

# Reward-agnostic Fine-tuning: Provable Statistical Benefits of Hybrid Reinforcement Learning

Gen Li<sup>\*†</sup>    Wenhao Zhan<sup>\*‡</sup>    Jason D. Lee<sup>‡</sup>    Yuejie Chi<sup>§</sup>    Yuxin Chen<sup>†</sup>

May 17, 2023

## Abstract

This paper studies tabular reinforcement learning (RL) in the hybrid setting, which assumes access to both an offline dataset and online interactions with the unknown environment. A central question boils down to how to efficiently utilize online data to strengthen and complement the offline dataset and enable effective policy fine-tuning. Leveraging recent advances in reward-agnostic exploration and offline RL, we design a three-stage hybrid RL algorithm that beats the best of both worlds — pure offline RL and pure online RL — in terms of sample complexities. The proposed algorithm does not require any reward information during data collection. Our theory is developed based on a new notion called *single-policy partial concentrability*, which captures the trade-off between distribution mismatch and miscoverage and guides the interplay between offline and online data.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Main contributions . . . . .	3
1.2	Notation . . . . .	3
<b>2</b>	<b>Preliminaries and problem settings</b>	<b>4</b>
<b>3</b>	<b>Algorithm</b>	<b>6</b>
3.1	A three-stage algorithm . . . . .	7
3.2	Subroutine for solving the subproblem (21b) . . . . .	9
<b>4</b>	<b>Main results</b>	<b>11</b>
<b>5</b>	<b>Related works</b>	<b>13</b>
<b>6</b>	<b>Analysis of Theorem 1</b>	<b>15</b>
6.1	Step 1: establishing the proximity of $\widehat{d}^\pi$ (resp. $\widehat{d}^{\text{off}}$ ) and $d^\pi$ (resp. $d^{\text{off}}$ ) . . . . .	15
6.2	Step 2: showing that $\pi^{\text{imitate}}$ (resp. $\pi^{\text{explore}}$ ) covers $\widehat{d}^{\text{off}}$ (resp. $d^{\pi^*}$ ) adequately . . . . .	16
6.3	Step 3: establishing the performance of offline RL . . . . .	18
<b>7</b>	<b>Discussion</b>	<b>22</b>
<b>A</b>	<b>Useful algorithmic subroutines from prior works</b>	<b>27</b>
A.1	Subroutine: occupancy estimation for any policy $\pi$ . . . . .	27
A.2	Subroutine: reward-agnostic online exploration . . . . .	27
A.3	Subroutine: pessimistic model-based offline RL . . . . .	28

---

<sup>\*</sup>The first two authors contributed equally.

<sup>†</sup>Department of Statistics and Data Science, Wharton School, University of Pennsylvania, Philadelphia, PA 19104, USA.

<sup>‡</sup>Department of Electrical and Computer Engineering, Princeton University, Princeton, NJ 08544, USA.

<sup>§</sup>Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

<b>B Proof for the stopping criterion and the iteration complexity for solving (21b)</b>	<b>30</b>
<b>C Proofs of technical lemmas</b>	<b>33</b>
C.1 Proof of Lemma 2	33
C.2 Proof of Lemma 6	33

# 1 Introduction

As reinforcement learning (RL) shows promise in achieving super-human empirical success across diverse fields (e.g., games (Berner et al., 2019; Mnih et al., 2013; Silver et al., 2016; Vinyals et al., 2019), robotics (Brambilla et al., 2013), autonomous driving (Shalev-Shwartz et al., 2016)), theoretical understanding about RL has also been substantially expanded, with the aim of distilling fundamental principles that can inform and guide practice. Among all sorts of theoretical questions being pursued, how to make the best use of data emerges as a question of profound interest for problems with enormous dimensionality.

There are at least two mainstream mechanisms when it comes to data collection: online RL and offline RL. Let us briefly describe their attributes and differences as follows.

*Online RL.* In this setting, an agent learns how to maximize her cumulative reward through interaction with the unknown environment (by, say, executing a sequence of adaptively chosen actions and utilizing the instantaneous feedback of the environment). Given that all information about the environment is obtained through real-time data collection, the main challenge lies in how to (optimally) manage the trade-off between exploration and exploitation. Towards this, one popular approach advertises the principle of optimism in the face of uncertainty — e.g., employing upper confidence bounds during value estimation to guide exploration — whose effectiveness has been shown for both the tabular case (Auer and Ortner, 2006; Azar et al., 2017; Bai et al., 2019; Dann et al., 2017; Dong et al., 2019; Jaksch et al., 2010; Jin et al., 2018; Li et al., 2021b; Ménard et al., 2021b; Zhang et al., 2020b) and the case with function approximation (Du et al., 2021; Foster et al., 2021; Jin et al., 2021a, 2020b; Li et al., 2021a; Zanette et al., 2020; Zhou et al., 2021a).

*Offline RL.* In contrast, offline RL assumes access to a pre-collected dataset, without given permission to perform any further data collection. The feasibility of reliable offline RL depends heavily on the quality of the dataset at hand. A central challenge stems from the presence of distribution shift: the distribution of the offline dataset might differ significantly from that induced by the target policy. Another common challenge arises from insufficient data coverage: a nontrivial fraction of the state-action pairs might be inadequately visited in the available dataset, thus precluding one from faithfully evaluating many policies based solely on the offline dataset. To circumvent these obstacles, recent works proposed the principle of pessimism in the face of uncertainty, recommending caution when selecting poorly visited actions (Jin et al., 2021b; Kumar et al., 2020; Li et al., 2022; Liu et al., 2020; Rashidinejad et al., 2021; Shi et al., 2022; Uehara and Sun, 2021; Yin et al., 2021). Without requiring uniform coverage of all policies, the pessimism approach proves effective as long as the so-called *single-policy concentrability* is satisfied, which only assumes adequate coverage over the part of the state-action space reachable by the desirable policy.

In reality, however, both mechanisms above come with limitations. For instance, even the single-policy concentrability requirement might be too stringent (and hence fragile) for offline RL, as it is not uncommon for the historical dataset to miss a small yet essential part of the state-action space. Pure online RL might also be overly restrictive, given that there might be information from past data that could help initialize online exploration and mitigate the burden of further data collection.

All this motivates the studies of hybrid RL, a scenario where the agent has access to an offline dataset while, in the meantime, (limited) online data collection is permitted as well. Oftentimes, this scenario is practically not only feasible but also appealing: on the one hand, offline data provides useful information for policy pre-training, while further online exploration helps enrich existing data and allows for effective policy fine-tuning. As a matter of fact, multiple empirical works (Hester et al., 2018; Kalashnikov et al., 2018; Nair et al., 2020, 2018; Rajeswaran et al., 2017; Vecerik et al., 2017) indicated that combining online RL with offline datasets outperforms both pure online RL and pure offline RL. Nevertheless, theoretical pursuits about hybrid RL are lagging behind. Two recent works Ross and Bagnell (2012); Xie et al. (2021b) studied a restricted setting, where the agent is aware of a Markovian behavior policy (a policy that generates

offline data) and can either execute the behavior policy or any other adaptive choice to draw samples in each episode; in this case, Xie et al. (2021b) proved that under the single-policy concentrability assumption of the offline dataset, having perfect knowledge about the behavior policy does not improve online exploration in the minimax sense. Another strand of works Nakamoto et al. (2023); Song et al. (2022); Wagenmaker and Pacchiano (2022) looked at a more general offline dataset and investigated how to leverage offline data in online exploration. From the sample complexity viewpoint, Wagenmaker and Pacchiano (2022) studied the statistical benefits of hybrid RL in the presence of linear function approximation; the result therein, however, required strong assumptions on data coverage (i.e., all-policy concentrability) and fell short of unveiling provable gains in the tabular case (as we shall elucidate momentarily). In light of such theoretical inadequacy in previous works, this paper is motivated to pursue the following question:

*Does hybrid RL allow for improved sample complexity compared to pure online or offline RL in the tabular case?*

## 1.1 Main contributions

We deliver an affirmative answer to the above question. Further relaxing the single-policy concentrability assumption, we introduce a relaxed notation called single-policy *partial* concentrability (to be made precise in Definition 2), which (i) allows the dataset to miss a fraction of the state-action space visited by the optimal policy and (ii) captures the tradeoff between distribution mismatch and lack of coverage. Armed with this notion, our results reveal provable statistical benefits of hybrid RL compared with both pure online and offline RL. The main contributions are summarized below.

*A novel three-stage algorithm.* We design a new hybrid RL algorithm consisting of three stages. In the first stage, we obtain crude estimation of the occupancy distribution  $d^\pi$  w.r.t. any policy  $\pi$  as well as the data distribution  $d^{\text{off}}$  of the offline dataset. The second stage performs online exploration; in particular, we execute one exploration policy to imitate the offline dataset and another one to explore the inadequately visited part of the unknown environment, with both policies computed by approximately solving convex optimization sub-problems. Notably, these two stages do not count on the availability of reward information, and thus operate in a reward-agnostic manner. The final stage then invokes the state-of-the-art offline RL algorithm for policy learning, on the basis of all data we have available (including both online and offline data).

*Computationally efficient subroutines.* Throughout the first two stages of the algorithm, we need to solve a couple of convex sub-problems with exponentially large dimensions. In order to attain computational efficiency, we design efficient Frank-Wolfe-type paradigms to solve the sub-problems approximately, which run in polynomial time. This plays a crucial role in ensuring computational tractability of the proposed three-stage algorithm.

*Improved sample complexity.* We characterize the sample complexity of our algorithm (see Theorem 1), which provably improves upon both pure online and offline RL. On the one hand, hybrid RL achieves strictly enhanced performance compared to pure offline RL (assuming the same sample size) when the offline data falls short of covering all state-action pairs reachable by the desired policy. On the other hand, the sample size allocated to online exploration in our algorithm might only need to be proportional to the fraction  $\sigma$  of the state-action space uncovered by the offline dataset, thus resulting in sample size saving in general compared to pure online RL (a case with  $\sigma = 1$ ).

**Notation.** Let us also introduce several useful notation. For integer  $m > 0$ , we let  $[m]$  represent the set  $\{1, \dots, m\}$ . For any set  $\mathcal{B}$ , we denote by  $\mathcal{B}^c$  its complement. For any policy  $\pi_0$ , we let  $\mathbb{1}_{\pi_0} : \Pi \rightarrow \{0, 1\}$  be an indicator function such that  $\mathbb{1}_{\pi_0}(\pi) = 1$  if  $\pi = \pi_0$  and  $\mathbb{1}_{\pi_0}(\pi) = 0$  otherwise. For any finite set  $\mathcal{A}$ , we denote by  $\Delta(\mathcal{A})$  the probability simplex over  $\mathcal{A}$ . Letting  $\mathcal{X} := (S, A, H, \frac{1}{\epsilon}, \frac{1}{\delta})$ , we use the notation  $f(\mathcal{X}) = O(g(\mathcal{X}))$  or  $f(\mathcal{X}) \lesssim g(\mathcal{X})$  to indicate the existence of a universal constant  $C_1 > 0$  such that  $f \leq C_1 g$ , the notation  $f(\mathcal{X}) \gtrsim g(\mathcal{X})$  to indicate that  $g(\mathcal{X}) = O(f(\mathcal{X}))$ , and the notation  $f(\mathcal{X}) \asymp g(\mathcal{X})$  to mean that  $f(\mathcal{X}) \lesssim g(\mathcal{X})$  and  $f(\mathcal{X}) \gtrsim g(\mathcal{X})$  hold simultaneously. The notation  $\tilde{O}(\cdot)$  is defined in the same way as  $O(\cdot)$  except that it hides logarithmic factors.

## 2 Preliminaries and problem settings

**Episodic finite-horizon MDPs.** In this paper, we study episodic finite-horizon Markov decision processes with  $S$  states,  $A$  actions, and horizon length  $H$ . We shall employ  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, P = \{P_h\}_{h=1}^H, r = \{r_h\}_{h=1}^H)$  to represent such an MDP, where  $\mathcal{S} = [S]$  and  $\mathcal{A} = [A]$  represent the state space and the action space, respectively. For each step  $h \in [H]$ , we let  $P_h : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  represent the transition probability at this step, such that taking action  $a$  in state  $s$  at step  $h$  yields a transition to the next state drawn from the distribution  $P_h(\cdot | s, a)$ ; throughout the paper, we often employ the shorthand notation  $P_{h,s,a} := P_h(\cdot | s, a)$ . Another ingredient is the reward function specified by  $r_h : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  at step  $h$ ; namely, the agent will receive an immediate reward  $r_h(s, a)$  upon executing action  $a$  in state  $s$  at step  $h$ . It is assumed that the reward function is fully revealed upon completion of online data collection. Additionally, we assume throughout that each episode of the MDP starts from an initial state independently generated from some (unknown) initial state distribution  $\rho \in \Delta(\mathcal{S})$ .

A time-inhomogeneous Markovian policy is often denoted by  $\pi = \{\pi_h\}_{h=1}^H$  with  $\pi_h : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ , where  $\pi_h(\cdot | s)$  characterizes the (randomized) action selection probability of the agent in state  $s$  at step  $h$ . If  $\pi$  is a deterministic policy, then we often abuse the notation and let  $\pi_h(s)$  represent the action selected in state  $s$  at step  $h$ . We find it convenient to introduce the following notation:

$$\Pi := \text{the set of all deterministic policies.} \quad (1)$$

In this paper, we often need to cope with mixed deterministic policies (so that each realization of the policy is randomly drawn from a mixture of deterministic policies). A mixed deterministic policy  $\pi^{\text{mixed}}$  is often denoted by

$$\pi^{\text{mixed}} = \sum_{\pi \in \Pi} \mu(\pi) \pi = \mathbb{E}_{\pi \sim \mu}[\pi] \quad \text{for some } \mu \in \Delta(\Pi). \quad (2)$$

Moreover, for any policy  $\pi$ , we define its associated value function (resp. Q-function) as follows, representing the expected cumulative rewards conditioned on an initial state (resp. an initial state-action pair):

$$\begin{aligned} V_h^\pi(s) &:= \mathbb{E}_\pi \left[ \sum_{h'=h}^H r_{h'}(s, a) \middle| s_h = s \right], \quad \forall s \in \mathcal{S}; \\ Q_h^\pi(s, a) &:= \mathbb{E}_\pi \left[ \sum_{h'=h}^H r_{h'}(s, a) \middle| s_h = s, a_h = a \right], \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}. \end{aligned}$$

Here,  $\mathbb{E}_\pi[\cdot]$  indicates the expectation over the length- $H$  sample trajectory  $(s_1, a_1, s_2, a_2, \dots, s_H, a_H)$  when executing policy  $\pi$  in  $\mathcal{M}$ , where  $s_h$  (resp.  $a_h$ ) denotes the state (resp. action) at step  $h$  of this trajectory. When the initial state is drawn from  $\rho$ , we further augment the notation and denote

$$V_1^\pi(\rho) = \mathbb{E}_{s \sim \rho} [V_1^\pi(s)].$$

Importantly, there exists at least one deterministic policy, denoted by  $\pi^*$  throughout, that is able to maximize  $V_h^\pi(s)$  and  $Q_h^\pi(s, a)$  simultaneously for all  $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$ ; namely,

$$V_h^*(s) := V_h^{\pi^*}(s) = \max_{\pi} V_h^\pi(s), \quad Q_h^*(s, a) := Q_h^{\pi^*}(s, a) = \max_{\pi} Q_h^\pi(s, a), \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}.$$

The interested reader is referred to Bertsekas (2017) for more backgrounds about MDPs.

Moving beyond value functions and Q-functions, we would like to define, for each policy  $\pi$ , the associated state-action occupancy distribution  $d^\pi = [d_h^\pi]_{1 \leq h \leq H}$  such that

$$d_h^\pi(s, a) := \mathbb{P}(s_h = s, a_h = a | \pi), \quad \forall (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H];$$

in other words, this is the probability of the state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$  being visited by  $\pi$  at step  $h$ . We shall also overload  $d^\pi$  to represent the state occupancy distribution such that

$$d_h^\pi(s) := \sum_{a \in \mathcal{A}} d_h^\pi(s, a) = \mathbb{P}(s_h = s | \pi), \quad \forall (s, h) \in \mathcal{S} \times [H]. \quad (3)$$

Given that each episode always starts with a state drawn from  $\rho$ , it is easily seen that

$$d_1^\pi(s) = \rho(s) \quad (4)$$

holds for any policy  $\pi$  and any  $s \in \mathcal{S}$ .

**Sampling mechanism.** We consider a hybrid RL setting that assumes access to a historical dataset as well as the ability to further explore the environment via real-time sampling, as detailed below.

*Offline data.* Suppose that we have available a historical dataset (also called an offline dataset)

$$\mathcal{D}^{\text{off}} = \{\tau^{k,\text{off}}\}_{1 \leq k \leq K^{\text{off}}}, \quad (5)$$

containing  $K^{\text{off}}$  sample trajectories each of length  $H$ . Here, the  $k$ -th trajectory in  $\mathcal{D}^{\text{off}}$  is denoted by

$$\tau^{k,\text{off}} = (s_1^{k,\text{off}}, a_1^{k,\text{off}}, \dots, s_H^{k,\text{off}}, a_H^{k,\text{off}}), \quad (6)$$

where  $s_h^{k,\text{off}}$  and  $a_h^{k,\text{off}}$  indicate respectively the state and action at step  $h$  of this trajectory  $\tau^{k,\text{off}}$ . It is assumed that each trajectory  $\tau^{k,\text{off}}$  is drawn *independently* using policy  $\pi^{\text{off}}$ , which takes the form of a mixture of deterministic policies

$$\pi^{\text{off}} = \mathbb{E}_{\pi \sim \mu^{\text{off}}} [\pi] \quad \text{with } \mu^{\text{off}} \in \Delta(\Pi). \quad (7)$$

Note that the learner only has access to the data samples but not  $\pi^{\text{off}}$ . Throughout the paper, we use  $d^{\text{off}} = \{d_h^{\text{off}}\}_{1 \leq h \leq H}$  to represent the occupancy distribution of this offline dataset such that

$$d_h^{\text{off}}(s, a) := \mathbb{P}((s_h^{k,\text{off}}, a_h^{k,\text{off}}) = (s, a)), \quad \forall (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]. \quad (8)$$

*Online exploration.* In addition to the offline dataset, the learner is allowed to interact with the unknown environment and collect more data in real time, in the hope of compensating for the insufficiency of the pre-collected data at hand and fine-tuning the policy estimate. More specifically, the learner is able to sample  $K^{\text{on}}$  trajectories sequentially. In each sample trajectory,

- the initial state is generated independently from an (unknown) distribution  $\rho \in \Delta(\mathcal{S})$ ;
- the learner selects a policy to execute the MDP, obtaining a sample trajectory of length  $H$ .

The total number of sample trajectories is thus given by

$$K = K^{\text{off}} + K^{\text{on}}. \quad (9)$$

**Concentrability assumptions for the offline dataset.** To quantify the quality of the historical dataset, prior offline RL literature introduced the following single-policy concentrability coefficient based on certain density ratio of interest; see, e.g., Li et al. (2022); Rashidinejad et al. (2021).

**Definition 1** (Single-policy concentrability). *The single-policy concentrability coefficient  $C^*$  of the offline dataset  $\mathcal{D}^{\text{off}}$  is defined as*

$$C^* := \max_{(s,a,h) \in \mathcal{S} \times \mathcal{A} \times [H]} \frac{d_h^{\pi^*}(s, a)}{d_h^{\text{off}}(s, a)}. \quad (10)$$

In words,  $C^*$  employs the  $\ell_\infty$ -norm of the density ratio  $d^{\pi^*}/d^{\text{off}}$  to capture the shift of distributions between the occupancy distribution induced by the desired policy  $\pi^*$  and the data distribution at hand. The terminology ‘‘single-policy’’ underscores that Definition 1 only compares the offline data distribution against the one generated by a single policy  $\pi^*$ , which stands in stark contrast to other all-policy concentrability coefficients that are defined to account for all policies simultaneously.

One notable fact about Definition 1 is that: for  $C^*$  to be finite, the historical data distribution needs to cover all state-action-step tuples reachable by  $\pi^*$ . This requirement is, in general, inevitable if only the offline dataset is available; see the minimax lower bounds in Li et al. (2022); Rashidinejad et al. (2021) for more precise justifications. However, a requirement of this kind could be overly stringent for the hybrid setting considered herein, as the issue of incomplete coverage can potentially be overcome with the aid of online data collection. In light of this observation, we generalize Definition 1 to account for the trade-offs between distributional mismatch and partial coverage.

**Definition 2** (Single-policy partial concentrability). *For any  $\sigma \in [0, 1]$ , the single-policy partial concentrability coefficient  $C^*(\sigma)$  of the offline dataset  $\mathcal{D}^{\text{off}}$  is defined as*

$$C^*(\sigma) := \min \left\{ \max_{1 \leq h \leq H} \max_{(s,a) \in \mathcal{G}_h} \frac{d_h^{\pi^*}(s,a)}{d_h^{\text{off}}(s,a)} \mid \{\mathcal{G}_h\}_{1 \leq h \leq H} \subseteq \mathcal{G}(\sigma) \right\}, \quad (11)$$

where

$$\mathcal{G}(\sigma) := \left\{ \{\mathcal{G}_h\}_{1 \leq h \leq H} \subseteq \mathcal{S} \times \mathcal{A} \mid \frac{1}{H} \sum_{h=1}^H \sum_{(s,a) \notin \mathcal{G}_h} d_h^{\pi^*}(s,a) \leq \sigma \right\}. \quad (12)$$

In Definition 2, we allow a fraction of the state-action space reachable by  $\pi^*$  to be insufficiently covered (as reflected in the definition of  $\mathcal{G}(\sigma)$  measured by the state-action occupancy distribution) — hence the terminology “partial”. Intuitively,  $\mathcal{G}_h$  corresponds to a set of state-action pairs that undergo reasonable distribution shift (so that the corresponding density ratio does not rise above  $C^*(\sigma)$ ), whereas the total occupancy density of its complement subset  $\mathcal{G}_h^c$  induced by  $\pi^*$  is under control (i.e., no larger than  $\sigma$  when averaged across steps). As a self-evident fact,  $C^*(\sigma)$  is non-increasing in  $\sigma$ ; this means that as  $\sigma$  increases, we might incur a less severe distribution shift in a restricted part, at the price of less coverage. In this sense,  $C^*(\sigma)$  reflects certain tradeoffs between distribution shift and coverage. Clearly,  $C^*(\sigma)$  reduces to  $C^*$  in Definition 1 by taking  $\sigma = 0$ .

**Goal.** Given a historical dataset  $\mathcal{D}^{\text{off}}$  containing  $K^{\text{off}}$  sample trajectories, we would like to design an online exploration scheme, in conjunction with the accompanying policy learning algorithm, so as to achieve desirable policy learning (or policy fine-tuning) in a data-efficient manner. Ideally, we would expect a hybrid RL algorithm to harvest provable statistical benefits compared to both purely online RL and purely offline RL approaches.

### 3 Algorithm

In this section, we come up with a new algorithm to tackle the hybrid RL setting. Our algorithm design leverages recent ideas developed in offline RL and reward-agnostic online exploration to improve sample efficiency. The proposed algorithm consists of three stages to be described shortly; informally, the first two stages conduct reward-agnostic exploration to imitate and complement the offline dataset, whereas the third stage invokes a sample-optimal offline RL algorithm to compute a near-optimal policy.

In the sequel, we split the offline dataset  $\mathcal{D}^{\text{off}}$  into two halves:

$$\mathcal{D}^{\text{off},1} \quad \text{and} \quad \mathcal{D}^{\text{off},2}, \quad (13)$$

where  $\mathcal{D}^{\text{off},1}$  (resp.  $\mathcal{D}^{\text{off},2}$ ) consists of the first (resp. last)  $K^{\text{off}}/2$  independent trajectories from  $\mathcal{D}^{\text{off}}$ . As we shall also see momentarily, online exploration in the proposed algorithm — which collects  $K^{\text{on}}$  trajectories in total — can be divided into three parts, collecting  $K_{\text{prepare}}^{\text{on}}$ ,  $K_{\text{imitate}}^{\text{on}}$  and  $K_{\text{explore}}^{\text{on}}$  sample trajectories, respectively. Throughout this paper, for simplicity we choose

$$K_{\text{prepare}}^{\text{on}} = K_{\text{imitate}}^{\text{on}} = K_{\text{explore}}^{\text{on}} = K^{\text{on}}/3. \quad (14)$$

#### 3.1 A three-stage algorithm

We now elaborate on the three stages of the proposed algorithm.

**Stage 1: estimation of the occupancy distributions.** As a preparatory step for sample-efficient reward-agnostic exploration, we first attempt to estimate the occupancy distribution induced by any policy as well as the occupancy distribution  $d^{\text{off}}$  associated with the historical dataset, as described below.

- *Estimating  $d^\pi$  for any policy  $\pi$ .* In this step, we would like to sample the environment and collect a set of sample trajectories, in a way that allows for reasonable estimation of the occupancy distribution  $d^\pi$  induced by any policy  $\pi$ . For this purpose, we invoke the exploration strategy and the accompanying estimation scheme developed in Li et al. (2023). Working forward (i.e., from  $h = 1$  to  $H$ ), this approach collects, for each step  $h$ , a set of  $N$  sample trajectories in order to facilitate estimation of the occupancy distributions, which amounts to a total number of

$$NH =: K_{\text{prepare}}^{\text{on}} = K^{\text{on}}/3 \quad (15)$$

sample trajectories collected in this stage. See Algorithm 3 in Appendix A.1 for a precise description of this strategy. Noteworthy, while Algorithm 3 specifies how to estimate  $\widehat{d}^\pi$  for any policy  $\pi$ , we won't need to compute it explicitly unless this policy  $\pi$  is encountered during the subsequent steps of the algorithm; in other words,  $\widehat{d}^\pi$  should be viewed as a sort of "function handle" that will only be executed when called later.

- *Estimating  $d^{\text{off}}$  for the historical dataset  $\mathcal{D}^{\text{off}}$ .* In addition, we are in need of estimating the occupancy distribution  $d^{\text{off}}$ . Towards this end, we propose the following empirical estimate using the  $K^{\text{off}}/2$  sample trajectories from  $\mathcal{D}^{\text{off},1}$ :

$$\widehat{d}_h^{\text{off}}(s, a) = \frac{2N_h^{\text{off}}(s, a)}{K^{\text{off}}} \mathbb{1} \left( \frac{N_h^{\text{off}}(s, a)}{K^{\text{off}}} \geq c_{\text{off}} \left\{ \frac{\log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}} \right\} \right) \quad (16)$$

for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , where  $c_{\text{off}} > 0$  is some universal constant. Here,  $1 - \delta$  indicates the target success probability, and

$$N_h^{\text{off}}(s, a) = \sum_{k=1}^{K^{\text{off}}/2} \mathbb{1}(s_h^{k, \text{off}} = s, a_h^{k, \text{off}} = a), \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}. \quad (17)$$

In other words,  $\widehat{d}_h^{\text{off}}(s, a)$  is taken to be the empirical visitation frequency of  $(s, a)$  in  $\mathcal{D}^{\text{off},1}$  if  $(s, a)$  is adequately visited, and zero otherwise. The cutoff threshold  $c_{\text{off}} \left( \frac{\log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}} \right)$  will be made clear in our analysis.

**Stage 2: online exploration.** Armed with the above estimates of the occupancy distributions, we can readily proceed to compute the desired exploration policies and sample the environment. We seek to devise two exploration strategies, with one strategy selected to imitate the offline dataset, and the other one employed to explore the insufficiently visited territory. As a preliminary fact, if we have a dataset containing  $K$  independent trajectories — generated independently from a mixture of deterministic policies with occupancy distribution  $d^b$  — then it has been shown previously (see, e.g., Li et al. (2023, Section 3.3)) that the model-based offline approach is able to compute a policy  $\widehat{\pi}$  obeying

$$V^*(\rho) - V^{\widehat{\pi}}(\rho) \lesssim H \left[ \sum_h \sum_{s, a} \frac{d_h^{\pi^*}(s, a)}{1/H + K^{\text{on}} d_h^b(s, a)} \right]^{\frac{1}{2}}. \quad (18)$$

This upper bound in (18) provides a guideline regarding how to design a sample-efficient exploration scheme.

- *Imitating the offline dataset.* The offline dataset  $\mathcal{D}^{\text{off}}$  is most informative when it contains expert data, a scenario when the data distribution resembles the distribution induced by the optimal policy  $\pi^*$ . If this is the case, then it is desirable to find a policy similar to  $\pi^*$  in (7) (the mixed policy generating  $\mathcal{D}^{\text{off}}$ ) and employ it to collect new data, in order to retain and further strength the benefits of such offline data. To do so, we attempt to approximate  $d^{\pi^*}$  by  $\widehat{d}^{\text{off}}$  in (18) when attempting to minimize (18). In fact, we would like to compute a mixture of deterministic policies by (approximately) solving the following optimization problem:

$$\mu^{\text{imitate}} \approx \arg \min_{\mu \in \Delta(\Pi)} \sum_{h=1}^H \sum_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi' \sim \mu} [\widehat{d}_h^{\pi'}(s, a)]}, \quad (19)$$

which is clearly equivalent to

$$\mu^{\text{imitate}} \approx \arg \min_{\mu \in \Delta(\Pi)} \max_{\pi: \mathcal{S} \times [H] \rightarrow \Delta(\mathcal{A})} \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi' \sim \mu} [\widehat{d}_h^{\pi'}(s, a)]} \right]. \quad (20)$$

In order to solve this minimax problem (20) (note that its objective function is convex in  $\mu$ ), we resort to the Follow-The-Regularized-Leader (FTRL) strategy from the online learning literature (Shalev-Shwartz, 2012); more specifically, we perform the following updates iteratively for  $t = 1, \dots, T_{\max}$ :

$$\pi_h^{t+1}(\cdot|s) \propto \exp \left( \eta \sum_{k=1}^t \frac{\widehat{d}_h^{\text{off}}(s, \cdot)}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi' \sim \mu^k} [\widehat{d}_h^{\pi'}(s, \cdot)]} \right), \quad \forall s \in \mathcal{S}, \quad (21a)$$

$$\mu^{t+1} \approx \arg \min_{\mu \in \Delta(\Pi)} \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi' \sim \mu} [\widehat{d}_h^{\pi'}(s, a)]} \right], \quad (21b)$$

where  $\eta$  denotes the learning rate to be specified later. We shall discuss how to solve the optimization sub-problem (21b) in Section 3.2. The output of this step is a mixture of deterministic policies taking the following form:

$$\pi^{\text{imitate}} = \mathbb{E}_{\pi \sim \mu^{\text{imitate}}} [\pi] \quad \text{with} \quad \mu^{\text{imitate}} = \frac{1}{T_{\max}} \sum_{t=1}^{T_{\max}} \mu^t. \quad (22)$$

- *Exploring the unknown environment.* In addition to mimicking the behavior of the historical dataset, we shall also attempt to explore the environment in a way that complements pre-collected data. Towards this end, it suffices to invoke the reward-agnostic online exploration scheme proposed in Li et al. (2023), whose precise description will be provided in Algorithm 5 in Appendix A.2 to make the paper self-contained. The resulting policy mixture is denoted by

$$\pi^{\text{explore}} = \mathbb{E}_{\pi \sim \mu^{\text{explore}}} [\pi], \quad (23)$$

with  $\mu^{\text{explore}} \in \Delta(\Pi)$  representing the associated weight vector.

With the above two exploration policies (22) and (23) in place, we execute the MDP to obtain sample trajectories as follows:

- 1) Execute the MDP  $K_{\text{imitate}}^{\text{on}}$  times using policy  $\pi^{\text{imitate}}$  to obtain a dataset containing  $K_{\text{imitate}}^{\text{on}} = K^{\text{on}}/3$  independent sample trajectories, denoted by  $\mathcal{D}_{\text{imitate}}^{\text{on}}$ ;
- 2) Execute the MDP  $K_{\text{explore}}^{\text{on}}$  times using policy  $\pi^{\text{explore}}$  to obtain a dataset containing  $K_{\text{explore}}^{\text{on}} = K^{\text{on}}/3$  independent sample trajectories, denoted by  $\mathcal{D}_{\text{explore}}^{\text{on}}$ .

**Stage 3: policy learning via offline RL.** Once the above online exploration process is completed, we are positioned to compute a near-optimal policy on the basis of the data in hand. More precisely,

- Let us look at the following dataset

$$\mathcal{D} = \mathcal{D}^{\text{off},2} \cup \mathcal{D}_{\text{imitate}}^{\text{on}} \cup \mathcal{D}_{\text{explore}}^{\text{on}}. \quad (24)$$

In light of the complicated statistical dependency between  $\mathcal{D}^{\text{off},1}$  and  $\mathcal{D}_{\text{imitate}}^{\text{on}} \cup \mathcal{D}_{\text{explore}}^{\text{on}}$ , we only include the second half  $\mathcal{D}^{\text{off},2}$  of the offline dataset  $\mathcal{D}^{\text{off}}$ , so as to exploit the fact that  $\mathcal{D}^{\text{off},2}$  is statistically independent from  $\mathcal{D}_{\text{imitate}}^{\text{on}} \cup \mathcal{D}_{\text{explore}}^{\text{on}}$ .

- We invoke the pessimistic model-based offline RL algorithm proposed in Li et al. (2022) to compute the final policy estimate  $\widehat{\pi}$ ; see Algorithm 6 in Appendix A.3 for more details.



---

**Algorithm 1:** The proposed hybrid RL algorithm.

---

- 1 **Input:** offline dataset  $\mathcal{D}^{\text{off}}$  (containing  $K^{\text{off}}$  trajectories), parameters  $N, K^{\text{on}}, T_{\max}$ , learning rate  $\eta$ .
- 2 **Initialize:**  $\pi_h^1(a | s) = 1/A$  for any  $(s, a, h)$ ;  $K = K^{\text{off}} + K^{\text{on}}$ ; split  $\mathcal{D}^{\text{off}}$  into two halves  $\mathcal{D}^{\text{off},1}$  and  $\mathcal{D}^{\text{off},2}$ .  
 /\* Estimation of occupancy distributions for any policy  $\pi$ . \*/
- 3 Call Algorithm 3, which allows one to specify  $\widehat{d}_h^\pi(s, a)$  for any deterministic policy  $\pi$  and any  $(s, a, h)$ .  
 /\* Estimation of occupancy distributions of the historical data. \*/
- 4 Use the dataset  $\mathcal{D}^{\text{off},1}$  to compute

$$\widehat{d}_h^{\text{off}}(s, a) = \frac{2N_h^{\text{off}}(s, a)}{K^{\text{off}}} \mathbb{1} \left( \frac{N_h^{\text{off}}(s, a)}{K^{\text{off}}} \geq c_{\text{off}} \left\{ \frac{\log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}} \right\} \right)$$

for any  $(s, a, h)$ , where  $N_h^{\text{off}}(s, a) = \sum_{k=1}^{K^{\text{off}}/2} \mathbb{1}(s_h^k = s, a_h^k = a)$  and  $c_{\text{off}} > 0$  is some absolute constant.

- /\* Compute a general sample-efficient online exploration scheme. \*/
- 5 Call Algorithm 5 with estimators  $\widehat{d}^\pi$  to compute policy  $\pi^{\text{explore}}$  and the associated weight  $\mu^{\text{explore}}$ .  
 /\* Compute an online exploration scheme tailored to the offline dataset. \*/
- 6 **for**  $t = 1, \dots, T_{\max}$  **do**
- 7     Compute  $\mu^t$  using Algorithm 2.
- 8     Update  $\pi_h^{t+1}(a | s)$  for all  $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$  such that:

$$\pi_h^{t+1}(a | s) = \frac{\exp \left( \eta \sum_{k=1}^t \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi' \sim \mu^k} [\widehat{d}_h^{\pi'}(s, a)]} \right)}{\sum_{a' \in \mathcal{A}} \exp \left( \eta \sum_{k=1}^t \frac{\widehat{d}_h^{\text{off}}(s, a')}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi' \sim \mu^k} [\widehat{d}_h^{\pi'}(s, a')] } \right)},$$

- 9 Set  $\mu^{\text{imitate}} = \frac{1}{T_{\max}} \sum_{t=1}^{T_{\max}} \mu^t$  and  $\pi^{\text{imitate}} = \mathbb{E}_{\pi \sim \mu^{\text{imitate}}}[\pi]$ .  
 /\* Sampling using the above two exploration policies. \*/
  - 10 Collect  $K_{\text{imitate}}^{\text{on}}$  (resp.  $K_{\text{explore}}^{\text{on}}$ ) sample trajectories using  $\pi^{\text{imitate}}$  (resp.  $\pi^{\text{explore}}$ ) to form a dataset  $\mathcal{D}_{\text{imitate}}^{\text{on}}$  (resp.  $\mathcal{D}_{\text{explore}}^{\text{on}}$ ).  
 /\* Run the model-based offline RL algorithm. \*/
  - 11 Apply Algorithm 6 to the dataset  $\mathcal{D} = \mathcal{D}^{\text{off},2} \cup \mathcal{D}_{\text{imitate}}^{\text{on}} \cup \mathcal{D}_{\text{explore}}^{\text{on}}$  to compute a policy  $\widehat{\pi}$ .
  - 12 **Output:** policy  $\widehat{\pi}$ .
- 

### 3.2 Subroutine for solving the subproblem (21b)

While (21b) is a convex optimization subproblem, it involves optimization over a parameter space with exponentially large dimensions. In order to solve it in a computationally feasible manner, we propose a tailored subroutine based on the Frank-Wolfe algorithm (Bubeck, 2015).

Before proceeding, recall that when specifying  $\widehat{d}^\pi$  in Algorithm 3, we draw  $N$  independent trajectories  $\{s_1^{n,h}, a_1^{n,h}, \dots, s_{h+1}^{n,h}\}_{1 \leq n \leq N}$  and compute an empirical estimate  $\widehat{P}_h$  of the probability transition kernel at step  $h$  such that

$$\widehat{P}_h(s' | s, a) = \frac{\mathbb{1}(N_h(s, a) > \xi)}{\max \{N_h(s, a), 1\}} \sum_{n=1}^N \mathbb{1}(s_h^{n,h} = s, a_h^{n,h} = a, s_{h+1}^{n,h} = s'), \quad \forall (s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}, \quad (25)$$

where  $N_h(s, a) = \sum_{n=1}^N \mathbb{1}\{s_h^{n,h} = s, a_h^{n,h} = a\}$ .

**The proposed Frank-Wolfe-type algorithm.** With this set of notation in place and with an initial guess taken to be the indicator function  $\mu^{(1)} = \mathbb{1}_{\pi_{\text{init}}}$  for an arbitrary policy  $\pi_{\text{init}} \in \Pi$ , the  $k$ -th iteration of our iterative procedure for solving (21b) can be described as follows.

- *Computing a search direction.* The search direction of the Frank-Wolfe algorithm is typically taken to be a feasible direction that maximizes its correlation with the gradient of the objective function (Bubeck, 2015). When specialized to the current sub-problem (21b), the search direction can be taken to be the Dirac measure  $\delta_{\pi^{(k)}}$ , where

$$\pi^{(k)} = \arg \max_{\pi \in \Pi} f(\pi, \mu^{(k)}) := \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^\pi(s, a) \widehat{d}_h^{\text{off}}(s, a)}{\left( \frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi' \sim \mu^{(k)}} [\widehat{d}_h^{\pi'}(s, a)] \right)^2} \right]. \quad (26)$$

As it turns out, this optimization problem (26) can be efficiently solved by applying dynamic programming (Bertsekas, 2017) to an augmented MDP  $\mathcal{M}^{\text{off}}$  constructed as follows.

- Introduce an augmented finite-horizon MDP  $\mathcal{M}^{\text{off}} = (\mathcal{S} \cup \{s_{\text{aug}}\}, \mathcal{A}, H, \widehat{P}^{\text{aug}}, r^{\text{off}})$ , where  $s_{\text{aug}}$  is an augmented state. We choose the reward function to be

$$r_h^{\text{off}}(s, a) = \begin{cases} \frac{\pi_h^{t+1}(a|s) \widehat{d}_h^{\text{off}}(s, a)}{\left( \frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi' \sim \mu^{(k)}} [\widehat{d}_h^{\pi'}(s, a)] \right)^2}, & \text{if } (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H], \\ 0, & \text{if } (s, a, h) \in \{s_{\text{aug}}\} \times \mathcal{A} \times [H], \end{cases} \quad (27)$$

and the probability transition kernel as

$$\widehat{P}_h^{\text{aug}}(s' | s, a) = \begin{cases} \widehat{P}_h(s' | s, a), & \text{if } s' \in \mathcal{S} \\ 1 - \sum_{s' \in \mathcal{S}} \widehat{P}_h(s' | s, a), & \text{if } s' = s_{\text{aug}} \end{cases} \quad \text{for all } (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H], \quad (28a)$$

$$\widehat{P}_h^{\text{aug}}(s' | s_{\text{aug}}, a) = \mathbb{1}(s' = s_{\text{aug}}) \quad \text{for all } (a, h) \in \mathcal{A} \times [H]. \quad (28b)$$

- *Frank-Wolfe updates.* We then update the iterate  $\mu^{(k+1)}$  as a convex combination of the current iterate and the direction found in the previous step:

$$\mu^{(k+1)} = (1 - \alpha) \mu^{(k)} + \alpha \mathbb{1}_{\pi^{(k)}}, \quad (29)$$

where the stepsize is chosen to be

$$\alpha = \frac{S}{(K^{\text{on}}H)^3}. \quad (30)$$

**Stopping rule.** It is also necessary to specify the stopping rule of the above iterative procedure. Throughout this paper, the above subroutine will terminate as long as

$$\sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi' \sim \mu^{(k)}} [\widehat{d}_h^{\pi'}(s, a)]} \right] \leq 108SH, \quad (31)$$

with the final output taken to be  $\mu^{t+1} = \mu^{(k)}$ . We shall justify the feasibility of this stopping rule (namely, the fact that this stopping criterion can be met by some mixed policy) in Section B.

**Iteration complexity.** Encouragingly, the above subroutine in conjunction with the stopping rule (31) leads to an iteration complexity no larger than

$$(\text{iteration complexity}) \quad O\left(\frac{(K^{\text{on}}H)^4}{S^2}\right) \quad (32)$$

The proof of this claim is postponed to Section B.

## 4 Main results

As it turns out, the proposed procedure in Algorithm 1 is capable of achieving provable sample efficiency, as demonstrated in the following theorem. Here and below, we recall the notation

$$K = K^{\text{off}} + K^{\text{on}}. \quad (36)$$

---

**Algorithm 2:** Subroutine for solving the sub-problem (21b).

---

1 **Initialize:**  $\mu^{(1)} = \mathbb{1}_{\pi_{\text{init}}}$  for an arbitrary policy  $\pi_{\text{init}} \in \Pi$ .

2 **for**  $k = 1, 2, \dots$  **do**

3     Exit for-loop if the following condition is met: // **stopping criterion**

4

$$\sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi' \sim \mu^{(k)}} [\widehat{d}_h^{\pi'}(s, a)]} \right] \leq 108SH. \quad (33)$$

/\* Find the search direction \*/

5     Compute the optimal deterministic policy  $\pi^{(k), \text{aug}}$  of the MDP

$\mathcal{M}_{\text{off}} = (\mathcal{S} \cup \{s_{\text{aug}}\}, \mathcal{A}, H, \widehat{P}^{\text{aug}}, r_{\text{off}})$ , where  $s_{\text{aug}}$  is an augmented state,

$$r_h^{\text{off}}(s, a) = \begin{cases} \frac{\pi_h^{t+1}(a|s) \widehat{d}_h^{\text{off}}(s, a)}{\left(\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi' \sim \mu^{(k)}} [\widehat{d}_h^{\pi'}(s, a)]\right)^2}, & \text{if } (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H], \\ 0, & \text{if } (s, a, h) \in \{s_{\text{aug}}\} \times \mathcal{A} \times [H], \end{cases} \quad (34)$$

and the augmented probability transition kernel is given by

$$\widehat{P}_h^{\text{aug}}(s' | s, a) = \begin{cases} \widehat{P}_h(s' | s, a), & \text{if } s' \in \mathcal{S} \\ 1 - \sum_{s' \in \mathcal{S}} \widehat{P}_h(s' | s, a), & \text{if } s' = s_{\text{aug}} \end{cases} \quad \text{for all } (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]; \quad (35a)$$

$$\widehat{P}_h^{\text{aug}}(s' | s_{\text{aug}}, a) = \mathbb{1}(s' = s_{\text{aug}}) \quad \text{for all } (a, h) \in \mathcal{A} \times [H]. \quad (35b)$$

Let  $\pi^{(k)}$  be the corresponding optimal deterministic policy of  $\pi^{(k), \text{aug}}$  in the original state space.

6     Update // **Frank-Wolfe update**

7

$$\mu^{(k+1)} = (1 - \alpha)\mu^{(k)} + \alpha \mathbb{1}_{\pi^{(k)}}, \quad \text{where } \alpha = \frac{S}{(K^{\text{on}}H)^3}.$$

8 **Output:** the policy mixture  $\mu^{t+1} = \mu^{(k)}$ .

---

**Theorem 1.** Consider  $\delta \in (0, 1)$  and  $\varepsilon \in (0, H]$ . Choose the algorithmic parameters such that

$$\eta = \sqrt{\frac{\log A}{2T_{\max}(K^{\text{on}}H)^2}} \quad \text{and} \quad T_{\max} \geq 2(K^{\text{on}}H)^2 \log A.$$

Suppose that

$$K^{\text{on}} + K^{\text{off}} \geq c_1 \frac{H^3 SC^*(\sigma)}{\varepsilon^2} \log^2 \frac{K}{\delta} \quad (37a)$$

$$K^{\text{on}} \geq c_1 \frac{H^3 SA \min\{H\sigma, 1\}}{\varepsilon^2} \log \frac{K}{\delta} \quad (37b)$$

for some large enough constant  $c_1 > 0$ . Then with probability at least  $1 - \delta$ , the policy  $\widehat{\pi}$  returned by Algorithm 1 satisfies

$$V_1^*(\rho) - V^{\widehat{\pi}}(\rho) \leq \varepsilon,$$

provided that  $K^{\text{on}}$  and  $K^{\text{off}}$  both exceed some polynomial  $\text{poly}(H, S, A, C^*(\sigma), \log \frac{K}{\delta})$  (independent of  $\varepsilon$ ).

In a nutshell, Theorem 1 uncovers that our algorithm yields  $\varepsilon$ -accuracy as long as

$$K^{\text{on}} + K^{\text{off}} \gtrsim \frac{H^3 SC^*(\sigma)}{\varepsilon^2} \log^2 \frac{K}{\delta}, \quad (38a)$$

$$K^{\text{on}} \gtrsim \frac{H^3 SA \min\{H\sigma, 1\}}{\varepsilon^2} \log \frac{K}{\delta}, \quad (38b)$$

ignoring lower-order terms. Several implications of this result are as follows.

**Sample complexity benefits compared with pure online or pure offline RL.** In order to make apparent its advantage compared with both pure offline RL and pure online RL, we make the following comparisons:

- *Sample complexity with balanced online and offline data.* For the ease of presentation, let us look at a simple case where  $K^{\text{off}} = K^{\text{on}} = K/2$ . The the sample complexity bound (38) in this case simplifies to

$$\tilde{O} \left( \min_{\sigma \in [0,1]} \left\{ \frac{H^3 SA \min\{H\sigma, 1\}}{\varepsilon^2} + \frac{H^3 SC^*(\sigma)}{\varepsilon^2} \right\} \right) =: \tilde{O} \left( \min_{\sigma \in [0,1]} f_{\text{mixed}}(\sigma) \right). \quad (39)$$

- *Comparisons with pure online RL.* We now look at pure online RL, corresponding to the case where  $K = K^{\text{on}}$  (so that all sample episodes are collected via online exploration). In this case, the minimax-optimal sample complexity for computing an  $\varepsilon$ -optimal policy is known to be (Azar et al., 2017; Li et al., 2023)

$$\tilde{O} \left( \frac{H^3 SA}{\varepsilon^2} \right) = \tilde{O}(f_{\text{mixed}}(1)) \quad (40)$$

assuming that  $\varepsilon$  is sufficiently small, which is clearly worse than (39). For instance, if there exists some very small  $\sigma \ll 1/H$  obeying  $C^*(\sigma) \lesssim 1$ , then the ratio of (39) to (40) is at most

$$H\sigma + \frac{1}{A} \ll 1, \quad (41)$$

thus resulting in substantial sample size savings.

- *Comparisons with pure offline RL.* In contrast, in the pure offline case where  $K = K^{\text{off}}$ , the minimax sample complexity is known to be (Li et al., 2022)

$$\tilde{O} \left( \frac{H^3 SC^*(0)}{\varepsilon^2} \right) = \tilde{O}(f_{\text{mixed}}(0)) \quad (42)$$

for any target accuracy level  $\varepsilon$ , which is apparently larger than (39) in general. In particular, recognizing that  $C^*(0) = \infty$  in the presence of incomplete coverage of the state-action space reachable by  $\pi^*$ , we might harvest enormous sample size benefits (by exploiting the ability of online RL to visit the previously uncovered state-action-step tuples).

**Comparison with Wagenmaker and Pacchiano (2022).** It is worth noting that Wagenmaker and Pacchiano (2022) also considered policy fine-tuning and proposed a method called FTPedel to tackle linear MDPs. The results therein, however, were mainly instance-dependent, thus making it difficult to compare in general. That being said, we would like to clarify two points:

- Wagenmaker and Pacchiano (2022) imposed all-policy concentrability assumptions, requiring the combined dataset (i.e., the offline and online data altogether) to cover certain feature vectors for all linear softmax policies (see Wagenmaker and Pacchiano (2022, Definition 4.1)). In contrast, our results only assume single-policy (partial) concentrability, which is much weaker than the all-policy counterpart.
- When specializing Wagenmaker and Pacchiano (2022, Corollary 1) to the tabular cases, the sample complexity therein becomes  $\tilde{O}(H^7 S^2 A^2 / \varepsilon^2)$ , which could be much larger than our result.

**Miscellaneous properties of the proposed algorithm.** In addition to the sample complexity advantages, the proposed hybrid RL enjoys several attributes that could be practically appealing.

- *Adaptivity to unknown optimal  $\sigma$ .* While we have introduced the parameter  $\sigma$  to capture incomplete coverage, our algorithm does not rely on any knowledge of  $\sigma$ . Take the balanced case described around (39) for instance: our algorithm automatically identifies the optimal  $\sigma$  that minimizes the function  $f_{\text{mixed}}(\sigma)$  over all  $\sigma \in [0, 1]$ . In other words, Algorithm 1 is able to automatically identify the optimal trade-offs between distribution mismatch and inadequate coverage.
- *Reward-agnostic data collection.* It is noteworthy that the online exploration procedure employed in Algorithm 1 does not require any prior information about the reward function. In other words, it is mainly designed to improve coverage of the state-action space, a property independent from the reward function. In truth, the reward function is only queried at the last step to output the learned policy. This enables us to perform hybrid RL in a reward-agnostic manner, which is particularly intriguing in practice, as there is no shortage of scenarios where the reward functions might be engineered subsequently to meet different objectives.
- *Strengthening behavior cloning.* Another notable feature is that our algorithm does not rely on prior knowledge about the policies generating the offline dataset  $\mathcal{D}^{\text{off}}$ ; in fact, it is capable of finding a mixed exploration policy  $\pi^{\text{imitate}}$  that inherits the advantages of the unknown behavior policy  $\pi^{\text{off}}$ . This could be of particular interest for behavior cloning, where the offline dataset  $\mathcal{D}^{\text{off}}$  is generated by an expert policy, with  $C^* = C^*(0) \approx 1$ , i.e. the expert policy covers the optimal one. In this situation, the supplement of online data collection improves behavior cloning by lowering the statistical error from  $\sqrt{\frac{H^3 SC^*}{K_{\text{off}}}}$  to  $\sqrt{\frac{H^3 SC^*}{K_{\text{off}} + K_{\text{on}}}}$ , together with an executable learned policy  $\pi^{\text{imitate}}$ .

## 5 Related works

In this section, we briefly discuss a small set of additional prior works related to the current paper.

**(Reward-aware) online RL.** In online RL, an agent seeks to find a near-optimal policy by sequentially and adaptively interacting with the unknown environment, without having access to any additional offline dataset. The extensive studies of online RL gravitate around how to optimally trade off exploration against exploitation, for which the principle of optimism in the face of uncertainty plays a crucial role (Auer and Ortner, 2006; Azar et al., 2017; Bai et al., 2019; Dann et al., 2017; Dong et al., 2019; Jaksch et al., 2010; Jin et al., 2018; Li et al., 2021b; Ménard et al., 2021b; Zhang et al., 2020b). Information-theoretic lower bounds have been established by Domingues et al. (2021); Jin et al. (2018), matching existing sample complexity upper bounds when the target accuracy level  $\varepsilon$  is sufficiently small. A further strand of works extended these studies to the case with function approximation, including both linear function approximation (Jin et al., 2020b; Li et al., 2021a; Zanette et al., 2020; Zhou et al., 2021a) and other more general families of function approximation (Du et al., 2021; Foster et al., 2021; Jin et al., 2021a).

**Offline RL.** In contrast to online RL, offline RL assumes access to a pre-collected offline dataset and precludes active interactions with the environment. Given the absence of further data collection, the sample complexity of pure offline RL depends heavily upon the quality of the offline dataset at hand, which has often been characterized via some sorts of concentrability coefficients in prior works (Rashidinejad et al., 2021; Zhan et al., 2022). Earlier works (Chen and Jiang, 2019; Munos and Szepesvári, 2008) typically operated under the assumption of all-policy concentrability — namely, the assumption that the dataset covers the visited state-action pairs of all possible policies — thus imposing a stringent requirement for the offline dataset to be highly explorative. To circumvent this stringent assumption, Jin et al. (2021b); Kumar et al. (2020); Li et al. (2022); Liu et al. (2020); Rashidinejad et al. (2021); Shi et al. (2022); Uehara and Sun (2021); Yin et al. (2021) incorporated the pessimism principle amid uncertainty into the algorithm designs and, as a result, required only single-policy concentrability (so that the dataset only needs to cover the part of the state-action space reachable by the optimal policy). With regards to the basic tabular case, Li et al. (2022) proved that the pessimistic model-based offline algorithm is capable of achieving minimax-optimal sample

complexity for the full  $\varepsilon$ -range, accommodating both the episodic finite-horizon case and the discounted infinite-horizon analog. Moving beyond single-agent tabular settings, a recent line of works investigated offline RL in the presence of general function approximation (Jin et al., 2020c; Xie et al., 2021a; Zhan et al., 2022), environment shift (Shi and Chi, 2022; Zhou et al., 2021b), and in the context of zero-sum Markov games (Cui and Du, 2022; Yan et al., 2022).

**Hybrid RL.** While there were a number of empirical works (Hester et al., 2018; Kalashnikov et al., 2018; Nair et al., 2020, 2018; Rajeswaran et al., 2017; Vecerik et al., 2017) suggesting the performance gain of combining online RL with offline datasets (compared to pure online or offline learning), rigorous theoretical evidence remained highly limited. Ross and Bagnell (2012); Xie et al. (2021b) attempted to develop theoretical understanding by looking at one special hybrid scenario, where the agent can perform either of the following in each episode: (i) collecting a new online episode; and (ii) executing a prescribed and fixed reference policy to generate a sample episode. In this setting, Xie et al. (2021b) showed that in the minimax sense, combining online learning with samples generated by such a reference policy is not advantageous in comparison with pure online or offline RL. Note that our results do not contradict with the lower bound in Xie et al. (2021b), given that we exploit “partial” single-policy concentrability that implies additional structure except for the worst case. Akin to the current paper, Song et al. (2022); Wagenmaker and Pacchiano (2022) studied online RL with additional access to an offline dataset. Nevertheless, Song et al. (2022) mainly focused on the issue of computational efficiency, and the algorithm proposed therein does not come with improved sample complexity. In contrast, Wagenmaker and Pacchiano (2022) focused attention on statistical efficiency, although the sample complexity derived therein is highly suboptimal when specialized to the tabular setting.

**Reward-free and task-agnostic exploration.** Reward-free and task-agnostic exploration, which refer to the scenario where the agent first collects online sample trajectories without guidance of any information about the reward function(s), has garnered much recent attention (Brafman and Tenenholz, 2002; Huang et al., 2022; Jin et al., 2020a; Zhang et al., 2020a, 2021b). Focusing on the tabular case, the earlier work Jin et al. (2020a) put forward a reward-free exploration scheme that achieves minimax optimality in terms of the dependency on  $S$ ,  $A$  and  $1/\varepsilon$ , with the horizon dependency further improved by subsequent works (Kaufmann et al., 2021; Li et al., 2023; Ménard et al., 2021a). In particular, the exploration scheme proposed in Li et al. (2023) was shown to achieve minimax-optimal sample complexity when there exist a polynomial number of pre-determined but unseen reward functions of interest, which inspires the algorithm design of the present paper. Moreover, reward-free RL has been extended to account for function approximation, including both linear (Agarwal et al., 2020; Qiao and Wang, 2022; Wagenmaker et al., 2022; Wang et al., 2020; Zhang et al., 2021a) and nonlinear function classes (Chen et al., 2022).

## 6 Analysis of Theorem 1

In this section, we present the proof for our main result in Theorem 1. Throughout the proof, we let  $\{\mathcal{G}_h\}_{1 \leq h \leq H}$  denote a sequence of subsets obeying

$$\max_{1 \leq h \leq H} \max_{(s,a) \in \mathcal{G}_h} \frac{d_h^{\pi^*}(s,a)}{d_h^{\text{off}}(s,a)} = C^*(\sigma) \quad \text{and} \quad \frac{1}{H} \sum_{h=1}^H \sum_{(s,a) \notin \mathcal{G}_h} d_h^{\pi^*}(s,a) \leq \sigma, \quad (43)$$

as motivated by Definition 2. As it turns out, if  $K^{\text{on}} \geq c_1 \frac{H^3 SA}{\varepsilon^2} \log \frac{K}{\delta}$  for some large enough constant  $c_1 > 0$ , then the claimed result in Theorem 1 follows immediately from the main theory in Li et al. (2023) developed for pure online exploration. As a result, it suffices to prove the theorem by replacing Condition (37b) with

$$K^{\text{on}} \geq c_1 \frac{H^4 SA \sigma}{\varepsilon^2} \log \frac{K}{\delta} \quad (44)$$

throughout this section.

On a high level, our proof comprises the following three steps:

- Establish the proximity of  $\widehat{d}^{\text{off}}$  (resp.  $\widehat{d}^\pi$ ) and  $d^{\text{off}}$  (resp.  $d^\pi$ ).
- Show that the mixed policy  $\pi^{\text{imitate}}$  is able to mimic and strengthen the offline dataset  $\mathcal{D}^{\text{off}}$ , while the mixed policy  $\pi^{\text{explore}}$  is capable of exploring the part of the state-action space that has not been adequately visited by  $\mathcal{D}^{\text{off}}$ .
- Derive the sub-optimality of the policy returned by the offline RL algorithm (i.e., Algorithm 6) when applied to the hybrid dataset  $\mathcal{D} = \mathcal{D}^{\text{off},2} \cup \mathcal{D}_{\text{imitate}}^{\text{on}} \cup \mathcal{D}_{\text{explore}}^{\text{on}}$ .

In the sequel, we shall elaborate on these three steps.

## 6.1 Step 1: establishing the proximity of $\widehat{d}^\pi$ (resp. $\widehat{d}^{\text{off}}$ ) and $d^\pi$ (resp. $d^{\text{off}}$ )

To begin with, the goodness of the occupancy distribution estimators  $\widehat{d}^\pi$  (cf. Algorithm 3) has been analyzed in Li et al. (2023, Lemma 4), which come with the following performance guarantees.

**Lemma 1** (Li et al. (2023)). *Recall that  $\xi = c_\xi H^3 S^3 A^3 \log \frac{HSA}{\delta}$  for some large enough constant  $c_\xi > 0$ . With probability at least  $1 - \delta$ , the estimated occupancy distributions specified in Algorithm 3 satisfy*

$$\frac{1}{2} \widehat{d}_h^\pi(s, a) - \frac{\xi}{4N} \leq d_h^\pi(s, a) \leq 2\widehat{d}_h^\pi(s, a) + 2e_h^\pi(s, a) + \frac{\xi}{4N} \quad (45)$$

simultaneously for all  $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$  and all deterministic policy  $\pi \in \Pi$ , provided that

$$K^{\text{on}} \geq C_N H^{18} S^{14} A^{14} \log^2 \frac{HSA}{\delta} \quad (46)$$

for some large enough constant  $C_N > 0$ . Here,  $\{e_h^\pi(s, a)\}$  is some non-negative sequence satisfying

$$\sum_{s,a} e_h^\pi(s, a) \leq \frac{2SA}{K^{\text{on}}} + \frac{13SAH\xi}{N} \quad \text{for all } h \in [H] \text{ and all deterministic Markov policy } \pi. \quad (47)$$

We now turn to the estimator  $\widehat{d}^{\text{off}}$  (cf. (16)) for the occupancy distribution of the offline dataset, for which we begin with the following lemma concerning the proximity of  $d_h^{\text{off}}$  and  $\widehat{d}_h^{\text{off}}$ . The proof of this lemma is deferred to Section C.1.

**Lemma 2.** *Suppose that  $c_{\text{off}} \geq 48$ . With probability at least  $1 - \delta/3$ , one has*

$$\frac{1}{3} \widehat{d}_h^{\text{off}}(s, a) \leq d_h^{\text{off}}(s, a) \leq \widehat{d}_h^{\text{off}}(s, a) + 5c_{\text{off}} \left\{ \frac{\log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}} \right\} \quad (48)$$

simultaneously for all  $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ .

This lemma implies that: when  $d_h^{\text{off}}(s, a) \lesssim \frac{\log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}}$ , the estimator  $\widehat{d}_h^{\text{off}}(s, a)$  might be unable to track  $d_h^{\text{off}}(s, a)$  in a faithful manner. This motivates us to single out the following two subsets of state-action pairs for which  $\widehat{d}_h^{\text{off}}(s, a)$  might become problematic at step  $h$ :

- the set  $\mathcal{G}_h^c$  (see (12) for the definition of  $\mathcal{G}_h$ ), which corresponds to the set of optimal state-action pairs that even the true data distribution  $d_h^{\text{off}}$  cannot cover adequately;
- another set  $\mathcal{T}_h^{\text{small}}$  defined as

$$\mathcal{T}_h^{\text{small}} := \left\{ (s, a) : d_h^{\text{off}}(s, a) \leq 10c_{\text{off}} \left( \frac{\log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}} \right) \right\}, \quad (49)$$

comprising those state-action pairs for which  $\widehat{d}_h^{\text{off}}(s, a)$  might not be a faithful estimator of  $d_h^{\text{off}}(s, a)$ .

In what follow, we shall adopt the notation:

$$\mathcal{T}_h := \mathcal{G}_h^c \cup \mathcal{T}_h^{\text{small}}. \quad (50)$$

It is straightforward to demonstrate that:

- For any  $(s, a) \notin \mathcal{T}_h^{\text{small}}$ , it is seen from Lemma 2 that

$$d_h^{\text{off}}(s, a) \leq \widehat{d}_h^{\text{off}}(s, a) + \frac{1}{2}d_h^{\text{off}}(s, a) \iff d_h^{\text{off}}(s, a) \leq 2\widehat{d}_h^{\text{off}}(s, a). \quad (51)$$

- For any  $(s, a) \in \mathcal{G}_h$ , Condition (43) tells us that

$$d_h^{\pi^*}(s, a) \leq C^*(\sigma)d_h^{\text{off}}(s, a). \quad (52)$$

As a consequence, any  $(s, a) \notin \mathcal{T}_h$  necessarily obeys

$$d_h^{\pi^*}(s, a) \leq C^*(\sigma)d_h^{\text{off}}(s, a) \leq 2C^*(\sigma)\widehat{d}_h^{\text{off}}(s, a). \quad (53)$$

Another useful observation that we can readily make is as follows:

$$\begin{aligned} \sum_{h=1}^H \sum_{(s,a) \in \mathcal{T}_h} d_h^{\pi^*}(s, a) &\leq \sum_{h=1}^H \sum_{(s,a) \notin \mathcal{G}_h} d_h^{\pi^*}(s, a) + \sum_{h=1}^H \sum_{(s,a) \in \mathcal{G}_h \cap \mathcal{T}_h^{\text{small}}} d_h^{\pi^*}(s, a) \\ &\leq H\sigma + \sum_{h=1}^H \sum_{(s,a) \in \mathcal{G}_h \cap \mathcal{T}_h^{\text{small}}} d_h^{\pi^*}(s, a) \\ &\leq H\sigma + C^*(\sigma) \sum_{h=1}^H \sum_{(s,a) \in \mathcal{T}_h^{\text{small}}} d_h^{\text{off}}(s, a) \mathbb{1}(a = \pi^*(s)) \\ &\leq H\sigma + C^*(\sigma)HS \cdot 10c_{\text{off}} \left( \frac{\log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}} \right) \\ &\leq H\sigma + 10c_{\text{off}} \left( \frac{C^*(\sigma)HS \log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{4C^*(\sigma)H^6 S^5 A^4 \log \frac{HSA}{\delta}}{K^{\text{on}}} \right) =: \widehat{\sigma}. \end{aligned} \quad (54)$$

Here, the second and the third lines arise from Condition (43), the penultimate line invokes the definition (49) of  $\mathcal{T}_h^{\text{small}}$ , whereas the last line is valid since  $N = K^{\text{on}}/(3H)$  (see (15)).

## 6.2 Step 2: showing that $\pi^{\text{imitate}}$ (resp. $\pi^{\text{explore}}$ ) covers $\widehat{d}^{\text{off}}$ (resp. $d^{\pi^*}$ ) adequately

In this step, we aim to demonstrate the quality of the two exploration policies  $\pi^{\text{imitate}}$  and  $\pi^{\text{explore}}$ , designed for different purposes.

**Goodness of  $\pi^{\text{imitate}}$ .** We begin by assessing the quality of the exploration policy  $\pi^{\text{imitate}}$ . Towards this, we first make note of the following crude bound:

$$\frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi^{\sim} \mu^t}[\widehat{d}_h^{\pi}(s, a)]} \leq \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}}H}} \leq K^{\text{on}}H =: L.$$

In view of the convergence guarantees for FTRL (Shalev-Shwartz, 2012, Corollary 2.16), we see that: if  $\eta = \sqrt{\frac{\log A}{2T_{\max}L^2}} = \sqrt{\frac{\log A}{2T_{\max}(K^{\text{on}}H)^2}}$ , then running FTRL for  $T_{\max}$  iterations results in

$$\max_{a \in \mathcal{A}} \frac{1}{T_{\max}} \sum_{t=1}^{T_{\max}} \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi^{\sim} \mu^t}[\widehat{d}_h^{\pi}(s, a)]} - \frac{1}{T_{\max}} \sum_{t=1}^{T_{\max}} \sum_{a \in \mathcal{A}} \pi_h^t(a | s) \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi^{\sim} \mu^t}[\widehat{d}_h^{\pi}(s, a)]}$$



$$\leq K^{\text{on}} H \sqrt{\frac{2 \log A}{T_{\max}}} \quad (55)$$

for all  $s \in \mathcal{S}$  and  $1 \leq h \leq H$ . Therefore, recalling that  $\mu^{\text{imitate}} = \frac{1}{T_{\max}} \sum_{t=1}^{T_{\max}} \mu^t$  and applying Jensen's inequality yield

$$\begin{aligned} & \sum_{h=1}^H \sum_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}} H} + \mathbb{E}_{\pi \sim \mu^{\text{imitate}}} [\widehat{d}_h^{\pi}(s, a)]} \\ & \leq \sum_{h=1}^H \sum_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \frac{1}{T_{\max}} \sum_{t=1}^{T_{\max}} \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}} H} + \mathbb{E}_{\pi \sim \mu^t} [\widehat{d}_h^{\pi}(s, a)]} \\ & \leq \sum_{h=1}^H \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} \frac{1}{T_{\max}} \sum_{t=1}^{T_{\max}} \pi_h^t(a | s) \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}} H} + \mathbb{E}_{\pi \sim \mu^t} [\widehat{d}_h^{\pi}(s, a)]} + K^{\text{on}} H^2 S \sqrt{\frac{2 \log A}{T_{\max}}}, \end{aligned} \quad (56)$$

where the second inequality results from (55). In addition, it follows from the stopping rule (31) that

$$\sum_{h=1}^H \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} \pi_h^t(a | s) \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}} H} + \mathbb{E}_{\pi \sim \mu^t} [\widehat{d}_h^{\pi}(s, a)]} \leq 108SH. \quad (57)$$

As a consequence, combining (56) and (57) yields

$$\sum_{h \in [H]} \sum_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}} H} + \mathbb{E}_{\pi \sim \mu^{\text{imitate}}} [\widehat{d}_h^{\pi}(s, a)]} \leq 108SH + K^{\text{on}} H^2 S \sqrt{\frac{2 \log A}{T_{\max}}} \leq 109SH, \quad (58)$$

provided that  $T_{\max} \geq 2(K^{\text{on}} H)^2 \log A$ . The fact that the left-hand side of (58) is well-controlled suggests that  $\pi^{\text{imitate}}$  is able to cover  $\widehat{d}^{\text{off}}$  adequately, a crucial fact we shall rely on in the subsequent analysis.

**Goodness of  $\pi^{\text{explore}}$ .** Next, we turn attention to the other exploration policy  $\pi^{\text{explore}}$ , computed via Algorithm 5. The following performance guarantees have been established in Li et al. (2023, Section 3.2).

**Lemma 3.** *The distribution  $\mu^{\text{explore}} \in \Delta(\Pi)$  returned by Algorithm 5 satisfies*

$$\max_{\pi} \sum_{h=1}^H \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} \frac{\widehat{d}_h^{\pi}(s, a)}{\frac{1}{K^{\text{on}} H} + \mathbb{E}_{\pi' \sim \mu^{\text{explore}}} [\widehat{d}_h^{\pi'}(s, a)]} \leq 2HSA.$$

In light of the performance bound (18) for the subsequent offline RL approach, Lemma 3 suggests that  $\pi^{\text{explore}}$  is able to explore well with regards to the visitation of any policy  $\pi$  — including the optimal policy  $\pi^*$ .

### 6.3 Step 3: establishing the performance of offline RL

Now, we can readily proceed to analyze the performance of the model-based offline procedure described in Algorithm 6. In this subsection, we abuse the notation  $\widehat{P}$  to represent the empirical transition kernel constructed within the offline subroutine in Algorithm 7. Additionally, we introduce a  $S$ -dimensional vector  $d_h^{\pi^*} := [d_h^{\pi^*}(s)]_{s \in \mathcal{S}}$ .

#### 6.3.1 Step 3.1: error decomposition

To begin with, we convert the sub-optimality gap of the policy estimate  $\widehat{\pi}$  into several terms that shall be controlled separately. The following two preliminary facts, which have been established in Li et al. (2022), prove useful for this purpose.

**Lemma 4.** *With probability exceeding  $1 - \delta/3$ , one has*

$$N_h^{\text{main}}(s, a) \geq N_h^{\text{trim}}(s, a), \quad \forall (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$$

and

$$\langle d_j^{\pi^*}, V_j^* - V_j^{\widehat{\pi}} \rangle \leq 2 \sum_{h:h \geq j} \sum_{s,a} d_h^{\pi^*}(s, a) b_h(s, a), \quad 1 \leq h \leq H,$$

where  $b_h(s, a)$  is defined in line 6 of Algorithm 7.

In view of Lemma 4, we can derive, for all  $j \in [H]$ ,

$$\begin{aligned} \langle d_j^{\pi^*}, V_j^* - V_j^{\widehat{\pi}} \rangle &\leq 2 \sum_{h:h \geq j} \sum_{s,a} d_h^{\pi^*}(s, a) b_h(s, a) = 2 \sum_{h:h \geq j} \sum_s d_h^{\pi^*}(s, \pi_h^*(s)) b_h(s, \pi_h^*(s)) \\ &\leq 2 \sum_{h:h \geq j} \sum_{s: (s, \pi_h^*(s)) \notin \mathcal{T}_h} \sqrt{2d_h^{\pi^*}(s, \pi_h^*(s)) C^*(\sigma) \widehat{d}_h^{\text{off}}(s, \pi_h^*(s))} b_h(s, \pi_h^*(s)) + 2 \sum_{h:h \geq j} \sum_{(s,a) \in \mathcal{T}_h} d_h^{\pi^*}(s, a) b_h(s, a) \\ &\leq 2 \sum_{h:h \geq j} \sum_{s: (s, \pi_h^*(s)) \notin \mathcal{T}_h} \sqrt{2d_h^{\pi^*}(s, \pi_h^*(s)) C^*(\sigma) \widehat{d}_h^{\text{off}}(s, \pi_h^*(s))} b_h(s, \pi_h^*(s)) \\ &\quad + 4 \sum_{h:h \geq j} \sum_{(s,a) \in \mathcal{T}_h} \widehat{d}_h^{\pi^*}(s, a) b_h(s, a) + \frac{8H^2SA}{K^{\text{on}}} + \frac{53c_\xi H^6 S^4 A^4}{N} \log \frac{HSA}{\delta}. \end{aligned}$$

Here, the second line comes from (53), whereas the third line is due to Lemma 1 and the basic fact that  $b_h(s, a) \leq H$  (see line 6 of Algorithm 7). Substituting the definition of  $b_h$  (see line 6 of Algorithm 7) into the above display and applying Lemma 4, we arrive at

$$\begin{aligned} \langle d_j^{\pi^*}, V_j^* - V_j^{\widehat{\pi}} \rangle &\leq \sum_{h:h \geq j} \sum_s \max_{a: (s,a) \notin \mathcal{T}_h} \left\{ \sqrt{8d_h^{\pi^*}(s, a) C^*(\sigma) \widehat{d}_h^{\text{off}}(s, a)} \right. \\ &\quad \left. \min \left\{ \sqrt{\frac{c_b \log \frac{K}{\delta}}{N_h^{\text{trim}}(s, a)} \text{Var}_{\widehat{P}_h(\cdot|s,a)}(\widehat{V}_{h+1})} + \frac{c_b H \log \frac{K}{\delta}}{N_h^{\text{trim}}(s, a)}, H \right\} \right\} \\ &\quad + 4H \sum_{h:h \geq j} \sum_{(s,a) \in \mathcal{T}_h} \widehat{d}_h^{\pi^*}(s, a) \sqrt{\frac{c_b \log \frac{K}{\delta}}{N_h^{\text{trim}}(s, a) + 1} + \frac{8H^2SA}{K^{\text{on}}} + \frac{53c_\xi H^6 S^4 A^4}{N} \log \frac{K}{\delta}}, \quad (59) \end{aligned}$$

where we recall that  $c_b > 0$  is also an absolute constant used to specify  $b_h(s, a)$ .

It is worth noting that the right-hand side of (59) involves a variance term  $\text{Var}_{\widehat{P}_h(\cdot|s,a)}(\widehat{V}_{h+1})$  w.r.t. the empirical model  $\widehat{P}$ . As it turns out, the following lemma established in Li et al. (2022, Lemma 8) makes apparent the intimate connection between  $\text{Var}_{\widehat{P}_h(\cdot|s,a)}(\widehat{V}_{h+1})$  and  $\text{Var}_{P_h(\cdot|s,a)}(\widehat{V}_{h+1})$ .

**Lemma 5.** *With probability exceeding  $1 - \delta/3$ , we have, for all  $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ ,*

$$\text{Var}_{\widehat{P}_h(\cdot|s,a)}(\widehat{V}_{h+1}) \leq 2\text{Var}_{P_h(\cdot|s,a)}(\widehat{V}_{h+1}) + \frac{10H^2 \log \frac{K}{\delta}}{3N_h^{\text{trim}}(s, a)}.$$

Substituting the result of Lemma 5 into (59) leads to

$$\langle d_j^{\pi^*}, V_j^* - V_j^{\widehat{\pi}} \rangle \leq \gamma_1 + \gamma_2 + \frac{8H^2SA}{K^{\text{on}}} + \frac{53c_\xi H^6 S^4 A^4}{N} \log \frac{K}{\delta}, \quad (60)$$

where

$$\begin{aligned} \gamma_1 &:= \sum_{h:h \geq j} \sum_s \max_{a: (s,a) \notin \mathcal{T}_h} \left\{ 2\sqrt{2d_h^{\pi^*}(s, a) C^*(\sigma) \widehat{d}_h^{\text{off}}(s, a)} \min \left\{ \sqrt{\frac{2c_b \log \frac{K}{\delta}}{N_h^{\text{trim}}(s, a)} \text{Var}_{P_h(\cdot|s,a)}(\widehat{V}_{h+1})} + \frac{4c_b H \log \frac{K}{\delta}}{N_h^{\text{trim}}(s, a)}, H \right\} \right\}; \\ \gamma_2 &:= 4H \sum_{h:h \geq j} \sum_{(s,a) \in \mathcal{T}_h} \widehat{d}_h^{\pi^*}(s, a) \sqrt{\frac{c_b \log \frac{K}{\delta}}{N_h^{\text{trim}}(s, a) + 1}}. \end{aligned}$$

This leaves us with two terms to bound, which we shall accomplish separately in the ensuing two steps.

### 6.3.2 Step 3.2: controlling $\gamma_1$ in (60)

Regarding the first term  $\gamma_1$  on the right-hand side of (60), let us first define the set  $\mathcal{I}_h$  as follows:

$$\mathcal{I}_h := \left\{ (s, a) : \mathbb{E}_{\pi \sim \mu^{\text{imitate}}} [\widehat{d}_h^\pi(s, a)] \geq \frac{\xi}{N} \right\}, \quad (61)$$

where we remind the reader that  $\xi = c_\xi H^3 S^3 A^3 \log \frac{HSA}{\delta}$  for some constant  $c_\xi > 0$ . Armed with this set, we can deduce that

$$\begin{aligned} & \sum_{h:h \geq j} \sum_s \max_{a:(s,a) \notin \mathcal{I}_h \cup \mathcal{T}_h} \sqrt{2d_h^{\pi^*}(s, a) C^*(\sigma) \widehat{d}_h^{\text{off}}(s, a)} \\ & \leq \sum_{h:h \geq j} \sum_s 2 \max_{a:(s,a) \notin \mathcal{I}_h} C^*(\sigma) \widehat{d}_h^{\text{off}}(s, a) \\ & \leq 2C^*(\sigma) \left( \frac{1}{K^{\text{on}}H} + \frac{\xi}{N} \right) \sum_{h:h \geq j} \sum_s \max_a \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi \sim \mu^{\text{imitate}}} [\widehat{d}_h^\pi(s, a)]} \\ & \leq 218HSC^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{K^{\text{on}}H} \right), \end{aligned}$$

where the first inequality arises from (53), the penultimate line utilizes the definition (61) of  $\mathcal{I}_h$ , and the last line comes from (58). This in turn allows us to upper bound  $\gamma_1$  as follows:

$$\begin{aligned} \gamma_1 & \leq \sum_{h:h \geq j} \sum_s 2 \max_{a:(s,a) \in \mathcal{I}_h} \sqrt{2d_h^{\pi^*}(s, a) C^*(\sigma) \widehat{d}_h^{\text{off}}(s, a)} \min \left\{ \sqrt{\frac{2c_b \log \frac{HK}{\delta}}{N_h^{\text{trim}}(s, a)} \text{Var}_{P_h(\cdot|s,a)}(\widehat{V}_{h+1})} + \frac{4c_b H \log \frac{K}{\delta}}{N_h^{\text{trim}}(s, a)}, H \right\} \\ & \quad + \sum_{h:h \geq j} \sum_s 2 \max_{a:(s,a) \notin \mathcal{I}_h \cup \mathcal{T}_h} \sqrt{2d_h^{\pi^*}(s, a) C^*(\sigma) \widehat{d}_h^{\text{off}}(s, a)} \cdot H \\ & \leq \sum_{h:h \geq j} \sum_s 2 \max_{a:(s,a) \in \mathcal{I}_h} \sqrt{2d_h^{\pi^*}(s, a) C^*(\sigma) \widehat{d}_h^{\text{off}}(s, a)} \min \left\{ \sqrt{\frac{2c_b \log \frac{HK}{\delta}}{N_h^{\text{trim}}(s, a)} \text{Var}_{P_h(\cdot|s,a)}(\widehat{V}_{h+1})} + \frac{4c_b H \log \frac{K}{\delta}}{N_h^{\text{trim}}(s, a)}, H \right\} \\ & \quad + 436H^2 SC^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{K^{\text{on}}H} \right) \\ & \leq 16c_b \sum_{h:h \geq j} \sum_s \max_{a:(s,a) \in \mathcal{I}_h} \sqrt{2d_h^{\pi^*}(s, a) C^*(\sigma) \widehat{d}_h^{\text{off}}(s, a)} \sqrt{\frac{\text{Var}_{P_h(\cdot|s,a)}(\widehat{V}_{h+1}) + H}{N_h^{\text{trim}}(s, a) + 1/H} \log^2 \frac{K}{\delta}} \\ & \quad + 436H^2 SC^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{K^{\text{on}}H} \right), \quad (62) \end{aligned}$$

where the last line makes use of the elementary fact that  $\min \left\{ \frac{x}{y}, \frac{u}{w} \right\} \leq \frac{x+u}{y+w}$  for any  $x, y, u, w > 0$ .

In addition, note that for any  $s$  obeying  $\mathbb{E}_{\pi \sim \mu^{\text{imitate}}} [d_h^\pi(s)] \geq \xi/N$ , we have

$$\begin{aligned} \mathbb{E}[N_h^{\text{aux}}(s)] & = \frac{1}{4} K^{\text{off}} d_h^{\text{off}}(s) + \frac{1}{6} K^{\text{on}} \mathbb{E}_{\pi \sim \mu^{\text{imitate}}} [d_h^\pi(s)] + \frac{1}{6} K^{\text{on}} \mathbb{E}_{\pi \sim \mu^{\text{explore}}} [d_h^\pi(s)] \\ & \geq \frac{1}{6} K^{\text{on}} \mathbb{E}_{\pi \sim \mu^{\text{imitate}}} [d_h^\pi(s)] \geq \frac{1}{6} K^{\text{on}} \cdot \frac{\xi}{N} = \frac{1}{2} c_\xi H^4 S^3 A^3 \log \frac{HSA}{\delta}, \end{aligned}$$

where the last line invokes the definition of  $\mathcal{I}_h$  and the choice  $NH = \frac{1}{3} K^{\text{on}}$ . It can then be straightforwardly justified using elementary concentration inequalities (see, e.g., Alon and Spencer (2016, Appendix A.1)) that: with probability exceeding  $1 - \delta/10$ ,

$$N_h^{\text{aux}}(s) \geq \frac{1}{2} \mathbb{E}[N_h^{\text{aux}}(s)] \geq \frac{1}{4} c_\xi H^4 S^3 A^3 \log \frac{HSA}{\delta}$$

holds simultaneously for all  $(s, h) \in \mathcal{S} \times [H]$ , and as a result,

$$N_h^{\text{trim}}(s) \geq N_h^{\text{aux}}(s) - 10\sqrt{N_h^{\text{aux}}(s) \log \frac{HS}{\delta}} \geq \frac{1}{2}N_h^{\text{aux}}(s) \geq \frac{1}{4}\mathbb{E}[N_h^{\text{aux}}(s)] \geq \frac{1}{24}K^{\text{on}}\mathbb{E}_{\pi \sim \mu^{\text{imitate}}}[d_h^\pi(s)].$$

Moreover, for any  $(s, a) \in \mathcal{I}_h$  (cf. (61)), one can invoke Lemma 2 to obtain

$$\mathbb{E}_{\pi \sim \mu^{\text{imitate}}}[d_h^\pi(s)] \geq \frac{1}{3}\mathbb{E}_{\pi \sim \mu^{\text{imitate}}}[\widehat{d}_h^\pi(s)] \geq \frac{\xi}{3N}.$$

Applying the same concentration of measurement argument as above further reveals that:

$$N_h^{\text{trim}}(s, a) \geq \frac{1}{24}K^{\text{on}}\mathbb{E}_{\pi \sim \mu^{\text{imitate}}}[d_h^\pi(s, a)] \geq \frac{1}{72}K^{\text{on}}\mathbb{E}_{\pi \sim \mu^{\text{imitate}}}[\widehat{d}_h^\pi(s, a)]$$

any  $(s, a) \in \mathcal{I}_h$ . Substitution into (62) then gives

$$\begin{aligned} \gamma_1 \leq & 16c_b \sum_{h:h \geq j} \sum_s \max_{a:(s,a) \in \mathcal{I}_h} \sqrt{2d_h^{\pi^*}(s, a)C^*(\sigma)\widehat{d}_h^{\text{off}}(s, a)} \sqrt{\frac{\text{Var}_{P_h(\cdot|s,a)}(\widehat{V}_{h+1}) + H}{1/H + \frac{1}{72}K^{\text{on}}\mathbb{E}_{\pi \sim \mu^{\text{imitate}}}[\widehat{d}_h^\pi(s, a)]}} \log^2 \frac{K}{\delta} \\ & + 436H^2SC^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{K^{\text{on}}H} \right). \end{aligned} \quad (63)$$

By virtue of the Cauchy-Schwarz inequality, we can further derive

$$\begin{aligned} & \sum_{h:h \geq j} \sum_s \max_a \sqrt{d_h^{\pi^*}(s, a)C^*(\sigma)\widehat{d}_h^{\text{off}}(s, a)} \sqrt{\frac{\text{Var}_{P_h(\cdot|s,a)}(\widehat{V}_{h+1}) + H}{1/H + \frac{1}{72}K^{\text{on}}\mathbb{E}_{\pi \sim \mu^{\text{imitate}}}[\widehat{d}_h^\pi(s, a)]}} \\ & \leq \sqrt{\sum_{h:h \geq j} \sum_{s,a} d_h^{\pi^*}(s, a) \left( \text{Var}_{P_h(\cdot|s,a)}(\widehat{V}_{h+1}) + H \right)} \cdot \sqrt{\sum_{h:h \geq j} \sum_s \max_a \frac{C^*(\sigma)\widehat{d}_h^{\text{off}}(s, a)}{1/H + \frac{1}{72}K^{\text{on}}\mathbb{E}_{\pi \sim \mu^{\text{imitate}}}[\widehat{d}_h^\pi(s, a)]}}. \end{aligned} \quad (64)$$

To further control this term, we resort to the following lemma, whose proof is deferred to Section C.2.

**Lemma 6.** *With probability at least  $1 - \delta/6$ , we have, for all  $j \in [H]$ ,*

$$\sum_{h:h \geq j} \sum_{s,a} d_h^{\pi^*}(s, a) \text{Var}_{P_h(\cdot|s,a)}(\widehat{V}_{h+1}) \leq 5H^2,$$

provided that

$$\begin{aligned} K^{\text{on}} & \geq c_{11} (H^{18}S^{14}A^{14} + H^5S^4A^3C^*(\sigma)) \log^2 \frac{K}{\delta} \\ K^{\text{off}} & \geq c_{11}HS(C^*(\sigma) + A) \log \frac{K}{\delta} \end{aligned}$$

for some sufficiently large constant  $c_{11} > 0$ .

Putting Lemma 6 together with (58) and (64), we obtain

$$\sum_{h:h \geq j} \sum_s \max_a \sqrt{d_h^{\pi^*}(s, a)C^*(\sigma)\widehat{d}_h^{\text{off}}(s, a)} \sqrt{\frac{\text{Var}_{P_h(\cdot|s,a)}(\widehat{V}_{h+1}) + H}{1/H + \frac{1}{72}K^{\text{on}}\mathbb{E}_{\pi \sim \mu^{\text{imitate}}}[\widehat{d}_h^\pi(s, a)]}} \lesssim \sqrt{\frac{H^3SC^*(\sigma)}{K^{\text{on}}}}. \quad (65)$$

Substitution into (63) results in

$$\gamma_1 \leq \sqrt{\frac{H^3SC^*(\sigma) \log^2 \frac{K}{\delta}}{K^{\text{on}}}} + H^2SC^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{K^{\text{on}}H} \right). \quad (66)$$

Akin to (63) and (66), we can also focus on the offline dataset and obtain

$$\gamma_1 \lesssim \sqrt{\frac{H^3 SC^*(\sigma)}{K^{\text{off}}} \log^2 \frac{K}{\delta}} + H^2 SC^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{K^{\text{off}} H} \right). \quad (67)$$

Combine (66) and (67) to arrive at

$$\begin{aligned} \gamma_1 &\lesssim \min \left\{ \sqrt{\frac{H^3 SC^*(\sigma)}{K^{\text{on}}} \log^2 \frac{K}{\delta}}, \sqrt{\frac{H^3 SC^*(\sigma)}{K^{\text{off}}} \log^2 \frac{K}{\delta}} \right\} + H^2 SC^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{\min\{K^{\text{on}}, K^{\text{off}}\} H} \right) \\ &\lesssim \sqrt{\frac{H^3 SC^*(\sigma)}{\max\{K^{\text{on}}, K^{\text{off}}\}} \log^2 \frac{K}{\delta}} + H^2 SC^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{\min\{K^{\text{on}}, K^{\text{off}}\} H} \right) \\ &\lesssim \sqrt{\frac{H^3 SC^*(\sigma)}{K^{\text{on}} + K^{\text{off}}} \log^2 \frac{K}{\delta}} + H^2 SC^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{\min\{K^{\text{on}}, K^{\text{off}}\} H} \right). \end{aligned} \quad (68)$$

### 6.3.3 Step 3.3: controlling $\gamma_2$ in (60)

We now turn attention to the term  $\gamma_2$  on the right-hand side of (60). Akin to (63), we can deduce that

$$\gamma_2 \leq 16c_b H \sum_{h:h \geq j} \sum_{(s,a) \in \mathcal{T}_h} \widehat{d}_h^{\pi^*}(s,a) \sqrt{\frac{\log \frac{K}{\delta}}{1 + \frac{1}{72} K^{\text{on}} \mathbb{E}_{\pi \sim \mu^{\text{explore}}} [\widehat{d}_h^{\pi}(s,a)]}} + 436H^2 SA \left( \frac{\xi}{N} + \frac{1}{K^{\text{on}} H} \right). \quad (69)$$

The Cauchy-Schwarz inequality then tells us that

$$\begin{aligned} &\sum_{h:h \geq j} \sum_{(s,a) \in \mathcal{T}_h} \widehat{d}_h^{\pi^*}(s,a) \sqrt{\frac{1}{1 + \frac{1}{72} K^{\text{on}} \mathbb{E}_{\pi \sim \mu^{\text{explore}}} [\widehat{d}_h^{\pi}(s,a)]}} \\ &\leq \sqrt{\sum_{h:h \geq j} \sum_{(s,a) \in \mathcal{T}_h} \frac{\widehat{d}_h^{\pi^*}(s,a)}{1 + \frac{1}{72} K^{\text{on}} \mathbb{E}_{\pi \sim \mu^{\text{explore}}} [\widehat{d}_h^{\pi}(s,a)]}} \cdot \sqrt{\sum_{h:h \geq j} \sum_{(s,a) \in \mathcal{T}_h} \widehat{d}_h^{\pi^*}(s,a)} \\ &\leq 6 \sqrt{\frac{2HSA}{K^{\text{on}}}} \cdot \sqrt{\sum_{h:h \geq j} \sum_{(s,a) \in \mathcal{T}_h} \widehat{d}_h^{\pi^*}(s,a)} \\ &\leq 6 \sqrt{\frac{2HSA(2\widehat{\sigma} + HS\xi/N)}{K^{\text{on}}}}, \end{aligned}$$

where  $\widehat{\sigma}$  is defined in (54). Here, the penultimate line invokes Lemma 3, and the last line is valid since (according to Lemma 1 and (54))

$$\begin{aligned} \sum_{h:h \geq j} \sum_{(s,a) \in \mathcal{T}_h} \widehat{d}_h^{\pi^*}(s,a) &= \sum_{h:h \geq j} \sum_{s:(s,\pi^*(s)) \in \mathcal{T}_h} \widehat{d}_h^{\pi^*}(s,\pi^*(s)) \\ &\leq 2 \sum_{h:h \geq j} \sum_{s:(s,\pi^*(s)) \in \mathcal{T}_h} d_h^{\pi^*}(s,\pi^*(s)) + \frac{HS\xi}{N} \\ &\leq 2\widehat{\sigma} + \frac{HS\xi}{N}. \end{aligned}$$

Substitution of the above inequality into (69) yields

$$\gamma_2 \leq 96c_b \sqrt{\frac{2H^3 SA(2\widehat{\sigma} + HS\xi/N)}{K^{\text{on}}}} \log \frac{HK}{\delta} + 436H^2 SA \left( \frac{\xi}{N} + \frac{1}{K^{\text{on}} H} \right). \quad (70)$$

### 6.3.4 Step 3.4: putting all pieces together

To finish up, combining (60),(68) and (70) reveals that: with probability at least  $1 - \delta$ , one has

$$\begin{aligned}
V_1^*(\rho) - V^{\hat{\pi}}(\rho) &= \langle d_1^{\pi^*}, V_1^* - V^{\hat{\pi}} \rangle \\
&\lesssim \sqrt{\frac{H^3 S C^*(\sigma) \log^2 \frac{K}{\delta}}{K^{\text{on}} + K^{\text{off}}}} + \sqrt{\frac{H^4 S A \sigma \log \frac{K}{\delta}}{K^{\text{on}}}} + \sqrt{\frac{H^4 S^2 A C^*(\sigma) \log^2 \frac{K}{\delta}}{K^{\text{off}} K^{\text{on}}}} \\
&\quad + \sqrt{\frac{H^8 S^6 A^5 C^*(\sigma) \log^2 \frac{K}{\delta}}{N K^{\text{on}}}} + \sqrt{\frac{H^4 S^3 A^2 C^*(\sigma) \log \frac{K}{\delta}}{K K^{\text{on}}}} \\
&\quad + \frac{H^6 S^4 A^4 + H^5 S^4 A^3 C^*(\sigma)}{N} \log \frac{K}{\delta} + \frac{H^2 S(C^*(\sigma) + A)}{\min\{K^{\text{on}}, K^{\text{off}}\}} \\
&\lesssim \sqrt{\frac{H^3 S C^*(\sigma) \log^2 \frac{K}{\delta}}{K^{\text{on}} + K^{\text{off}}}} + \sqrt{\frac{H^4 S A \sigma \log \frac{K}{\delta}}{K^{\text{on}}}} + \frac{H^6 S^4 A^4 + H^5 S^4 A^3 C^*(\sigma)}{K^{\text{on}}} \log^2 \frac{K}{\delta} + \frac{H^2 S(C^*(\sigma) + A)}{K^{\text{off}}},
\end{aligned} \tag{71}$$

where the last inequality holds true as long as  $\min\{K^{\text{off}}, K^{\text{on}}\} \gtrsim HSA$ . Taking the right-hand side of (71) to be no larger than  $\varepsilon$ , we immediately establish Theorem 1 under the sample complexity assumption in this theorem.

## 7 Discussion

In this paper, we have studied the policy fine-tuning problem of practical interest, where one is allowed to exploit pre-collected historical data to facilitate and improve online RL. We have proposed a three-stage algorithm tailored to the tabular setting, which attains provable sample size savings compared with both pure online RL and pure offline RL algorithms. Our algorithm design has leveraged key insights from recent advances in both model-based offline RL and reward-agnostic online RL.

While the proposed algorithm achieves provable sample efficiency, this cannot be guaranteed unless the sample size already surpasses a fairly large threshold (in other words, the algorithm imposes a high burn-in cost). It would be of great interest to see whether one can achieve sample optimality for the entire  $\varepsilon$ -range. Another issue arises from the computation side: even though the proposed algorithm can be implemented in polynomial time, the computational complexity of the Frank-Wolfe-type subroutine might already be too expensive for solving problems with enormous dimensionality. Can we hope to further accelerate it to make it practically more appealing? Finally, it might also be interesting to study sample-efficient hybrid RL in the presence of low-complexity function approximation, in the hope of further reducing sample complexity.

## Acknowledgements

Y. Chen is supported in part by the Alfred P. Sloan Research Fellowship, the Google Research Scholar Award, the AFOSR grant FA9550-22-1-0198, the ONR grant N00014-22-1-2354, and the NSF grants CCF-2221009 and CCF-1907661. Y. Chi are supported in part by the grants ONR N00014-19-1-2404, NSF CCF-2106778, DMS-2134080 and CNS-2148212. J. D. Lee acknowledges support of the ARO under MURI Award W911NF-11-1-0304, the Sloan Research Fellowship, the NSF grants CCF 2002272, IIS 2107304 and CIF 2212262, the ONR Young Investigator Award, and the NSF CAREER Award 2144994.

## References

- Agarwal, A., Kakade, S., Krishnamurthy, A., and Sun, W. (2020). Flambe: Structural complexity and representation learning of low rank mdps. *arXiv preprint arXiv:2006.10814*.
- Alon, N. and Spencer, J. H. (2016). *The probabilistic method*. John Wiley & Sons.

- Auer, P. and Ortner, R. (2006). Logarithmic online regret bounds for undiscounted reinforcement learning. *Advances in neural information processing systems*, 19.
- Azar, M. G., Osband, I., and Munos, R. (2017). Minimax regret bounds for reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 263–272. JMLR.org.
- Bai, Y., Xie, T., Jiang, N., and Wang, Y.-X. (2019). Provably efficient  $q$ -learning with low switching cost. In *Advances in Neural Information Processing Systems*, volume 32, pages 8002–8011.
- Berner, C., Brockman, G., Chan, B., Cheung, V., Debiak, P., Dennison, C., Farhi, D., Fischer, Q., Hashme, S., Hesse, C., et al. (2019). Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*.
- Bertsekas, D. P. (2017). *Dynamic programming and optimal control (4th edition)*. Athena Scientific.
- Brafman, R. I. and Tennenholtz, M. (2002). R-max-a general polynomial time algorithm for near-optimal reinforcement learning. *Journal of Machine Learning Research*, 3(Oct):213–231.
- Brambilla, M., Ferrante, E., Birattari, M., and Dorigo, M. (2013). Swarm robotics: a review from the swarm engineering perspective. *Swarm Intelligence*, 7(1):1–41.
- Bubeck, S. (2015). Convex optimization: Algorithms and complexity. *Foundations and Trends<sup>®</sup> in Machine Learning*, 8(3-4):231–357.
- Chen, J. and Jiang, N. (2019). Information-theoretic considerations in batch reinforcement learning. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pages 1042–1051.
- Chen, J., Modi, A., Krishnamurthy, A., Jiang, N., and Agarwal, A. (2022). On the statistical efficiency of reward-free exploration in non-linear rl. *arXiv preprint arXiv:2206.10770*.
- Cui, Q. and Du, S. S. (2022). When is offline two-player zero-sum Markov game solvable? *arXiv preprint arXiv:2201.03522*.
- 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., Ménard, P., Kaufmann, E., and Valko, M. (2021). Episodic reinforcement learning in finite mdps: Minimax lower bounds revisited. In *Algorithmic Learning Theory*, pages 578–598. PMLR.
- Dong, K., Wang, Y., Chen, X., and Wang, L. (2019). Q-learning with UCB exploration is sample efficient for infinite-horizon MDP. *arXiv preprint arXiv:1901.09311*.
- Du, S. S., Kakade, S. M., Lee, J. D., Lovett, S., Mahajan, G., Sun, W., and Wang, R. (2021). Bilinear classes: A structural framework for provable generalization in rl.
- Foster, D. J., Kakade, S. M., Qian, J., and Rakhlin, A. (2021). The statistical complexity of interactive decision making. *arXiv preprint arXiv:2112.13487*.
- Hester, T., Vecerik, M., Pietquin, O., Lanctot, M., Schaul, T., Piot, B., Horgan, D., Quan, J., Sendonaris, A., Osband, I., et al. (2018). Deep Q-learning from demonstrations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.
- Huang, R., Yang, J., and Liang, Y. (2022). Safe exploration incurs nearly no additional sample complexity for reward-free RL. *arXiv preprint arXiv:2206.14057*.
- Jaksch, T., Ortner, R., and Auer, P. (2010). Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*.
- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. (2018). Is Q-learning provably efficient? In *Advances in Neural Information Processing Systems*, pages 4863–4873.

- Jin, C., Krishnamurthy, A., Simchowitz, M., and Yu, T. (2020a). Reward-free exploration for reinforcement learning. In *International Conference on Machine Learning*, pages 4870–4879. PMLR.
- Jin, C., Liu, Q., and Miryoosefi, S. (2021a). Bellman eluder dimension: New rich classes of RL problems, and sample-efficient algorithms. *Advances in Neural Information Processing Systems*, 34.
- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. (2020b). Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR.
- Jin, Y., Yang, Z., and Wang, Z. (2020c). Is pessimism provably efficient for offline rl? *arXiv preprint arXiv:2012.15085*.
- Jin, Y., Yang, Z., and Wang, Z. (2021b). Is pessimism provably efficient for offline RL? In *International Conference on Machine Learning*, pages 5084–5096. PMLR.
- Kalashnikov, D., Irpan, A., Pastor, P., Ibarz, J., Herzog, A., Jang, E., Quillen, D., Holly, E., Kalakrishnan, M., and Vanhoucke, V. (2018). Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on Robot Learning*, pages 651–673. PMLR.
- Kaufmann, E., Ménard, P., Domingues, O. D., Jonsson, A., Leurent, E., and Valko, M. (2021). Adaptive reward-free exploration. In *Algorithmic Learning Theory*, pages 865–891. PMLR.
- Kumar, A., Zhou, A., Tucker, G., and Levine, S. (2020). Conservative q-learning for offline reinforcement learning. In *Conference on Neural Information Processing Systems (NeurIPS)*.
- Li, G., Chen, Y., Chi, Y., Gu, Y., and Wei, Y. (2021a). Sample-efficient reinforcement learning is feasible for linearly realizable MDPs with limited revisiting. *Advances in Neural Information Processing Systems*, 34:16671–16685.
- Li, G., Shi, L., Chen, Y., Chi, Y., and Wei, Y. (2022). Settling the sample complexity of model-based offline reinforcement learning. *arXiv preprint arXiv:2204.05275*.
- Li, G., Shi, L., Chen, Y., Gu, Y., and Chi, Y. (2021b). Breaking the sample complexity barrier to regret-optimal model-free reinforcement learning. *Advances in Neural Information Processing Systems*, 34:17762–17776.
- Li, G., Yan, Y., Chen, Y., and Fan, J. (2023). Minimax-optimal reward-agnostic exploration in reinforcement learning. *arXiv preprint arXiv:2304.0727*.
- Liu, Y., Swaminathan, A., Agarwal, A., and Brunskill, E. (2020). Provably good batch reinforcement learning without great exploration. *arXiv preprint arXiv:2007.08202*.
- Ménard, P., Domingues, O. D., Jonsson, A., Kaufmann, E., Leurent, E., and Valko, M. (2021a). Fast active learning for pure exploration in reinforcement learning. In *International Conference on Machine Learning*, pages 7599–7608. PMLR.
- Ménard, P., Domingues, O. D., Shang, X., and Valko, M. (2021b). UCB momentum Q-learning: Correcting the bias without forgetting. In *International Conference on Machine Learning*, pages 7609–7618.
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. (2013). Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*.
- Munos, R. and Szepesvári, C. (2008). Finite-time bounds for fitted value iteration. In *Journal of Machine Learning Research*, volume 9, pages 815–857.
- Nair, A., Gupta, A., Dalal, M., and Levine, S. (2020). Awac: Accelerating online reinforcement learning with offline datasets. *arXiv preprint arXiv:2006.09359*.
- Nair, A., McGrew, B., Andrychowicz, M., Zaremba, W., and Abbeel, P. (2018). Overcoming exploration in reinforcement learning with demonstrations. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 6292–6299. IEEE.



- Nakamoto, M., Zhai, Y., Singh, A., Mark, M. S., Ma, Y., Finn, C., Kumar, A., and Levine, S. (2023). Cal-QL: Calibrated offline RL pre-training for efficient online fine-tuning. *arXiv preprint arXiv:2303.05479*.
- Qiao, D. and Wang, Y.-X. (2022). Near-optimal deployment efficiency in reward-free reinforcement learning with linear function approximation. *arXiv preprint arXiv:2210.00701*.
- Rajeswaran, A., Kumar, V., Gupta, A., Vezzani, G., Schulman, J., Todorov, E., and Levine, S. (2017). Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. *arXiv preprint arXiv:1709.10087*.
- Rashidinejad, P., Zhu, B., Ma, C., Jiao, J., and Russell, S. (2021). Bridging offline reinforcement learning and imitation learning: A tale of pessimism. *Advances in Neural Information Processing Systems*, 34:11702–11716.
- Ross, S. and Bagnell, J. A. (2012). Agnostic system identification for model-based reinforcement learning. *International Conference on Machine Learning*.
- Shalev-Shwartz, S. (2012). Online learning and online convex optimization. *Foundations and Trends® in Machine Learning*, 4(2):107–194.
- Shalev-Shwartz, S., Shammah, S., and Shashua, A. (2016). Safe, multi-agent, reinforcement learning for autonomous driving. *arXiv preprint arXiv:1610.03295*.
- Shi, L. and Chi, Y. (2022). Distributionally robust model-based offline reinforcement learning with near-optimal sample complexity. *arXiv preprint arXiv:2208.05767*.
- Shi, L., Li, G., Wei, Y., Chen, Y., and Chi, Y. (2022). Pessimistic Q-learning for offline reinforcement learning: Towards optimal sample complexity. In *International Conference on Machine Learning*, pages 19967–20025. PMLR.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. (2016). Mastering the game of Go with deep neural networks and tree search. *nature*, 529(7587):484–489.
- Song, Y., Zhou, Y., Sekhari, A., Bagnell, J. A., Krishnamurthy, A., and Sun, W. (2022). Hybrid RL: Using both offline and online data can make RL efficient. *arXiv preprint arXiv:2210.06718*.
- Uehara, M. and Sun, W. (2021). Pessimistic model-based offline RL: PAC bounds and posterior sampling under partial coverage. In *arXiv preprint arXiv:2107.06226*.
- Vecerik, M., Hester, T., Scholz, J., Wang, F., Pietquin, O., Piot, B., Heess, N., Rothörl, T., Lampe, T., and Riedmiller, M. (2017). Leveraging demonstrations for deep reinforcement learning on robotics problems with sparse rewards. *arXiv preprint arXiv:1707.08817*.
- Vinyals, O., Babuschkin, I., Czarnecki, W. M., Mathieu, M., Dudzik, A., Chung, J., Choi, D. H., Powell, R., Ewalds, T., Georgiev, P., et al. (2019). Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354.
- Wagenmaker, A. and Pacchiano, A. (2022). Leveraging offline data in online reinforcement learning. *arXiv preprint arXiv:2211.04974*.
- Wagenmaker, A. J., Chen, Y., Simchowitz, M., Du, S., and Jamieson, K. (2022). Reward-free RL is no harder than reward-aware RL in linear Markov decision processes. In *International Conference on Machine Learning*, pages 22430–22456.
- Wang, R., Du, S. S., Yang, L., and Salakhutdinov, R. R. (2020). On reward-free reinforcement learning with linear function approximation. *Advances in neural information processing systems*, 33:17816–17826.
- Xie, T., Cheng, C.-A., Jiang, N., Mineiro, P., and Agarwal, A. (2021a). Bellman-consistent pessimism for offline reinforcement learning. *arXiv preprint arXiv:2106.06926*.

- Xie, T., Jiang, N., Wang, H., Xiong, C., and Bai, Y. (2021b). Policy finetuning: Bridging sample-efficient offline and online reinforcement learning. *arXiv preprint arXiv:2106.04895*.
- Yan, Y., Li, G., Chen, Y., and Fan, J. (2022). Model-based reinforcement learning is minimax-optimal for offline zero-sum Markov games. *arXiv preprint arXiv:2206.04044*.
- Yin, M., Bai, Y., and Wang, Y.-X. (2021). Near-optimal offline reinforcement learning via double variance reduction. *Advances in neural information processing systems*, 34:7677–7688.
- Zanette, A., Lazaric, A., Kochenderfer, M., and Brunskill, E. (2020). Learning near optimal policies with low inherent bellman error. In *International Conference on Machine Learning*, pages 10978–10989. PMLR.
- Zhan, W., Huang, B., Huang, A., Jiang, N., and Lee, J. (2022). Offline reinforcement learning with realizability and single-policy concentrability. In *Conference on Learning Theory*, pages 2730–2775. PMLR.
- Zhang, W., Zhou, D., and Gu, Q. (2021a). Reward-free model-based reinforcement learning with linear function approximation. *Advances in Neural Information Processing Systems*, 34:1582–1593.
- Zhang, X., Ma, Y., and Singla, A. (2020a). Task-agnostic exploration in reinforcement learning. *Advances in Neural Information Processing Systems*, 33:11734–11743.
- Zhang, Z., Du, S., and Ji, X. (2021b). Near optimal reward-free reinforcement learning. In *International Conference on Machine Learning*, pages 12402–12412. PMLR.
- Zhang, Z., Zhou, Y., and Ji, X. (2020b). Model-free reinforcement learning: from clipped pseudo-regret to sample complexity. *arXiv preprint arXiv:2006.03864*.
- Zhou, D., Gu, Q., and Szepesvari, C. (2021a). Nearly minimax optimal reinforcement learning for linear mixture markov decision processes. In *Conference on Learning Theory*, pages 4532–4576. PMLR.
- Zhou, Z., Zhou, Z., Bai, Q., Qiu, L., Blanchet, J., and Glynn, P. (2021b). Finite-sample regret bound for distributionally robust offline tabular reinforcement learning. In *International Conference on Artificial Intelligence and Statistics*, pages 3331–3339. PMLR.

## A Useful algorithmic subroutines from prior works

In this section, we provide precise descriptions of several useful algorithmic subroutines that have been developed in recent works. The algorithm procedures are directly quoted from these prior works, with slight modification.

### A.1 Subroutine: occupancy estimation for any policy $\pi$

The first subroutine we'd like to describe is concerned with estimating the occupancy distribution  $d^\pi$  induced by any policy  $\pi$ , based on a carefully designed exploration strategy. This algorithm, proposed by Li et al. (2023), seeks to estimate  $\{d_h^\pi\}$  step by step (i.e., from  $h = 1, \dots, H$ ). For each  $h$ , it computes an appropriate exploration policy  $\pi^{\text{explore},h}$  to adequately explore what happens between step  $h$  and step  $h + 1$ , and then collect  $N$  sample trajectories using  $\pi^{\text{explore},h}$ . These turns allow us to estimate the occupancy distribution  $d_{h+1}^\pi$  for step  $h + 1$ . See Algorithm 3 for a precise description.

---

**Algorithm 3:** Subroutine for estimating occupancy distributions for any policy  $\pi$  (Li et al., 2023).

---

- 1 **Input:** target success probability  $1 - \delta$ , threshold  $\xi = c_\xi H^3 S^3 A^3 \log(HSA/\delta)$ .  
*/\* Estimate occupancy distributions for step 1. \*/*
- 2 Draw  $N$  independent episodes (using arbitrary policies), whose initial states are i.i.d. drawn from  $s_1^{n,0} \stackrel{\text{i.i.d.}}{\sim} \rho$  ( $1 \leq n \leq N$ ). Define the following functions

$$\widehat{d}_1^\pi(s) = \frac{1}{N} \sum_{n=1}^N \mathbb{1}\{s_1^{n,0} = s\}, \quad \widehat{d}_1^\pi(s, a) = \widehat{d}_1^\pi(s) \pi_1(a | s) \quad (72)$$

for any deterministic policy  $\pi : \mathcal{S} \times [H] \rightarrow \Delta(\mathcal{A})$  and any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . (Note that these functions are defined for future use and not computed for the moment, as we have not specified policy  $\pi$ .)

*/\* Estimate occupancy distributions for steps 2, ..., H. \*/*

- 3 **for**  $h = 1$  **to**  $H - 1$  **do**
  - 4 */\* Collect  $N$  sample trajectories using a suitable exploration policy. \*/*  
Call Algorithm 4 to compute an exploration policy  $\pi^{\text{explore},h}$  and compute an estimate  $\widehat{P}_h$  of the true transition kernel  $P_h$ .  
*/\* Specify how to compute  $\widehat{d}_{h+1}^\pi$  for any policy  $\pi$ . \*/*
  - 5 For any deterministic policy  $\pi : \mathcal{S} \times [H] \rightarrow \Delta(\mathcal{A})$  and any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , define
$$\widehat{d}_{h+1}^\pi(s) = \langle \widehat{P}_h(s | \cdot, \cdot), \widehat{d}_h^\pi(\cdot, \cdot) \rangle, \quad \widehat{d}_{h+1}^\pi(s, a) = \widehat{d}_{h+1}^\pi(s) \pi_{h+1}(a | s). \quad (73)$$

---

We note, however, that Algorithm 3 requires another subroutine to compute a suitable exploration policy  $\pi^{\text{explore},h}$ . As it turns out, this can be accomplished by approximately solving the following problem

$$\widehat{\mu}^h \approx \arg \max_{\mu \in \Delta(\Pi)} \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \log \left[ \frac{1}{KH} + \mathbb{E}_{\pi \sim \mu} [\widehat{d}_h^\pi(s, a)] \right] \quad (74)$$

via the Frank-Wolfe algorithm and returning  $\pi^{\text{explore},h} = \mathbb{E}_{\pi \sim \widehat{\mu}^h} [\pi]$ . See Algorithm 4 for details.

### A.2 Subroutine: reward-agnostic online exploration

Based on the estimated occupancy distributions specified in Algorithm 3, Li et al. (2023) proposed a reward-independent online exploration scheme that proves useful in exploring an unknown environment. In a nutshell, this scheme computes a desired exploration policy by approximately solving the following optimization sub-problem again using the Frank-Wolfe algorithm:

$$\mu^{\text{explore}} \approx \arg \max_{\mu \in \Delta(\Pi)} \left\{ \sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \log \left[ \frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi \sim \mu} [\widehat{d}_h^\pi(s, a)] \right] \right\}. \quad (75)$$

---

**Algorithm 4:** Subroutine for computing the exploration policy for step  $h$  in occupancy estimation (Li et al., 2023).

---

1 **Initialize:**  $\mu^{(0)} = \mathbb{1}_{\pi_{\text{init}}}$  for an arbitrary policy  $\pi_{\text{init}} \in \Pi$ ,  $T_{\text{max}} = \lfloor 50SA \log(KH) \rfloor$ .

2 **for**  $t = 0$  **to**  $T_{\text{max}}$  **do**

/\* find the optimal policy \*/

3 Compute the optimal deterministic policy  $\pi^{(t),b}$  of the augmented MDP

$\mathcal{M}_b^h = (\mathcal{S} \cup \{s_{\text{aug}}\}, \mathcal{A}, H, \widehat{P}^{\text{aug},h}, r_b^h)$ , where  $s_{\text{aug}}$  is an augmented state,

$$r_{b,j}^h(s, a) = \begin{cases} \frac{1}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi \sim \mu^{(t)}}[\widehat{d}_h^\pi(s, a)]}, & \text{if } (s, a, j) \in \mathcal{S} \times \mathcal{A} \times \{h\}, \\ 0, & \text{if } s = s_{\text{aug}} \text{ or } j \neq h, \end{cases} \quad (76)$$

and the augmented probability transition kernel is defined as

$$\widehat{P}_j^{\text{aug},h}(s' | s, a) = \begin{cases} \widehat{P}_j(s' | s, a), & \text{if } s' \in \mathcal{S} \\ 1 - \sum_{s' \in \mathcal{S}} \widehat{P}_j(s' | s, a), & \text{if } s' = s_{\text{aug}} \end{cases} \quad \text{for all } (s, a, j) \in \mathcal{S} \times \mathcal{A} \times [h]; \quad (77a)$$

$$\widehat{P}_j^{\text{aug},h}(s' | s, a) = \mathbb{1}(s' = s_{\text{aug}}) \quad \text{if } s = s_{\text{aug}} \text{ or } j > h. \quad (77b)$$

Let  $\pi^{(t)}$  be the corresponding optimal deterministic policy of  $\pi^{(t),b}$  in the original state space.

4 Compute // choose the stepsize

$$\alpha_t = \frac{\frac{1}{SA}g(\pi^{(t)}, \widehat{d}, \mu^{(t)}) - 1}{g(\pi^{(t)}, \widehat{d}, \mu^{(t)}) - 1}, \quad \text{where } g(\pi, \widehat{d}, \mu) = \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{\frac{1}{K^{\text{on}}H} + \widehat{d}_h^\pi(s, a)}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi \sim \mu}[\widehat{d}_h^\pi(s, a)]}.$$

Here,  $\widehat{d}_h^\pi(s, a)$  is computed via (72) for  $h = 1$ , and (73) for  $h \geq 2$ .

6 If  $g(\pi^{(t)}, \widehat{d}, \mu^{(t)}) \leq 2SA$  then exit for-loop. // stopping rule

7 Update // Frank-Wolfe update

8

$$\mu^{(t+1)} = (1 - \alpha_t) \mu^{(t)} + \alpha_t \mathbb{1}_{\pi^{(t)}}.$$

9 Set  $\pi^{\text{explore},h} = \mathbb{E}_{\pi \sim \mu^{(t)}}[\pi]$  with  $\widehat{\mu}^h = \mu^{(t)}$ . // The final exploration policy for step  $h$ .

/\* Draw samples using  $\pi^{\text{explore},h}$  to estimate the transition kernel. \*/

10 Draw  $N$  independent trajectories  $\{s_1^{n,h}, a_1^{n,h}, \dots, s_{h+1}^{n,h}\}_{1 \leq n \leq N}$  using policy  $\pi^{\text{explore},h}$  and compute

$$\widehat{P}_h(s' | s, a) = \frac{\mathbb{1}(N_h(s, a) > \xi)}{\max\{N_h(s, a), 1\}} \sum_{n=1}^N \mathbb{1}(s_h^{n,h} = s, a_h^{n,h} = a, s_{h+1}^{n,h} = s'), \quad \forall (s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S},$$

where  $N_h(s, a) = \sum_{n=1}^N \mathbb{1}\{s_h^{n,h} = s, a_h^{n,h} = a\}$ .

11 **Output:** the exploration policy  $\pi^{\text{explore},h}$ , the weight  $\widehat{\mu}^h$ , and the estimated kernel  $\widehat{P}_h$ .

---

The resulting policy takes the form of a mixture of deterministic policies, as given by  $\pi^{\text{explore}} = \mathbb{E}_{\pi \sim \mu^{\text{explore}}}[\pi]$ . This exploration policy is then employed to execute the MDP for a number of times in order to collect enough information about the unknowns. See Algorithm 5 for the whole procedure.

### A.3 Subroutine: pessimistic model-based offline RL

Given a historical dataset containing a collection of statistically independent sample trajectories, Li et al. (2022) came up with a model-based approach that enjoys provable minimax optimality. This approach first employs a two-fold subsampling trick in order to decouple the statistical dependency across different steps of a single trajectory. After this subsampling step, this approach resorts to the principle of pessimism in the face of uncertainty, which employs value iteration but penalizes the updates via proper variance-aware

---

**Algorithm 5:** Subroutine for computing the desired online exploration policy (Li et al., 2023).

---

1 **Initialize:**  $\mu_b^{(0)} = \delta_{\pi_{\text{init}}}$  for an arbitrary policy  $\pi_{\text{init}} \in \Pi$ ,  $T_{\text{max}} = \lfloor 50SAH \log(KH) \rfloor$ .

2 **for**  $t = 0$  **to**  $T_{\text{max}}$  **do**

/\* find the optimal policy \*/

3 Compute the optimal deterministic policy  $\pi^{(t),b}$  of the MDP  $\mathcal{M}_b = (\mathcal{S} \cup \{s_{\text{aug}}\}, \mathcal{A}, H, \widehat{P}^{\text{aug}}, r_b)$ , where  $s_{\text{aug}}$  is an augmented state,

$$r_{b,h}(s, a) = \begin{cases} \frac{1}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi \sim \mu_b^{(t)}}[\widehat{d}_h^\pi(s, a)]}, & \text{if } (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H], \\ 0, & \text{if } (s, a, h) \in \{s_{\text{aug}}\} \times \mathcal{A} \times [H], \end{cases} \quad (78)$$

and the augmented probability transition kernel is given by

$$\widehat{P}_h^{\text{aug}}(s' | s, a) = \begin{cases} \widehat{P}_h(s' | s, a), & \text{if } s' \in \mathcal{S} \\ 1 - \sum_{s' \in \mathcal{S}} \widehat{P}_h(s' | s, a), & \text{if } s' = s_{\text{aug}} \end{cases} \quad \text{for all } (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]; \quad (79a)$$

$$\widehat{P}_h^{\text{aug}}(s' | s_{\text{aug}}, a) = \mathbf{1}(s' = s_{\text{aug}}) \quad \text{for all } (a, h) \in \mathcal{A} \times [H]. \quad (79b)$$

Let  $\pi^{(t)}$  be the corresponding optimal deterministic policy of  $\pi^{(t),b}$  in the original state space.

4 Compute // choose the stepsize

$$\alpha_t = \frac{\frac{1}{SAH}g(\pi^{(t)}, \widehat{d}, \mu_b^{(t)}) - 1}{g(\pi^{(t)}, \widehat{d}, \mu_b^{(t)}) - 1}, \quad \text{where} \quad g(\pi, \widehat{d}, \mu) = \sum_{h=1}^H \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{\frac{1}{K^{\text{on}}H} + \widehat{d}_h^\pi(s, a)}{\frac{1}{K^{\text{on}}H} + \mathbb{E}_{\pi \sim \mu}[\widehat{d}_h^\pi(s, a)]}.$$

Here,  $\widehat{d}_h^\pi(s, a)$  is computed via (72) for  $h = 1$ , and (73) for  $h \geq 2$ .

6 If  $g(\pi^{(t)}, \widehat{d}, \mu_b^{(t)}) \leq 2HSA$  then exit for-loop. // stopping rule

7 Update // Frank-Wolfe update

8

$$\mu_b^{(t+1)} = (1 - \alpha_t) \mu_b^{(t)} + \alpha_t \mathbf{1}_{\pi^{(t)}}.$$

9 **Output:** the exploration policy  $\pi^{\text{explore}} = \mathbb{E}_{\pi \sim \mu_b^{(t)}}[\pi]$  and the associated weight  $\mu^{\text{explore}} = \mu_b^{(t)}$ .

---

penalization (i.e., Bernstein-style lower confidence bounds). Details can be found in Algorithm 6.

---

**Algorithm 6:** A pessimistic model-based offline RL algorithm (Li et al., 2022).

---

1 **Input:** a dataset  $\mathcal{D}$ ; reward function  $r$ . Let  $K_0$  denote the number of sample trajectories in  $\mathcal{D}$ .

2 **Subsampling:** run the following procedure to generate the subsampled dataset  $\mathcal{D}^{\text{trim}}$ .

- 1) *Data splitting.* Split  $\mathcal{D}$  into two halves:  $\mathcal{D}^{\text{main}}$  (which contains the first  $K_0/2$  trajectories), and  $\mathcal{D}^{\text{aux}}$  (which contains the remaining  $K_0/2$  trajectories); we let  $N_h^{\text{main}}(s)$  (resp.  $N_h^{\text{aux}}(s)$ ) denote the number of sample transitions in  $\mathcal{D}^{\text{main}}$  (resp.  $\mathcal{D}^{\text{aux}}$ ) that transition from state  $s$  at step  $h$ .
- 2) *Lower bounding*  $\{N_h^{\text{main}}(s)\}$  using  $\mathcal{D}^{\text{aux}}$ . For each  $s \in \mathcal{S}$  and  $1 \leq h \leq H$ , compute
$$N_h^{\text{trim}}(s) := \max \left\{ N_h^{\text{aux}}(s) - 10 \sqrt{N_h^{\text{aux}}(s) \log \frac{HS}{\delta}}, 0 \right\}; \quad (80)$$
- 3) *Random subsampling.* Let  $\mathcal{D}^{\text{main}'}$  be the set of all sample transitions (i.e., the quadruples taking the form  $(s, a, h, s')$ ) from  $\mathcal{D}^{\text{main}}$ . Subsample  $\mathcal{D}^{\text{main}'}$  to obtain  $\mathcal{D}^{\text{trim}}$ , such that for each  $(s, h) \in \mathcal{S} \times [H]$ ,  $\mathcal{D}^{\text{trim}}$  contains  $\min\{N_h^{\text{trim}}(s), N_h^{\text{main}}(s)\}$  sample transitions randomly drawn from  $\mathcal{D}^{\text{main}'}$ . (We shall also let  $N_h^{\text{trim}}(s, a)$  denote the number of samples that visits  $(s, a, h)$  in  $\mathcal{D}^{\text{trim}}$ .)

3 **Run VI-LCB:** set  $\mathcal{D}_0 = \mathcal{D}^{\text{trim}}$ ; run Algorithm 7 to compute a policy  $\widehat{\pi}$ .

---

---

**Algorithm 7:** Offline value iteration with lower confidence bounds (Li et al., 2022).

---

1 **Input:** dataset  $\mathcal{D}_0$ ; reward function  $r$ ; target success probability  $1 - \delta$ .

2 **Initialization:**  $\widehat{V}_{H+1} = 0$ .

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

4     compute the empirical transition kernel  $\widehat{P}_h$  as

$$\widehat{P}_h(s' | s, a) = \begin{cases} \frac{1}{N_h(s, a)} \sum_{i=1}^N \mathbb{1}\{(s_i, a_i, h_i, s'_i) = (s, a, h, s')\}, & \text{if } N_h(s, a) > 0, \\ \frac{1}{S}, & \text{else,} \end{cases} \quad (81)$$

where  $N_h(s, a) := \sum_{i=1}^N \mathbb{1}\{(s_i, a_i, h_i) = (s, a, h)\}$  and  $N_h(s) := \sum_{i=1}^N \mathbb{1}\{(s_i, h_i) = (s, h)\}$ .

5     **for**  $s \in \mathcal{S}, a \in \mathcal{A}$  **do**

6         compute the penalty term  $b_h(s, a)$  as

$$\forall (s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]: b_h(s, a) = \min \left\{ \sqrt{\frac{c_b \log \frac{K}{\delta}}{N_h(s, a)} \text{Var}_{\widehat{P}_h(\cdot | s, a)}(\widehat{V}_{h+1})} + c_b H \frac{\log \frac{K}{\delta}}{N_h(s, a)}, H \right\}$$

for some universal constant  $c_b > 0$  (e.g.,  $c_b = 16$ ); set

$$\widehat{Q}_h(s, a) = \max \{r_h(s, a) + \widehat{P}_{h, s, a} \widehat{V}_{h+1} - b_h(s, a), 0\}.$$

7     **for**  $s \in \mathcal{S}$  **do**

8         set  $\widehat{V}_h(s) = \max_a \widehat{Q}_h(s, a)$  and  $\widehat{\pi}_h(s) \in \arg \max_a \widehat{Q}_h(s, a)$ .

9 **Output:**  $\widehat{\pi} = \{\widehat{\pi}_h\}_{1 \leq h \leq H}$ .

---

## B Proof for the stopping criterion and the iteration complexity for solving (21b)

**Feasibility of the stopping rule (31).** We first demonstrate that the stopping rule (31) can be satisfied by some mixed policy, namely,

$$\min_{\mu \in \Delta(\Pi)} \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot | s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi \sim \mu}[\widehat{d}_h^\pi(s, a)]} \right] \leq 108SH. \quad (82)$$

Towards this end, we focus attention on analyzing a specific choice of the mixed policy  $\mu^{\text{off}}$  — the one that represents the mixed policy that generates the offline dataset. Making use of the definition (16) of  $\widehat{d}^{\text{off}}$  gives

$$\begin{aligned} \widehat{d}_h^{\text{off}}(s, a) &= \frac{2N_h^{\text{off}}(s, a)}{K^{\text{off}}} \mathbb{1} \left( \frac{N_h^{\text{off}}(s, a)}{K^{\text{off}}} \geq c_{\text{off}} \left\{ \frac{\log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K} \right\} \right) \\ &\leq 3d_h^{\text{off}}(s, a) \mathbb{1} \left( \frac{3}{2}d_h^{\text{off}}(s, a) \geq c_{\text{off}} \left\{ \frac{\log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K} \right\} \right) \\ &\leq 3d_h^{\text{off}}(s, a) \mathbb{1} \left( d_h^{\text{off}}(s, a) \geq \frac{2}{3}c_{\text{off}} \left\{ \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K} \right\} \right), \end{aligned} \quad (83)$$

where the second line relies on (95). This combined with Lemma 1 results in

$$\begin{aligned} &\sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} \pi_h^t(a | s) \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi \sim \mu^{\text{off}}}[\widehat{d}_h^\pi(s, a)]} \\ &\leq \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} \pi_h^t(a | s) \frac{3d_h^{\text{off}}(s, a) \mathbb{1} \left( d_h^{\text{off}}(s, a) \geq \frac{2}{3}c_{\text{off}} \left\{ \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K} \right\} \right)}{\frac{1}{KH} + \frac{1}{2}\mathbb{E}_{\pi \sim \mu^{\text{off}}}[d_h^\pi(s, a) - 2e_h^\pi(s, a) - \frac{\xi}{4N}]}. \end{aligned} \quad (84)$$

Moreover, inequality (47) tells us that: when  $d_h^{\text{off}}(s, a) \geq \frac{2}{3}c_{\text{off}}\left(\frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}}\right)$  for some large enough constant  $c_{\text{off}} > 0$ , we have

$$\mathbb{E}_{\pi \sim \mu^{\text{off}}} \left[ d_h^\pi(s, a) - 2e_h^\pi(s, a) - \frac{\xi}{4N} \right] \geq \mathbb{E}_{\pi \sim \mu^{\text{off}}} \left[ d_h^\pi(s, a) - \frac{4SA}{K^{\text{on}}} - \frac{27c_\xi S^4 A^4 H^4 \log \frac{HSA}{\delta}}{N} \right] \geq \frac{1}{2}d_h^{\text{off}}(s, a). \quad (85)$$

In turn, this implies that

$$\begin{aligned} & \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \pi_h^t(a|s) \frac{3d_h^{\text{off}}(s, a) \mathbb{1}\left(d_h^{\text{off}}(s, a) \geq \frac{2}{3}c_{\text{off}}\left\{\frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K}\right\}\right)}{\frac{1}{KH} + \frac{1}{2}\mathbb{E}_{\pi \sim \mu^{\text{off}}} \left[ d_h^\pi(s, a) - 2e_h^\pi(s, a) - \frac{\xi}{4N} \right]} \\ & \leq \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \pi_h^t(a|s) \frac{3d_h^{\text{off}}(s, a) \mathbb{1}\left(d_h^{\text{off}}(s, a) \geq \frac{2}{3}c_{\text{off}}\left\{\frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K}\right\}\right)}{\frac{1}{KH} + \frac{1}{4}d_h^{\text{off}}(s, a)} \\ & \leq \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \pi_h^t(a|s) \frac{3d_h^{\text{off}}(s, a)}{\frac{1}{4}d_h^{\text{off}}(s, a)} = 12S. \end{aligned} \quad (86)$$

Consequently, combine (84) and (86) to yield

$$\sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi \sim \mu^{\text{off}}} \left[ \widehat{d}_h^\pi(s, a) \right]} \right] \leq 12SH, \quad (87)$$

which clearly validates the claim (82) (with an even better pre-constant).

Before moving forward, we single out one useful property that arises from the above arguments:

$$\mathbb{E}_{\pi \sim \mu^{\text{off}}} \left[ \widehat{d}_h^\pi(s, a) \right] \geq \frac{1}{12} \widehat{d}_h^{\text{off}}(s, a). \quad (88)$$

To prove the validity of this claim (88), it suffices to make the following two observations:

- When  $d_h^{\text{off}}(s, a) \geq \frac{2}{3}c_{\text{off}}\left(\frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}}\right)$ , it has been shown in (84) and (85) in conjunction with (83) that

$$\mathbb{E}_{\pi \sim \mu^{\text{off}}} \left[ \widehat{d}_h^\pi(s, a) \right] \geq \frac{1}{4}d_h^{\text{off}}(s, a) \geq \frac{1}{12} \widehat{d}_h^{\text{off}}(s, a). \quad (89)$$

- When  $d_h^{\text{off}}(s, a) < \frac{2}{3}c_{\text{off}}\left(\frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}}\right)$ , one sees from (83) that  $\widehat{d}_h^{\text{off}}(s, a) = 0$ , and hence (88) holds true trivially.

**Iteration complexity.** Suppose that the stopping criterion (31) is not yet met in the  $k$ -th iteration, namely,

$$\sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}} \left[ \widehat{d}_h^{\pi'}(s, a) \right]} \right] > 108SH. \quad (90)$$

It can be easily seen that

$$\begin{aligned} & \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\left(\frac{1}{KH} + \widehat{d}_h^{\pi^{(k)}}(s, a)\right) \widehat{d}_h^{\text{off}}(s, a)}{\left(\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}} \left[ \widehat{d}_h^{\pi'}(s, a) \right]\right)^2} \right] \\ & \geq \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\left(\frac{1}{KH} + \mathbb{E}_{\pi \sim \mu^{\text{off}}} \left[ \widehat{d}_h^\pi(s, a) \right]\right) \widehat{d}_h^{\text{off}}(s, a)}{\left(\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}} \left[ \widehat{d}_h^{\pi'}(s, a) \right]\right)^2} \right] \\ & \geq \frac{1}{12} \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\left(\widehat{d}_h^{\text{off}}(s, a)\right)^2}{\left(\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}} \left[ \widehat{d}_h^{\pi'}(s, a) \right]\right)^2} \right] \end{aligned}$$

$$\begin{aligned}
&\geq 3 \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]} \right] \mathbf{1} \left( \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]} > 36 \right) \\
&= 3 \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]} \right] \\
&\quad - 3 \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]} \right] \mathbf{1} \left( \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]} \leq 36 \right) \\
&\geq 3 \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]} \right] - 108SH \\
&\geq 2 \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]} \right], \tag{91}
\end{aligned}$$

where the first inequality follows from the choice (26) of  $\pi^{(k)}$ , the second inequality is a consequence of the relation (88), and the last line makes use of (90). This in turn allows one to demonstrate that

$$\begin{aligned}
&\sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{(\widehat{d}_h^{\pi^{(k)}}(s, a) - \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]) \widehat{d}_h^{\text{off}}(s, a)}{\left(\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]\right)^2} \right] \\
&= \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\left(\frac{1}{KH} + \widehat{d}_h^{\pi^{(k)}}(s, a)\right) \widehat{d}_h^{\text{off}}(s, a)}{\left(\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]\right)^2} \right] - \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\left(\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]\right) \widehat{d}_h^{\text{off}}(s, a)}{\left(\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]\right)^2} \right] \\
&\geq \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]} \right] > 108SH,
\end{aligned}$$

where the last line results from (91) and the condition (90). We can then readily derive

$$\begin{aligned}
&\sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]} \right] - \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k+1)}}[\widehat{d}_h^{\pi'}(s, a)]} \right] \\
&= \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{(\mathbb{E}_{\pi' \sim \mu^{(k+1)}}[\widehat{d}_h^{\pi'}(s, a)] - \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]) \widehat{d}_h^{\text{off}}(s, a)}{\left(\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]\right) \left(\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k+1)}}[\widehat{d}_h^{\pi'}(s, a)]\right)} \right] \\
&= \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\alpha (\widehat{d}_h^{\pi^{(k)}}(s, a) - \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]) \widehat{d}_h^{\text{off}}(s, a)}{\left(\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]\right) \left(\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k+1)}}[\widehat{d}_h^{\pi'}(s, a)]\right)} \right] \\
&= \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\alpha (\widehat{d}_h^{\pi^{(k)}}(s, a) - \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]) \widehat{d}_h^{\text{off}}(s, a)}{\left(\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]\right)^2} \right] \\
&\quad - \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\alpha^2 (\widehat{d}_h^{\pi^{(k)}}(s, a) - \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)])^2 \widehat{d}_h^{\text{off}}(s, a)}{\left(\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k)}}[\widehat{d}_h^{\pi'}(s, a)]\right)^2 \left(\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu^{(k+1)}}[\widehat{d}_h^{\pi'}(s, a)]\right)} \right] \\
&\geq 108\alpha SH - \alpha^2 (K^3 H^4) = \frac{108S^2}{K^3 H^2} - \frac{S^2}{K^3 H^2} = \frac{107S^2}{K^3 H^2}, \tag{92}
\end{aligned}$$

where the third line relies on the update rule (29), and the last line utilizes the choice (30) of  $\alpha$ .

In summary, the above argument reveals that: before the stopping criterion is met, each iteration is able to make progress at least as large as in (92). Recognizing the crude bound

$$0 \leq \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} \left[ \frac{\widehat{d}_h^{\text{off}}(s, a)}{\frac{1}{KH} + \mathbb{E}_{\pi' \sim \mu}[\widehat{d}_h^{\pi'}(s, a)]} \right] \leq KH \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mathbb{E}_{a \sim \pi_h^{t+1}(\cdot|s)} [\widehat{d}_h^{\text{off}}(s, a)] \leq KH^2$$

that holds for any  $\mu \in \Delta(\Pi)$ , one can combine this with (92) to conclude that: the proposed procedure terminates within  $O\left(\frac{K^4 H^4}{S^2}\right)$  iterations, as claimed.



## C Proofs of technical lemmas

### C.1 Proof of Lemma 2

The Bernstein inequality combined with the union bound tells us that, with probability at least  $1 - \delta/3$ ,

$$\begin{aligned} \left| \frac{2N_h^{\text{off}}(s, a)}{K^{\text{off}}} - d_h^{\text{off}}(s, a) \right| &\leq 6\sqrt{\frac{d_h^{\text{off}}(s, a) \log \frac{HSA}{\delta}}{K^{\text{off}}}} + \frac{6 \log \frac{HSA}{\delta}}{K^{\text{off}}} \\ &\leq \frac{d_h^{\text{off}}(s, a)}{2} + \frac{24 \log \frac{HSA}{\delta}}{K^{\text{off}}} \end{aligned} \quad (93)$$

holds simultaneously for all  $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times [H]$ , where the last line invokes the AM-GM inequality. This in turn reveals that

$$\frac{4N_h^{\text{off}}(s, a)}{3K^{\text{off}}} - \frac{16 \log \frac{HSA}{\delta}}{K^{\text{off}}} \leq d_h^{\text{off}}(s, a) \leq \frac{4N_h^{\text{off}}(s, a)}{K^{\text{off}}} + \frac{48 \log \frac{HSA}{\delta}}{K^{\text{off}}}. \quad (94)$$

As a result, we can show that:

- If  $\frac{N_h^{\text{off}}(s, a)}{K^{\text{off}}} \geq c_{\text{off}} \left( \frac{\log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}} \right)$  for some  $c_{\text{off}} \geq 48$ , then one has

$$\frac{2N_h^{\text{off}}(s, a)}{3K^{\text{off}}} \leq d_h^{\text{off}}(s, a) \leq \frac{6N_h^{\text{off}}(s, a)}{K^{\text{off}}} \quad \text{and} \quad \widehat{d}_h^{\text{off}}(s, a) = \frac{2N_h^{\text{off}}(s, a)}{K^{\text{off}}} \quad (95a)$$

$$\implies \frac{1}{3} \widehat{d}_h^{\text{off}}(s, a) \leq d_h^{\text{off}}(s, a) \leq 3 \widehat{d}_h^{\text{off}}(s, a). \quad (95b)$$

- If instead  $\frac{N_h^{\text{off}}(s, a)}{K^{\text{off}}} < c_{\text{off}} \left( \frac{\log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}} \right)$ , then one has  $\widehat{d}_h^{\text{off}}(s, a) = 0$ , and therefore,

$$\begin{aligned} d_h^{\text{off}}(s, a) &\geq 0 = \frac{1}{3} \widehat{d}_h^{\text{off}}(s, a), \\ d_h^{\text{off}}(s, a) &\leq \frac{4N_h^{\text{off}}(s, a)}{K^{\text{off}}} + \frac{48 \log \frac{HSA}{\delta}}{K^{\text{off}}} < 5c_{\text{off}} \left\{ \frac{\log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}} \right\} \\ &= \widehat{d}_h^{\text{off}}(s, a) + 5c_{\text{off}} \left\{ \frac{\log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}} \right\}. \end{aligned}$$

Taken collectively, these inequalities demonstrate that

$$\frac{1}{3} \widehat{d}_h^{\text{off}}(s, a) \leq d_h^{\text{off}}(s, a) \leq \widehat{d}_h^{\text{off}}(s, a) + 5c_{\text{off}} \left\{ \frac{\log \frac{HSA}{\delta}}{K^{\text{off}}} + \frac{H^4 S^4 A^4 \log \frac{HSA}{\delta}}{N} + \frac{SA}{K^{\text{on}}} \right\}, \quad (96)$$

provided that  $c_{\text{off}} > 0$  is sufficiently large.

### C.2 Proof of Lemma 6

Before embarking on the proof, let us introduce several notation. For each  $1 \leq h \leq H$ , define  $P_h^{\pi^*} \in \mathbb{R}^{S \times S}$  and  $d_h^{\pi^*} \in \mathbb{R}^S$  such that: for all  $s \in \mathcal{S}$ ,

$$P_h^{\pi^*}(s, \cdot) = P_h(\cdot | s, \pi^*(s)) \quad \text{and} \quad d_h^{\pi^*}(s) = d_h^{\pi^*}(s, \pi^*(s)). \quad (97)$$

To begin with, it can be easily seen from (60) and the basic fact  $\text{Var}_{P_h(\cdot | s, a)}(\widehat{V}_{h+1}) \leq H^2$  that

$$\langle d_j^{\pi^*}, V_j^* - V_j^{\widehat{\pi}} \rangle \leq 8H \underbrace{\sum_{h: h \geq j} \sum_s \max_{a: (s, a) \in \mathcal{I}_h} \sqrt{2d_h^*(s, a) C^*(\sigma) \widehat{d}_h^{\text{off}}(s, a)}}_{=: \gamma_3} \sqrt{\frac{c_b \log \frac{K}{\delta}}{N_h^{\text{trim}}(s, a) + 1}}$$

$$+ 4H \underbrace{\sum_{h:h \geq j} \sum_{(s,a) \in \mathcal{T}_h} \widehat{d}_h^{\pi^*}(s,a)}_{=: \gamma_2} \sqrt{\frac{c_b \log \frac{K}{\delta}}{N_h^{\text{trim}}(s,a) + 1}} + \frac{61c_\xi H^6 S^4 A^4}{N} \log \frac{K}{\delta} \quad (98)$$

holds for any  $j \in [H]$ . Note that we have bounded  $\gamma_2$  in (70). We then need to bound  $\gamma_3$ .

With regards to the term  $\gamma_3$ : invoking similar arguments as for (63) leads to

$$\begin{aligned} \gamma_3 &\leq 64c_b H \sum_{h:h \geq j} \sum_s \max_a \sqrt{2d_h^{\pi^*}(s,a) C^*(\sigma) \widehat{d}_h^{\text{off}}(s,a)} \sqrt{\frac{\log \frac{K}{\delta}}{1/H + \frac{1}{72} K^{\text{on}} \mathbb{E}_{\pi \in \mu^{\text{imitate}}} [\widehat{d}_h^{\pi}(s,a)]}} + 4H^2 S C^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{K^{\text{on}} H} \right) \\ &\leq 64c_b H \sqrt{2 \sum_{h:h \geq j} \sum_{s,a} d_h^{\pi^*}(s,a)} \cdot \sqrt{\sum_{h:h \geq j} \sum_s \max_a \frac{C^*(\sigma) \widehat{d}_h^{\text{off}}(s,a) \log \frac{K}{\delta}}{1/H + \frac{1}{72} K^{\text{on}} \mathbb{E}_{\pi \in \mu^{\text{imitate}}} [\widehat{d}_h^{\pi}(s,a)]}} + 4H^2 S C^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{K^{\text{on}} H} \right) \\ &\leq 768c_b \sqrt{\frac{13H^4 S C^*(\sigma)}{K^{\text{on}}} \log \frac{K}{\delta}} + 4H^2 S C^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{K^{\text{on}} H} \right), \end{aligned} \quad (99)$$

where the second step invokes the Cauchy-Schwartz inequality, and the last line comes from (58). Similarly, repeating the above argument but focusing on the offline dataset, we can derive (which we omit for the sake of brevity)

$$\begin{aligned} \gamma_3 &\leq 64c_b H \sum_{h:h \geq j} \sum_s \max_a \sqrt{2d_h^{\pi^*}(s,a) C^*(\sigma) \widehat{d}_h^{\text{off}}(s,a)} \sqrt{\frac{\log \frac{K}{\delta}}{1/H + \frac{1}{72} K^{\text{off}} d_h^{\text{off}}(s,a)}} + 4H^2 S C^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{K^{\text{off}} H} \right) \\ &\leq 64c_b H \sqrt{2 \sum_{h:h \geq j} \sum_{s,a} d_h^{\pi^*}(s,a)} \cdot \sqrt{\sum_{h:h \geq j} \sum_s \max_a \frac{C^*(\sigma) \widehat{d}_h^{\text{off}}(s,a) \log \frac{K}{\delta}}{1/H + \frac{1}{72} K^{\text{off}} d_h^{\text{off}}(s,a)}} + 4H^2 S C^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{K^{\text{off}} H} \right) \\ &\leq 768c_b \sqrt{\frac{6H^4 S C^*(\sigma)}{K^{\text{off}}} \log \frac{K}{\delta}} + 4H^2 S C^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{K^{\text{off}} H} \right), \end{aligned} \quad (100)$$

where the last line makes use of (87).

Combining (99), (100) and (70) with (98), we can show that

$$\begin{aligned} \langle d_j^{\pi^*}, V_j^* - V_j^{\widehat{\pi}} \rangle &\leq \min \left\{ 768c_b \sqrt{\frac{13H^4 S C^*(\sigma)}{K^{\text{on}}} \log \frac{K}{\delta}} + 4H^2 S C^*(\sigma) \left( \frac{\xi}{N} + \frac{1}{KH} \right), 768c_b \sqrt{\frac{6H^4 S C^*(\sigma)}{K^{\text{off}}} \log \frac{K}{\delta}} \right\} \\ &\quad + 96c_b \sqrt{\frac{2H^3 S A (2\widehat{\sigma} + H S \xi / N)}{K^{\text{on}}} \log \frac{HK}{\delta}} + 436H^2 S A \left( \frac{\xi}{N} + \frac{1}{\min\{K^{\text{on}}, K^{\text{off}}\} H} \right) \\ &\quad + \frac{61c_\xi H^6 S^4 A^4}{N} \log \frac{K}{\delta} \\ &\leq 1 \end{aligned} \quad (101)$$

for all  $1 \leq j \leq H$ , with the proviso that

$$\begin{aligned} \frac{H^7 S^5 A^4}{K^{\text{on}}} \log^2 \frac{K}{\delta} &\leq c_{10} \\ \frac{H^5 S^4 A^3 C^*(\sigma) \log \frac{K}{\delta}}{K^{\text{on}}} &\leq c_{10} \\ \frac{H S C^*(\sigma) \log \frac{K}{\delta}}{K^{\text{off}}} &\leq c_{10} \\ \frac{H S A \log \frac{K}{\delta}}{K^{\text{off}}} &\leq c_{10} \end{aligned}$$

for some sufficiently small constant  $c_{10} > 0$ . As a consequence, we can demonstrate that

$$\begin{aligned}
& \sum_{h:h \geq j} \sum_{s,a} d_h^{\pi^*}(s,a) \text{Var}_{P_h(\cdot|s,a)}(\widehat{V}_{h+1}) \\
& \leq \sum_{h:h \geq j} 2 \sum_{s,a} d_h^{\pi^*}(s,a) \left( \text{Var}_{P_h(\cdot|s,a)}(V_{h+1}^*) + \text{Var}_{P_h(\cdot|s,a)}(V_{h+1}^* - \widehat{V}_{h+1}) \right) \\
& \leq 4H^2 + H \sum_{h:h \geq j} \sum_s d_h^{\pi^*}(s, \pi^*(s)) \mathbb{E}_{P_h(\cdot|s, \pi^*(s))} [V_{h+1}^* - \widehat{V}_{h+1}] \\
& = 4H^2 + H \sum_{h:h \geq j} (d_h^{\pi^*})^\top P_h^{\pi^*} [V_{h+1}^* - \widehat{V}_{h+1}] \\
& = 4H^2 + H \sum_{h \geq j} (d_{h+1}^{\pi^*})^\top [V_{h+1}^* - \widehat{V}_{h+1}] \leq 5H^2
\end{aligned} \tag{102}$$

for all  $1 \leq j \leq H$ . Here, the third line in (102) applies the following fact

$$\begin{aligned}
& \sum_h \sum_{s,a} d_h^{\pi^*}(s,a) \text{Var}_{P_h(\cdot|s,a)}(V_{h+1}^*) \\
& = \sum_h \sum_s d_h^{\pi^*}(s, \pi^*(s)) \left[ \langle P_h^{\pi^*}(s, \cdot), V_{h+1}^* \circ V_{h+1}^* \rangle - (\langle P_h^{\pi^*}(s, \cdot), V_{h+1}^* \rangle)^2 \right] \\
& = \sum_h \sum_s d_h^{\pi^*}(s, \pi^*(s)) \left[ \langle P_h^{\pi^*}(s, \cdot), V_{h+1}^* \circ V_{h+1}^* \rangle - (V_h^*(s) - r_h^*(s, \pi^*(s)))^2 \right] \\
& \leq \sum_h \sum_s d_h^{\pi^*}(s, \pi^*(s)) \left[ \langle P_h^{\pi^*}(s, \cdot), V_{h+1}^* \circ V_{h+1}^* \rangle - (V_h^*(s))^2 + 2H \right] \\
& = \sum_h (d_h^{\pi^*})^\top P_h^{\pi^*} (V_{h+1}^* \circ V_{h+1}^*) - \sum_h (d_h^{\pi^*})^\top (V_h^* \circ V_h^*) + 2H^2 \\
& = \sum_h (d_{h+1}^{\pi^*})^\top (V_{h+1}^* \circ V_{h+1}^*) - \sum_h (d_h^{\pi^*})^\top (V_h^* \circ V_h^*) + 2H^2 \\
& \leq 2H^2,
\end{aligned}$$

where the second identity comes from the Bellman equation; the third relation uses the fact that  $V_h^*(s) \leq H$ , and the penultimate line holds since  $(d_h^{\pi^*})^\top P_h^{\pi^*} = (d_{h+1}^{\pi^*})^\top$ ; the penultimate step in (102) is due to (101); and the line in (102) holds true since  $(d_h^{\pi^*})^\top P_h^{\pi^*} = (d_{h+1}^{\pi^*})^\top$ .