\pi}(s_1^i) = \sum_{h=1}^H \sum_{s, a} d_h^{i,\pi}(s, a) \text{Var}_{\mathbf{P}_h^i}(\mathcal{V}_{h+1}^{i,\pi})(s, a).$$

F.2 Concentration inequalities

Lemma F.4 (Deviation inequality for categorical distribution, [Jonsson et al. 2020](#)). *Let $(X_t)_{t \in \mathbb{N}^*}$ be i.i.d. samples from a probability measure \mathbf{P} supported on $\{1, \dots, m\}$. We denote by $\hat{\mathbf{P}}_n$ the empirical vector of probabilities, i.e., for all $k \in \{1, \dots, m\}$,*

$$\hat{\mathbf{P}}_n(k) := \frac{1}{n} \sum_{\ell=1}^n \mathbf{1}_{\{k\}}(X_\ell).$$

For all \mathbf{P} and for all $\delta \in (0, 1)$,

$$\mathbb{P} \left(\exists n \in \mathbb{N}^*, n \text{KL}(\hat{\mathbf{P}}_n \| \mathbf{P}) > \log(1/\delta) + (m-1) \log(e(1 + n/(m-1))) \right) \leq \delta.$$

Lemma F.5 (Corollary 11 by [Talebi and Maillard 2018](#)). *Let \mathbf{P}, \mathbf{Q} two probability distributions on \mathcal{S} . For all functions $f : \mathcal{S} \mapsto [0, H]$,*

$$\mathbf{P}f - \mathbf{Q}f \leq \sqrt{2 \text{Var}_{\mathbf{Q}}(f) \text{KL}(\mathbf{P} \| \mathbf{Q})} + \frac{2}{3} H \text{KL}(\mathbf{P}, \mathbf{Q}).$$

where we have defined $\mathbf{P}f := \mathbb{E}_{s \sim \mathbf{P}}[f(s)]$.

Lemma F.6 (Lemma H.9 by [Tiapkin et al. 2023](#)). *For any two probability measures \mathbf{P}, \mathbf{Q} on \mathcal{S} , for $f, g : \mathcal{S} \mapsto [0, b]$ two functions defined on \mathcal{S} , we have that*

$$\text{Var}_{\mathbf{P}}(f) \leq 2 \text{Var}_{\mathbf{P}}(g) + 2b\mathbf{P}|f - g| \quad \text{and} \tag{61}$$

$$\text{Var}_{\mathbf{Q}}(f) \leq \text{Var}_{\mathbf{P}}(f) + 3b^2 \|\mathbf{P} - \mathbf{Q}\|_1 \tag{62}$$

where we denote the absolute operator by $|f|(s) = |f(s)|$ for all $s \in \mathcal{S}$.

Lemma F.7 (Theorem 4 by [Maurer and Pontil 2009](#)). . Consider any $\delta > 0$ and any integer $n \geq 2$. Let Y, Y_1, \dots, Y_n be a collection of i.i.d. random variables falling within $[0, 1]$. Define the empirical mean $\bar{Y} := \frac{1}{n} \sum_{i=1}^n Y_i$ and empirical variance $\hat{Y}_n := \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2$. Then we have

$$\mathbb{P} \left[\left| \mathbb{E}[Y] - \frac{1}{n} \sum_{i=1}^n Y_i \right| > \sqrt{\frac{2\hat{Y}_n \log(2/\delta)}{n-1}} + \frac{7 \log(2/\delta)}{3(n-1)} \right] \leq \delta$$

Below, we state the self-normalized Bernstein-type inequality by [Domingues et al. \(2021c\)](#). Let $(Y_t)_{t \in \mathbb{N}^*}, (w_t)_{t \in \mathbb{N}^*}$ be two sequences of random variables adapted to a filtration $(\mathcal{F}_t)_{t \in \mathbb{N}}$. We assume that the weights are in the unit interval $w_t \in [0, 1]$ and predictable, i.e. \mathcal{F}_{t-1} measurable. We also assume that the random variables Y_t are bounded $|Y_t| \leq b$ and centered $\mathbb{E}[Y_t | \mathcal{F}_{t-1}] = 0$. Consider the following quantities

$$S_t := \sum_{s=1}^t w_s Y_s, \quad V_t := \sum_{s=1}^t w_s^2 \cdot \mathbb{E}[Y_s^2 | \mathcal{F}_{s-1}],$$

and let $h(x) \triangleq (x+1) \log(x+1) - x$ be the Cramér transform of a Poisson distribution of parameter 1.

Theorem F.1 (Anytime Bernstein-type concentration inequality for martingales). For all $\delta > 0$,

$$\mathbb{P} \left[\exists t \geq 1, (V_t/b^2 + 1)h \left(\frac{b|S_t|}{V_t + b^2} \right) \geq \log(1/\delta) + \log(4e(2t+1)) \right] \leq \delta.$$

The previous inequality can be weakened to obtain a more explicit bound: if $b \geq 1$ with probability at least $1 - \delta$, for all $t \geq 1$,

$$|S_t| \leq \sqrt{2V_t \log(4e(2t+1)/\delta)} + 3b \log(4e(2t+1)/\delta).$$

Next, we apply this Bernstein inequality to a particular distribution. Let \mathcal{F}_t for $t \in \mathbb{N}$ be a filtration and $(X_t)_{t \in \mathbb{N}^*}$ be a sequence of Bernoulli random variables with $\mathbb{P}(X_t = 1 | \mathcal{F}_{t-1}) = P_t$ with P_t being \mathcal{F}_{t-1} -measurable and X_t being \mathcal{F}_t -measurable.

Corollary F.1. For all $\delta > 0$,

$$\mathbb{P} \left(\exists n : \left| \sum_{t=1}^n X_t - P_t \right| > \frac{1}{8} \sum_{t=1}^n P_t + 11 \log \left(\frac{4e(2n+1)}{\delta} \right) \right) \leq \delta.$$

Proof. Given a simplified version, we have with probability at least $1 - \delta$ by applying inequality $2ab \leq a^2 + b^2$ for $a, b \geq 0$

$$\left| \sum_{t=1}^n X_t - P_t \right| \leq \sqrt{2 \cdot \frac{V_n}{8} \cdot 8 \log(4e(2n+1)/\delta)} + 3 \log(4e(2n+1)/\delta) \leq \frac{1}{8} \sum_{t=1}^n P_t + 11 \log(4e(2n+1)/\delta).$$

□

F.3 Performance-difference Lemma

Lemma F.8 (Lemma 3 of [Russo 2019](#)). Let us consider two MDPs $\mathcal{M}_1 = (\mathcal{S}, \mathcal{A}, H, \mathbf{r}^{(1)}, \mathbf{P}^{(1)})$ and $\mathcal{M}_2 = (\mathcal{S}, \mathcal{A}, H, \mathbf{r}^{(2)}, \mathbf{P}^{(2)})$. Let $\mathcal{V}_1^{(1),\pi}(s)$ and $\mathcal{V}_1^{(2),\pi}(s)$ are values of a fixed policy π in MDP \mathcal{M}_1 and \mathcal{M}_2 respectively. Then it holds

$$\mathcal{V}_1^{(1),\pi}(s) - \mathcal{V}_1^{(2),\pi}(s) = \mathbb{E}_{\pi, \mathcal{M}_1} \left[\sum_{h=1}^H \left(\mathbf{r}_h^{(1)} - \mathbf{r}_h^{(2)} \right) (s_h, a_h) + \left(\mathbf{P}_h^{(1)} - \mathbf{P}_h^{(2)} \right) \mathcal{V}_{h+1}^{(2),\pi}(s_h, a_h) \right],$$

where expectation is taken over the trajectories $(s_1, a_1, \dots, s_H, a_H)$ generated by policy π in an MDP \mathcal{M}_1 .

Corollary F.2. *Let us consider two MDPs $\mathcal{M}_1 = (\mathcal{S}, \mathcal{A}, H, r^{(1)}, P^{(1)})$ and $\mathcal{M}_2 = (\mathcal{S}, \mathcal{A}, H, r^{(2)}, P^{(2)})$, such that $\forall (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H] : |r_h^{(1)}(s, a) - r_h^{(2)}(s, a)| \leq \varepsilon_r$, $|\mathcal{V}_h^{(1), \pi}(s)| \leq c$, $|\mathcal{V}_h^{(2), \pi}(s)| \leq c$, and $\|P_h^{(1)}(s, a) - P_h^{(2)}(s, a)\|_1 \leq \varepsilon_P$ where $c > 0$ is a positive constant and $\mathcal{V}_1^{(1), \pi}(s)$ and $\mathcal{V}_1^{(2), \pi}(s)$ are values of a fixed policy π in MDP \mathcal{M}_1 and \mathcal{M}_2 respectively. Then it holds*

$$\mathcal{V}_1^{(1), \pi}(s_1) - \mathcal{V}_1^{(2), \pi}(s_1) \leq \varepsilon_P cH + \varepsilon_r H.$$

Proof. Follows directly from combination of Lemma F.8, Holder's inequality and a fact that $\|\mathcal{V}_h^{(2), \pi}\|_1 \leq c$. \square

Inspired by a construction of Ross and Bagnell (2010), we can show that dependence H^2 in terms of ℓ_1 -distance between two models is non-improvable.

Lemma F.9. *There exist two MDPs $\mathcal{M}_1 = (\mathcal{S}, \mathcal{A}, H, r, P^1)$ and $\mathcal{M}_2 = (\mathcal{S}, \mathcal{A}, H, r, P^2)$ with the same reward function and different kernels, $H \geq 2$ such that $\forall (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H] : \|P_h^1(s, a) - P_h^2(s, a)\|_1 \leq \varepsilon_P$ for $0 < \varepsilon_P < 2/H$. Then there is a policy π such that values $\mathcal{V}_1^{1, \pi}(s)$ and $\mathcal{V}_1^{2, \pi}(s)$ in MDPs \mathcal{M}_1 and \mathcal{M}_2 satisfy*

$$\mathcal{V}_1^{1, \pi}(s_1) - \mathcal{V}_1^{2, \pi}(s_1) = \Omega(\varepsilon_P H^2).$$

Proof. Consider the problem with 2 states $\{s_1, s_2\}$ and 1 action $\{a\}$, the agent always starts at s_1 . The reward function satisfies $r_h(s_1, a) = 1, r_h(s_2, a) = 0$ for all $h \in H$. Finally, the transition kernels are the same for all h and are defined as

$$P_h^i(s_1|s_1, a) = 1 - p_i, \quad P_h^i(s_2|s_1, a) = p_i, \quad P_h^i(s_1|s_2, a) = 0, \quad P_h^i(s_2|s_2, a) = 1,$$

for $i \in \{1, 2\}$. In other words, the state s_2 is a sink with zero reward. Since there is only one action, the value is the same for any policy π . Let us take $p_1 = 0$ and $p_2 = \varepsilon_P$, then under the kernel P^1 the value $\mathcal{V}_1^{1, \pi}(s_1)$ is equal to H , whereas under the kernel P^2 , the value function $\mathcal{V}_1^{1, \pi}(s_1)$ it is equal to

$$\mathcal{V}_1^{2, \pi}(s_1) = 1 + (1 - \varepsilon_P) + (1 - \varepsilon_P)^2 + \dots + (1 - \varepsilon_P)^{H-1} = \frac{1 - (1 - \varepsilon_P)^H}{\varepsilon_P}.$$

Then we have

$$\mathcal{V}_1^{1, \pi}(s_1) - \mathcal{V}_1^{2, \pi}(s_1) = \frac{H\varepsilon_P - 1 + (1 - \varepsilon_P)^H}{\varepsilon_P}.$$

Now as $0 < \varepsilon_P < 2/H$, Bernoulli's inequality yields

$$\mathcal{V}_1^{1, \pi}(s_1) - \mathcal{V}_1^{2, \pi}(s_1) = \frac{H\varepsilon_P - 1 + (1 - \varepsilon_P)^{H/2}(1 - \varepsilon_P)^{H/2}}{\varepsilon_P} \geq \frac{H\varepsilon_P - 1 + (1 - H\varepsilon_P/2)(1 - H\varepsilon_P/2)}{\varepsilon_P} = \frac{\varepsilon_P H^2}{4},$$

where the first inequality comes from $(1 - x)^r \geq 1 - rx$ for $0 \leq x \leq 1$ and $r > 1$. \square