
Policy Finetuning in Reinforcement Learning via Design of Experiments using Offline Data

Anonymous Author(s)
Affiliation
Address
email

Abstract

1 In some applications of reinforcement learning, a dataset of pre-collected experi-
2 ence is already available but it is also possible to acquire some additional online
3 data to help improve the quality of the policy. However, it may be preferable to
4 gather additional data with a single, non-reactive exploration policy and avoid the
5 engineering costs associated with switching policies.

6 In this paper we propose an algorithm with provable guarantees that can leverage
7 an offline dataset to design a single non-reactive policy for exploration. We the-
8oretically analyze the algorithm and measure the quality of the final policy as a
9 function of the local coverage of the original dataset and the amount of additional
10 data collected.

11 1 Introduction

12 Reinforcement learning (RL) is a general framework for data-driven, sequential decision making
13 [Puterman, 1994, Sutton and Barto, 2018]. In RL, a common goal is to identify a near-optimal policy,
14 and there exist two main paradigms: *online* and *offline* RL.

15 Online RL is effective when the practical cost of a bad decision is low, such as in simulated environ-
16 ments (e.g., [Mnih et al., 2015, Silver et al., 2016]). In online RL, a well designed learning algorithm
17 starts from tabula rasa and implements a sequence of policies with a value that should approach that
18 of an optimal policy. When the cost of making a mistake is high, such as in healthcare [Gottesman
19 et al., 2018] and in self-driving [Kiran et al., 2021], an offline approach is preferred. In offline RL,
20 the agent uses a dataset of pre-collected experience to extract a policy that is as good as possible. In
21 this latter case, the quality of the policy that can be extracted from the dataset is limited by the quality
22 of the dataset.

23 Many applications, however, fall between these two opposite settings: for example, a company that
24 sells products online has most likely recorded the feedback that it has received from its customers, but
25 can also collect a small amount of additional strategic data in order to improve its recommendation
26 engine. While in principle an online exploration algorithm can be used to collect fresh data, in
27 practice there are a number of practical engineering considerations that require the policy to be
28 deployed to be **non-reactive**. We say that a policy is non-reactive, (or passive, memoryless) if it
29 chooses actions only according to the current state of the system. Most online algorithms are, by
30 design, reactive to the data being acquired.

31 An example of a situation where non-reactive policies may be preferred are those where a human
32 in the loop is required to validate each exploratory policy before they are deployed, to ensure they
33 are of high quality [Dann et al., 2019] and safe [Yang et al., 2021], as well as free of discriminatory
34 content [Koencke et al., 2020]. Other situations that may warrant non-reactive exploration are those
35 where the interaction with the user occurs through a distributed system with delayed feedback. In

36 recommendation systems, data collection may only take minutes, but policy deployment and updates
37 can span weeks [Afsar et al., 2022]. Similar considerations apply across various RL application
38 domains, including healthcare [Yu et al., 2021], computer networks [Xu et al., 2018], and new material
39 design [Raccuglia et al., 2016]. In all such cases, the engineering effort required to implement a
40 system that handles real-time policy switches may be prohibitive: deploying a single, non-reactive
41 policy is much preferred.

42 **Non-reactive exploration from offline data** Most exploration algorithms that we are aware of
43 incorporate policy switches when they interact with the environment [Dann and Brunskill, 2015,
44 Dann et al., 2017, Azar et al., 2017, Jin et al., 2018, Dann et al., 2019, Zanette and Brunskill, 2019,
45 Zhang et al., 2020b]. Implementing a sequence of non-reactive policies is necessary in order to achieve
46 near-optimal regret: the number of policy switches must be at least $\tilde{O}(H|\mathcal{S}||\mathcal{A}|\log\log K)$ where
47 $\mathcal{S}, \mathcal{A}, H, K$ are the state space, action space, horizon and the total number of episodes, respectively
48 [Qiao et al., 2022]. With no switches, i.e., when a fully non-reactive data collection strategy is
49 implemented, it is information theoretically impossible [Xiao et al., 2022] to identify a good policy
50 using a number of samples polynomial in the size of the state and action space.

51 However, these fundamental limits apply to the case where the agent learns from tabula rasa. In the
52 more common case where offline data is available, we demonstrate that it is possible to leverage the
53 dataset to design an effective non-reactive exploratory policy. More precisely, an available offline
54 dataset contains information (e.g., transitions) about a certain area of the state-action space, a concept
55 known as *partial coverage*. A dataset with partial coverage naturally identifies a ‘sub-region’ of the
56 original MDP—more precisely, a sub-graph—that is relatively well explored. We demonstrate that it
57 is possible to use the dataset to design a non-reactive policy that further explores such sub-region.
58 The additional data collected can be used to learn a near-optimal policy in such sub-region.

59 In other words, exploration with no policy switches can collect additional information and compete
60 with the best policy that is restricted to an area where the original dataset has sufficient information.
61 The value of such policy can be much higher than the one that can be computed using only the offline
62 dataset, and does not directly depend on a concentrability coefficient [Munos and Szepesvári, 2008,
63 Chen and Jiang, 2019].

64 Perhaps surprisingly, addressing the problem of reactive exploration in reinforcement learning
65 requires an approach that *combines both optimism and pessimism* in the face of uncertainty to
66 explore efficiently. While optimism drives exploration, pessimism ensures that the agent explores
67 conservatively, in a way that restricts its exploration effort to a region that it knows how to navigate,
68 and so our paper makes a technical contribution which can be of independent interest.

69 **Contributions** To the best of our knowledge, this is the first paper with theoretical rigor that considers
70 the problem of designing an experiment in reinforcement learning for online, passive exploration,
71 using a dataset of pre-collected experience. More precisely, our contributions are as follows:

- 72 • We introduce an algorithm that takes as input a dataset, uses it to design and deploy a non-reactive
73 exploratory policy, and then outputs a locally near-optimal policy.
- 74 • We introduce the concept of sparsified MDP, which is actively used by our algorithm to design the
75 exploratory policy, as well as to theoretically analyze the quality of the final policy that it finds.
- 76 • We rigorously establish a nearly minimax-optimal upper bound for the sample complexity needed
77 to learn a local ε -optimal policy using our algorithm.¹

78 2 Related Work

79 In this section we discuss some related literature. Our work is related to low-switching algorithms, but
80 unlike those, we focus on the limit case where *no-switches* are allowed. For more related work about
81 low-switching algorithms, offline RL, task-agnostic RL, and reward-free RL we refer to Appendix F.

82 **Low-switching RL** In reinforcement learning, [Bai et al., 2019] first proposed Q-learning with UCB2
83 exploration, proving an $O(H^3|\mathcal{S}||\mathcal{A}|\log K)$ switching cost. This was later improved by a factor of

¹More rigorously, our sample complexity matches the minimax lower bound when we have some degree of knowledge for the full MDP, see the discussion in Section section 5. The lower bound for the samples needed in reward-free case is proved in [Jin et al., 2020b], but their result applies to the non-homogeneous MDP. The sample complexity on homogeneous MDP should be shaved off by an H factor.

84 H by the UCBadvantage algorithm in [Zhang et al., 2020b]. Recently, [Qiao et al., 2022] generalized
 85 the policy elimination algorithm from [Cesa-Bianchi et al., 2013] and introduced APEVE, which
 86 attains an optimal $O(H |\mathcal{S}| |\mathcal{A}| \log \log K)$ switching cost. The reward-free version of their algorithm
 87 (which is not regret minimizing) has an $O(H |\mathcal{S}| |\mathcal{A}|)$ switching cost.

88 Similar ideas were soon applied in RL with linear function approximation [Gao et al., 2021, Wang
 89 et al., 2021, Qiao and Wang, 2022] and general function approximation [Qiao et al., 2023]. Addition-
 90 ally, numerous research efforts have focused on low-adaptivity in other learning domains, such as
 91 batched dueling bandits [Agarwal et al., 2022], batched convex optimization [Duchi et al., 2018],
 92 linear contextual bandits [Ruan et al., 2021], and deployment-efficient RL [Huang et al., 2022].

93 Our work was inspired by the problem of non-reactive policy design in linear contextual bandits.
 94 Given access to an offline dataset, [Zanette et al., 2021a] proposed an algorithm to output a single
 95 exploratory policy, which generates a dataset from which a near-optimal policy can be extracted.
 96 However, there are a number of additional challenges which arise in reinforcement learning, including
 97 the fact that the state space is only partially explored in the offline dataset. In fact, in reinforcement
 98 learning, [Xiao et al., 2022] established an exponential lower bound for any non-adaptive policy
 99 learning algorithm starting from tabula rasa.

100 3 Setup

101 Throughout this paper, we let $[n] = \{1, 2, \dots, n\}$. We adopt the big-O notation, where $\tilde{O}(\cdot)$ suppresses
 102 poly-log factors of the input parameters. We indicate the cardinality of a set \mathcal{X} with $|\mathcal{X}|$.

103 **Markov decision process** We consider time-homogeneous episodic Markov decision processes
 104 (MDPs). They are defined by a finite state space \mathcal{S} , a finite action space \mathcal{A} , a transition kernel \mathbb{P} , a
 105 reward function r and the episodic length H . The transition probability $\mathbb{P}(s' | s, a)$, which does not
 106 depend on the current time-step $h \in [H]$, denotes the probability of transitioning to state $s' \in \mathcal{S}$
 107 when taking action $a \in \mathcal{A}$ in the current state $s \in \mathcal{S}$. Typically we denote with s_1 the initial state.
 108 For simplicity, we consider deterministic reward functions $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$. A deterministic
 109 non-reactive (or memoryless, or passive) policy $\pi = \{\pi_h\}_{h \in [H]}$ maps a given state to an action.

110 The value function is defined as the expected cumulated reward. It depends on the state s under
 111 consideration, the transition \mathbb{P} and reward r that define the MDP as well as on the policy π being
 112 implemented. It is defined as $V_h(s; \mathbb{P}, r, \pi) = \mathbb{E}_{\mathbb{P}, \pi}[\sum_{i=h}^H r(s_i, a_i) | s_h = s]$, where $\mathbb{E}_{\mathbb{P}, \pi}$ denotes
 113 the expectation generated by \mathbb{P} and policy π . A closely related quantity is the state-action value
 114 function, or Q -function, defined as $Q_h(s, a; \mathbb{P}, r, \pi) = \mathbb{E}_{\mathbb{P}, \pi}[\sum_{i=h}^H r(s_i, a_i) | s_h = s, a_h = a]$.
 115 When it is clear from the context, we sometimes omit (\mathbb{P}, r) and simply write them as $V_h^\pi(s)$ and
 116 $Q_h^\pi(s, a)$. We denote an MDP defined by \mathcal{S}, \mathcal{A} and the transition matrix \mathbb{P} as $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathbb{P})$.

117 3.1 Interaction protocol

Algorithm 1 Design of experiments in reinforcement learning

Input: Offline dataset \mathcal{D}

1: *Offline phase:* use \mathcal{D} to compute the exploratory policy π_{ex}

2: *Online phase:* deploy π_{ex} to collect the online dataset \mathcal{D}'

3: *Planning phase:* receive the reward function r and use $\mathcal{D} \cup \mathcal{D}'$ to extract π_{final}

Output: Return π_{final}

118 In this paper we assume access to an *offline dataset* $\mathcal{D} = \{(s, a, s')\}$ where every state-action (s, a)
 119 is sampled in an i.i.d. fashion from some distribution μ and $s' \sim \mathbb{P}(\cdot | s, a)$, which is common in the
 120 offline RL literature [Xie et al., 2021a, Zhan et al., 2022, Rashidinejad et al., 2021, Uehara and Sun,
 121 2021]. We denote $N(s, a)$ and $N(s, a, s')$ as the number of (s, a) and (s, a, s') samples in the offline
 122 dataset \mathcal{D} , respectively. The interaction protocol considered in this paper consists of three distinct
 123 phases, which are displayed in algorithm 1. They are:

- 124 • the **offline phase**, where the learner uses an *offline dataset* \mathcal{D} of pre-collected experience to design
 125 the non-reactive exploratory policy π_{ex} ;

- 126 • the **online phase** where π_{ex} is deployed to generate the *online dataset* \mathcal{D}' ;
- 127 • the **planning phase** where the learner receives a reward function and uses all the data collected to
- 128 extract a good policy π_{final} with respect to that reward function.

129 The objective is to minimize the number of online episodic interactions needed to find a policy π_{final}
 130 whose value is as high as possible. Moreover, we focus on the reward-free RL setting [Jin et al.,
 131 2020a, Kaufmann et al., 2021, Li et al., 2023b], which is more general than reward-aware RL. In the
 132 offline and online phase, the data are generated without specific reward signals, and the entire reward
 133 information is then given in the planning phase. One of the primary advantages of using reward-free
 134 offline data is that it allows for the collection of data without the need for explicit reward signals.
 135 This can be particularly beneficial in environments where obtaining reward signals is costly, risky,
 136 ethically challenging, or where the reward functions are human-designed.

137 4 Algorithm: balancing optimism and pessimism for experimental design

138 In this section we outline our algorithm *Reward-Free Non-reactive Policy Design* (RF-NPD), which
 139 follows the high-level protocol described in algorithm 1. The technical novelty lies almost entirely
 140 in the design of the exploratory policy π_{ex} . In order to prepare the reader for the discussion of the
 141 algorithm, we first give some intuition in section 4.1 followed by the definition of sparsified MDP
 142 in section 4.2, a central concept of this paper, and then describe the implementation of line 1 in the
 143 protocol in algorithm 1 in section 4.3. We conclude by presenting the implementation of lines 2 and 3
 144 in the protocol in algorithm 1.

145 4.1 Intuition

146 In order to present the main intuition for this paper, in this section we assume that enough transitions
 147 are available in the dataset for every edge $(s, a) \rightarrow s'$, namely that the *critical condition*

$$N(s, a, s') \geq \Phi = \tilde{\Theta}(H^2) \quad (4.1)$$

148 holds for all tuples $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$ (the precise value for Φ will be given later in eq. (5.1)).
 149 Such condition is hardly satisfied everywhere in the state-action-state space, but assuming it in this
 150 section simplifies the presentation of one of the key ideas of this paper.

151 The key observation is that when eq. (4.1) holds for all (s, a, s') , we can use the empirical transition
 152 kernel to design an exploration policy π_{ex} to eventually extract a near-optimal policy π_{final} for any
 153 desired level of sub-optimality ε , despite eq. (4.1) being independent of ε . More precisely, let $\hat{\mathbb{P}}$ be
 154 the empirical transition kernel defined in the usual way $\hat{\mathbb{P}}(s' | s, a) = N(s, a, s')/N(s, a)$ for any
 155 tuple (s, a, s') . The intuition—which will be verified rigorously in the analysis of the algorithm—is
 156 the following:

157 *If eq. (4.1) holds for every (s, a, s') then $\hat{\mathbb{P}}$ can be used to design a non-reactive exploration policy*
 158 *π_{ex} which can be deployed on \mathcal{M} to find an ε -optimal policy π_{final} using $\asymp \frac{1}{\varepsilon^2}$ samples.*

159 We remark that even if the condition 4.1 holds for all tuples (s, a, s') , the empirical kernel $\hat{\mathbb{P}}$ is
 160 not accurate enough to extract an ε -optimal policy from the dataset \mathcal{D} without collecting further
 161 data. Indeed, the threshold $\Phi = \tilde{\Theta}(H^2)$ on the number of samples is independent of the desired
 162 sub-optimality $\varepsilon > 0$, while it is well known that at least $\sim \frac{1}{\varepsilon^2}$ offline samples are needed to find
 163 an ε -optimal policy. Therefore, directly implementing an offline RL algorithm to use the available
 164 offline dataset \mathcal{D} does not yield an ε -optimal policy. However, the threshold $\Phi = \tilde{\Theta}(H^2)$ is sufficient
 165 to *design* a non-reactive exploratory policy π_{ex} that can discover an ε -optimal policy π_{final} after
 166 collecting $\sim \frac{1}{\varepsilon^2}$ online data.

167 4.2 Sparsified MDP

168 The intuition in the prior section must be modified to work with heterogeneous datasets and dynamics
 169 where $N(s, a, s') \geq \Phi$ may fail to hold everywhere. For example, if $\mathbb{P}(s' | s, a)$ is very small
 170 for a certain tuple (s, a, s') , it is unlikely that the dataset contains $N(s, a, s') \geq \Phi$ samples for

171 that particular tuple. In a more extreme setting, if the dataset is empty, the critical condition in
 172 eq. (4.1) is violated for all tuples (s, a, s') , and in fact the lower bound of Xiao et al. [2022] states that
 173 finding ε -optimal policies by exploring with a non-reactive policy is not feasible with $\sim \frac{1}{\varepsilon^2}$ sample
 174 complexity. This suggests that in general it is not possible to output an ε -optimal policy using the
 175 protocol in algorithm 1.

176 However, a real-world dataset generally covers at least a portion of the state-action space, and so we
 177 expect the condition $N(s, a, s') \geq \Phi$ to hold somewhere; the sub-region of the MDP where it holds
 178 represents the connectivity graph of the *sparsified MDP*. This is the region that the agent knows how
 179 to navigate using the offline dataset \mathcal{D} , and so it is the one that the agent can explore further using π_{ex} .
 180 More precisely, the sparsified MDP is defined to have identical dynamics as the original MDP on the
 181 edges $(s, a) \rightarrow s'$ that satisfy the critical condition 4.1. When instead the edge $(s, a) \rightarrow s'$ fails to
 182 satisfy the critical condition 4.1, it is replaced with a transition $(s, a) \rightarrow s^\dagger$ to an absorbing state s^\dagger .

183 **Definition 4.1** (Sparsified MDP). *Let s^\dagger be an absorbing state, i.e., such that $\mathbb{P}(s^\dagger | s^\dagger, a) = 1$ and
 184 $r(s^\dagger, a) = 0$ for all $a \in \mathcal{A}$. The state space in the sparsified MDP \mathcal{M}^\dagger is defined as that of the
 185 original MDP with the addition of s^\dagger . The dynamics \mathbb{P}^\dagger of the sparsified MDP are defined as*

$$\mathbb{P}^\dagger(s' | s, a) = \begin{cases} \mathbb{P}(s' | s, a) & \text{if } N(s, a, s') \geq \Phi \\ 0 & \text{if } N(s, a, s') < \Phi, \end{cases} \quad \mathbb{P}^\dagger(s^\dagger | s, a) = \sum_{\substack{s' \neq s^\dagger \\ N(s, a, s') < \Phi}} \mathbb{P}(s' | s, a). \quad (4.2)$$

186 For any deterministic reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, the reward function on the sparsified MDP
 187 is defined as $r^\dagger(s, a) = r(s, a)$; for simplicity we only consider deterministic reward functions.

188 The *empirical sparsified MDP* $\widehat{\mathcal{M}}^\dagger = (\mathcal{S} \cup \{s^\dagger\}, \mathcal{A}, \widehat{\mathbb{P}}^\dagger)$ is defined in the same way but by using
 189 the empirical transition kernel in eq. (4.2). The empirical sparsified MDP is used by our algorithm
 190 to design the exploratory policy, while the (population) sparsified MDP is used for its theoretical
 191 analysis. They are two fundamental concepts in this paper. By formulating the sparsified MDP, we
 192 restrict the transitions and rewards within the area where we know how to navigate, embodying the
 193 principle of pessimism. Various forms of pessimistic regularization have been introduced to address
 194 the challenges of partially covered offline data. Examples include a pessimistic MDP [Kidambi et al.,
 195 2020] and limiting policies to those covered by offline data [Liu et al., 2020].

196 4.3 Offline design of experiments

Algorithm 2 RF-UCB ($\widehat{\mathcal{M}}^\dagger, K_{ucb}, \varepsilon, \delta$)

Input: $\delta \in (0, 1), \varepsilon > 0$, number of episode K_{ucb} , MDP $\widehat{\mathcal{M}}^\dagger$.

- 1: **Initialize** Counter $n^1(s, a) = n^1(s, a, s') = 0$ for any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$.
- 2: **for** $k = 1, 2, \dots, K_{ucb}$ **do**
- 3: **for** $h = H, H - 1, \dots, 1$ **do**
- 4: Set $U_h^k(s, a) = 0$ for any $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$.
- 5: **for** $(s, a) \in \mathcal{S} \times \mathcal{A}$ **do**
- 6: Calculate the empirical uncertainty $U_h^k(s, a)$ using eq. (4.4) where ϕ is from eq. (4.3)
- 7: **end for**
- 8: $\pi_h^k(s) := \arg \max_{a \in \mathcal{A}} U_h^k(s, a), \forall s \in \mathcal{S}$ and $\pi_h^k(s^\dagger) :=$ any action.
- 9: **end for**
- 10: Set initial state $s_1^k = s_1$.
- 11: **for** $h = 1, 2, \dots, H$ **do** Sample $a_h^k \sim \pi_h^k(s_h^k), s_{h+1}^k \sim \widehat{\mathbb{P}}^\dagger(s_h^k, a_h^k)$.
- 12: **end for**
- 13: $n^{k+1}(s, a) = n^k(s, a) + \sum_{h \in [H]} \mathbb{I}[(s, a) = (s_h^k, a_h^k)]$.
- 14: $n^{k+1}(s, a, s') = n^k(s, a, s') + \sum_{h \in [H]} \mathbb{I}[(s, a, s') = (s_h^k, a_h^k, s_{h+1}^k)]$.
- 15: **end for**

Output: $\pi_{ex} = \text{Uniform}\{\pi^k\}_{k \in [K_{ucb}]}$.

197 In this section we describe the main sub-component of the algorithm, namely the sub-routine that
 198 uses the offline dataset \mathcal{D} to compute the exploratory policy π_{ex} . The exploratory policy π_{ex} is a

199 mixture of the policies π^1, π^2, \dots produced by a variant of the reward-free exploration algorithm of
 200 [Kaufmann et al., 2021, Ménard et al., 2021]. Unlike prior literature, the reward-free algorithm is not
 201 interfaced with the real MDP \mathcal{M} , but rather *simulated on the empirical sparsified MDP* $\widehat{\mathcal{M}}^\dagger$. This
 202 avoids interacting with \mathcal{M} with a reactive policy, but it introduces some bias that must be controlled.
 203 The overall procedure is detailed in algorithm 2. To be clear, no real-world samples are collected by
 204 algorithm 2; instead we use the word ‘virtual samples’ to refer to those generated from $\widehat{\mathcal{M}}^\dagger$.

205 At a high level, algorithm 2 implements value iteration using the empirical transition kernel $\widehat{\mathbb{P}}^\dagger$,
 206 with the exploration bonus defined in eq. (4.3) that replaces the reward function. The exploration
 207 bonus can be seen as implementing the principle of optimism in the face of uncertainty; however, the
 208 possibility of transitioning to an absorbing state with zero reward (due to the use of the absorbing
 209 state in the definition of $\widehat{\mathbb{P}}^\dagger$) implements the principle of pessimism.

210 This delicate *interplay between optimism and pessimism is critical* to the success of the overall
 211 procedure: while optimism encourages exploration, pessimism ensures that the exploration efforts
 212 are directed to the region of the state space that the agent actually knows how to navigate, and
 213 prevents the agent from getting ‘trapped’ in unknown regions. In fact, these latter regions could have
 214 combinatorial structures [Xiao et al., 2022] which cannot be explored with non-reactive policies.

215 More precisely, at the beginning of the k -th virtual episode in algorithm 2, $n^k(s, a)$ and $n^k(s, a, s')$
 216 denote the counters for the number of virtual samples simulated from $\widehat{\mathcal{M}}^\dagger$ at each (s, a) and (s, a, s')
 217 tuple. We define the bonus function

$$\phi(x, \delta) = \frac{H}{x} [\log(6H |\mathcal{S}| |\mathcal{A}| / \delta) + |\mathcal{S}| \log(e(1 + x / |\mathcal{S}|))], \quad (4.3)$$

218 which is used to construct the *empirical uncertainty function* U_h^k , a quantity that serves as a proxy for
 219 the uncertainty of the value of any policy π on the sparsified MDP. Specifically, for the k -th virtual
 220 episode, we set $U_{H+1}^k(s, a) = 0$ and $s \in \mathcal{S}, a \in \mathcal{A}$. For $h \in [H]$, we further define:

$$U_h^k(s, a) = H \min\{1, \phi(n^k(s, a))\} + \widehat{\mathbb{P}}^\dagger(s, a)^\top (\max_{a'} U_{h+1}^k(\cdot, a')). \quad (4.4)$$

221 Note that, the above bonus function takes a similar form of the bonus function in [Ménard et al.,
 222 2021]. This order of $O(1/x)$ is set to achieve the optimal sample complexity, and other works have
 223 also investigated into other forms of bonus function [Kaufmann et al., 2021]. Finally, in line 10
 224 through line 12 the current policy π^k —which is the greedy policy with respect to U^k —is simulated
 225 on the empirical reference MDP $\widehat{\mathcal{M}}^\dagger$, and the virtual counters are updated. It is crucial to note that the
 226 simulation takes place entirely offline, by generating virtual transitions from $\widehat{\mathcal{M}}^\dagger$. Upon termination
 227 of algorithm 2, the uniform mixture of policies π^1, π^2, \dots form the non-reactive exploration policy
 228 π_{ex} , ensuring that the latter has wide ‘coverage’ over \mathcal{M}^\dagger .

229 4.4 Online and planning phase

230 Algorithm 2 implements line 1 of the procedure in algorithm 1 by finding the exploratory policy π_{ex} .
 231 After that, in line 2 of the interaction protocol the online dataset \mathcal{D}' is generated by deploying π_{ex} on
 232 the real MDP \mathcal{M} to generate K_{de} trajectories. Conceptually, the online dataset \mathcal{D}' and the offline
 233 dataset \mathcal{D} identify an updated empirical transition kernel $\widetilde{\mathbb{P}}$ and its sparsified version² $\widetilde{\mathbb{P}}^\dagger$. Finally, in
 234 line 3 a reward function r is received, and the value iteration algorithm (See Appendix E) is invoked
 235 with r as reward function and $\widetilde{\mathbb{P}}^\dagger$ as dynamics, and the near-optimal policy π_{final} is produced. The
 236 use of the (updated) empirical sparsified dynamics $\widetilde{\mathbb{P}}^\dagger$ can be seen as incorporating the principle of
 237 pessimism under uncertainty due to the presence of the absorbing state.

238 Our complete algorithm is reported in algorithm 3, and it can be seen as implementing the interaction
 239 protocol described in algorithm 1.

²For any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, we define $\widetilde{\mathbb{P}}^\dagger(s' | s, a) = \frac{m(s, a, s')}{m(s, a)}$ if $N(s, a, s') \geq \Phi$ and $\widetilde{\mathbb{P}}^\dagger(s' | s, a) = 0$ otherwise. Finally, for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, we have $\widetilde{\mathbb{P}}^\dagger(s^\dagger | s, a) = \frac{1}{m(s, a)} \sum_{s' \in \mathcal{S}, N(s, a, s') < \Phi} m(s, a, s')$ and for any $a \in \mathcal{A}$, we have $\widetilde{\mathbb{P}}^\dagger(s^\dagger | s^\dagger, a) = 1$. Here $N(s, a, s')$ is the counter of initial offline data and $m(\cdot, \cdot)$ is the counter of online data.

Algorithm 3 Reward-Free Non-reactive Policy Design (RF-NPD)

Input: Offline dataset \mathcal{D} , target suboptimality $\varepsilon > 0$, failure tolerance $\delta \in (0, 1]$.

- 1: Construct the empirical sparsified MDP $\widehat{\mathcal{M}}^\dagger$.
- 2: *Offline phase:* run RF-UCB($\widehat{\mathcal{M}}^\dagger, K_{ucb}, \varepsilon, \delta$) to obtain the exploratory policy π_{ex} .
- 3: *Online phase:* deploy π_{ex} on the MDP \mathcal{M} for K_{de} episodes to get the online dataset \mathcal{D}' .
- 4: *Planning phase:* receive the reward function r , construct $\widetilde{\mathcal{M}}^\dagger$ from the online dataset \mathcal{D}' , compute π_{final} (which is the optimal policy on $\widetilde{\mathcal{M}}^\dagger$) using value iteration (Appendix E).

Output: π_{final} .

240 5 Main Result

241 In this section, we present a performance bound on our algorithm, namely a bound on the sub-
242 optimality of the value of the final policy π_{final} when measured on the sparsified MDP \mathcal{M}^\dagger . The
243 sparsified MDP arises because it is generally not possible to directly compete with the optimal policy
244 using a non-reactive data collection strategy and a polynomial number of samples due to the lower
245 bound of Xiao et al. [2022]; more details are given in Appendix C.

246 In order to state the main result, we let $K = K_{ucb} = K_{de}$, where K_{ucb} and K_{de} are the number
247 of episodes for the offline simulation and online interaction, respectively. Let C be some universal
248 constant, and choose the threshold in the definition of sparsified MDP as

$$\Phi = 6H^2 \log(12H |\mathcal{S}|^2 |\mathcal{A}| / \delta). \quad (5.1)$$

249 **Theorem 5.1.** *For any $\varepsilon > 0$ and $0 < \delta < 1$, if we let the number of online episodes be*

$$K = \frac{CH^2 |\mathcal{S}|^2 |\mathcal{A}|}{\varepsilon^2} \text{polylog} \left(|\mathcal{S}|, |\mathcal{A}|, H, \frac{1}{\varepsilon}, \frac{1}{\delta} \right),$$

250 *then with probability at least $1 - \delta$, for any reward function r , the final policy π_{final} returned by*
251 *Algorithm 3 satisfies the bound*

$$\max_{\pi \in \Pi} V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi) - V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi_{final}) \leq \varepsilon. \quad (5.2)$$

252 The theorem gives a performance guarantee on the value of the policy π_{final} , which depends both
253 on the initial coverage of the offline dataset \mathcal{D} as well as on the number of samples collected in the
254 online phase. The dependence on the coverage of the offline dataset is implicit through the definition
255 of the (population) sparsified \mathcal{M}^\dagger , which is determined by the counts $N(\cdot, \cdot)$.

256 In order to gain some intuition, we examine some special cases as a function of the coverage of the
257 offline dataset.

258 **Empty dataset** Suppose that the offline dataset \mathcal{D} is empty. Then the sparsified MDP identifies a
259 *multi-armed bandit* at the initial state s_1 , where any action a taken from such state gives back the
260 reward $r(s_1, a)$ and leads to the absorbing state s^\dagger . In this case, our algorithm essentially designs an
261 allocation strategy π_{ex} that is uniform across all actions at the starting state s_1 . Given enough online
262 samples, π_{final} converges to the *action* with the highest instantaneous reward on the multi-armed
263 bandit induced by the start state. With no coverage from the offline dataset, the lower bound of Xiao
264 et al. [2022] for non-reactive policies precludes finding an ε -optimal policy on the original MDP \mathcal{M}
265 unless exponentially many samples are collected.

266 **Known connectivity graph** On the other extreme, assume that the offline dataset contains enough
267 information everywhere in the state-action space such that the critical condition 4.1 is satisfied for
268 all (s, a, s') tuples. Then the sparsified MDP and the real MDP coincide, i.e., $\mathcal{M} = \mathcal{M}^\dagger$, and so the
269 final policy π_{final} directly competes with the optimal policy π^* for any given reward function in
270 eq. (5.2). More precisely, the policy π_{final} is ε -suboptimal on \mathcal{M} if $\widetilde{O}(H^2 |\mathcal{S}|^2 |\mathcal{A}| / \varepsilon^2)$ trajectories
271 are collected in the online phase, a result that matches the lower bound for reward-free exploration of
272 Jin et al. [2020b] up to log factors. However, we achieve such result with a data collection strategy
273 that is completely passive, one that is computed with the help of an initial offline dataset whose size
274 $|\mathcal{D}| \approx \Phi \times |\mathcal{S}|^2 |\mathcal{A}| = \widetilde{O}(H^2 |\mathcal{S}|^2 |\mathcal{A}|)$ need *not depend on final accuracy* ε .

Partial coverage In more typical cases, the offline dataset has only *partial* coverage over the state-action space and the critical condition 4.1 may be violated in certain state-action-successor states. In this case, the connectivity graph of the sparsified MDP \mathcal{M}^\dagger is a sub-graph of the original MDP \mathcal{M} augmented with edges towards the absorbing state. The lack of coverage of the original dataset arises through the sparsified MDP in the guarantees that we present in theorem 5.1. In this section, we ‘translate’ such guarantees into guarantees on \mathcal{M} , in which case the ‘lack of coverage’ is naturally represented by the concentrability coefficient

$$C^* = \sup_{s,a} d_\pi(s, a) / \mu(s, a),$$

275 see for examples the papers [Munos and Szepesvári, 2008, Chen and Jiang, 2019] for background
 276 material on the concentrability factor. More precisely, we compute the sample complexity—in terms
 277 of online as well as offline samples—required for π_{final} to be ε -suboptimal with respect to any
 278 comparator policy π , and so in particular with respect to the optimal policy π_* on the “real” MDP \mathcal{M} .
 279 The next corollary is proved in appendix B.3.

280 **Corollary 5.2.** *Suppose that the offline dataset contains*

$$\tilde{O}\left(\frac{H^4 |\mathcal{S}|^2 |\mathcal{A}| C^*}{\varepsilon}\right),$$

281 *samples and that additional*

$$\tilde{O}\left(\frac{H^3 |\mathcal{S}|^2 |\mathcal{A}|}{\varepsilon^2}\right)$$

282 *online samples are collected during the online phase. Then with probability at least $1 - \delta$, for any*
 283 *reward function r , the policy π_{final} is ε -suboptimal with respect to any comparator policy π*

$$V_1(s_1; \mathbb{P}, r, \pi) - V_1(s_1; \mathbb{P}, r, \pi_{final}) \leq \varepsilon. \quad (5.3)$$

284 The online sample size is equivalent to the one that arises in the statement of theorem 5.1 (expressed
 285 as number of online trajectories), and does not depend on the concentrability coefficient. The
 286 dependence on the offline dataset in theorem 5.1 is implicit in the definition of sparsified MDP; here
 287 we have made it explicit using the notion of concentrability.

288 Corollary 5.2 can be used to compare the achievable guarantees of our procedure with that of an offline
 289 algorithm, such as the minimax-optimal procedure detailed in [Xie et al., 2021b]. The procedure
 290 described in [Xie et al., 2021b] achieves (5.3) with probability at least $1 - \delta$ by using

$$\tilde{O}\left(\frac{H^3 |\mathcal{S}| C^*}{\varepsilon^2} + \frac{H^{5.5} |\mathcal{S}| C^*}{\varepsilon}\right) \quad (5.4)$$

291 offline samples³. In terms of offline data, our procedure has a similar dependence on various factors,
 292 but it depends on the desired accuracy ε through $\tilde{O}(1/\varepsilon)$ as opposed to $\tilde{O}(1/\varepsilon^2)$ which is typical for
 293 an offline algorithm. This implies that in the small- ε regime, if sufficient online samples are collected,
 294 one can improve upon a fully offline procedure by collecting a number of additional online samples
 295 in a non-reactive way.

296 Finally, notice that one may improve upon an offline dataset by collecting more data from the
 297 distribution μ , i.e., without performing experimental design. Compared to this latter case, notice that
 298 our *online sample complexity does not depend on the concentrability coefficient*. Further discussion
 299 can be found in appendix B.

300 6 Proof

301 In this section we prove theorem 5.1, and defer the proofs of the supporting statements to the Appendix
 302 A.

³Technically, [Zhan et al., 2022] considers the non-homogeneous setting, and expresses their result in terms of number of trajectories. In obtaining eq. (5.4), we ‘removed’ an H factor due to our dynamics being homogeneous, and add it back to express the result in terms of number of samples. However, notice that [Zhan et al., 2022] consider the reward-aware setting, which is simpler than reward-free RL setting that we consider. This should add an additional $|\mathcal{S}|$ factor that is not accounted for in eq. (5.4), see the paper Jin et al. [2020b] for more details.

Let us define the comparator policy π_*^\dagger used for the comparison in eq. (5.2) to be the (deterministic) policy with the highest value function on the sparsified MDP:

$$\pi_*^\dagger := \arg \max_{\pi \in \Pi} V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi).$$

303 We can bound the suboptimality using the triangle inequality as

$$\begin{aligned} V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi_*^\dagger) - V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi_{final}) &\leq \left| V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi_*^\dagger) - V_1(s_1; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi_*^\dagger) \right| \\ &\quad + \underbrace{V_1(s_1; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi_*^\dagger) - V_1(s_1; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi_{final})}_{\leq 0} + \left| V_1(s_1; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi_{final}) - V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi_{final}) \right| \\ &\leq 2 \sup_{\pi \in \Pi, r} \left| V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi) - V_1(s_1; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi) \right|. \end{aligned}$$

304 The middle term after the first inequality is negative due to the optimality of π_{final} on $\tilde{\mathbb{P}}^\dagger$ and r^\dagger . It
305 suffices to prove that for any arbitrary policy π and reward function r the following statement holds
306 with probability at least $1 - \delta$

$$\underbrace{\left| V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi) - V_1(s_1; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi) \right|}_{\text{Estimation error}} \leq \frac{\varepsilon}{2}. \quad (6.1)$$

307 **Bounding the estimation error using the population uncertainty function** In order to prove
308 eq. (6.1), we first define the population *uncertainty function* X , which is a scalar function over the
309 state-action space. It represents the maximum estimation error on the value of any policy when it
310 is evaluated on $\tilde{\mathcal{M}}^\dagger$ instead of \mathcal{M} . For any $(s, a) \in \mathcal{S} \times \mathcal{A}$, the uncertainty function is defined as
311 $X_{H+1}(s, a) := 0$ and for $h \in [H]$,

$$X_h(s, a) := \min \left\{ H - h + 1; 9H\phi(m(s, a)) + \left(1 + \frac{1}{H} \right) \sum_{s'} \tilde{\mathbb{P}}^\dagger(s' | s, a) \left(\max_{a'} \{X_{h+1}(s', a')\} \right) \right\}.$$

312 We extend the definition to the absorbing state by letting $X_h(s^\dagger, a) = 0$ for any $h \in [H], a \in \mathcal{A}$. The
313 summation $\sum_{s'}$ used above is over $s' \in \mathcal{S} \cup \{s^\dagger\}$, but since $X_h(s^\dagger, a) = 0$ for any $h \in [H], a \in \mathcal{A}$,
314 it is equivalent to that over $s' \in \mathcal{S}$. Intuitively, $X_h(s, a)$ takes a similar form as Bellman optimality
315 equation. The additional $(1 + 1/H)$ factor and additional term $9H\phi(m(s, a))$ quantify the uncertainty
316 of the true Q function on the sparsified MDP and $9H\phi(m(s, a))$ will converge to zero when the
317 sample size goes to infinity. This definition of uncertainty function and the following lemma follow
318 closely from the uncertainty function defined in [Ménard et al., 2021].

319 The next lemma highlights the key property of the uncertainty function X , namely that for any reward
320 function and any policy π , we can upper bound the estimation error via the uncertainty function at
321 the initial times-step; it is proved in appendix A.2.1.

322 **Lemma 6.1.** *With probability $1 - \delta$, for any reward function r and any deterministic policy π , it*
323 *holds that*

$$\left| V_1(s_1; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi) - V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi) \right| \leq \max_a X_1(s_1, a) + C \sqrt{\max_a X_1(s_1, a)}. \quad (6.2)$$

324 The uncertainty function X contains the inverse number of *online* samples $1/m(s, a)$ through
325 $\phi(m(s, a))$, and so lemma 6.1 expresses the estimation error in eq. (6.1) as the maximum expected
326 size of the confidence intervals $\sup_{\pi} \mathbb{E}_{\tilde{\mathbb{P}}^\dagger, (s, a) \sim \pi} \sqrt{1/m(s, a)}$, a quantity that directly depends on the
327 number $m(\cdot, \cdot)$ of samples collected during the online phase.

328 **Leveraging the exploration mechanics** Throughout this section, C denotes some universal constant
329 and may vary from line to line. Recall that the agent greedily minimizes the *empirical uncertainty*
330 *function* U to compute the exploratory policy π_{ex} . The empirical uncertainty is defined as
331 $U_{H+1}^k(s, a) = 0$ for any $k, s \in \mathcal{S}, a \in \mathcal{A}$ and

$$U_h^k(s, a) = H \min\{1, \phi(n^k(s, a))\} + \hat{\mathbb{P}}^\dagger(s, a)^\top \left(\max_{a'} U_{h+1}^k(\cdot, a') \right), \quad (6.3)$$

332 where $n^k(s, a)$ is the counter of the times we encounter (s, a) until the beginning of the k -th virtual
 333 episode in the simulation phase. Note that, $U_h^k(s, a)$ takes a similar form as $X_h(s, a)$, except that
 334 $U_h^k(s, a)$ depends on the empirical transition probability $\widehat{\mathbb{P}}^\dagger$ while $X_h(s, a)$ depends on the true
 335 transition probability on the sparsified MDP. For the exploration scheme to be effective, X and U
 336 should be close in value, a concept which is at the core of this work and which we formally state
 337 below and prove in appendix A.2.2.

338 **Lemma 6.2** (Bounding uncertainty function with empirical uncertainty functions). *With probability*
 339 *at least $1 - \delta$, we have for any $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$,*

$$X_h(s, a) \leq \frac{C}{K} \sum_{k=1}^K U_h^k(s, a).$$

340 Notice that X_h is the population uncertainty after the online samples have been collected, while U_h^k
 341 is the corresponding empirical uncertainty which varies during the planning phase.

342 **Rate of decrease of the estimation error** Combining lemmas 6.1 and 6.2 shows that (a function
 343 of) the agent's uncertainty estimate U upper bounds the estimation error in eq. (6.1). In order to
 344 conclude, we need to show that U decreases on average at the rate $1/K$, a statement that we present
 345 below and prove in appendix A.2.3.

346 **Lemma 6.3.** *With probability at least $1 - \delta$, we have*

$$\frac{1}{K} \sum_{k=1}^K U_1^k(s, a) \leq \frac{H^2 |\mathcal{S}|^2 |\mathcal{A}|}{K} \text{polylog} \left(K, |\mathcal{S}|, |\mathcal{A}|, H, \frac{1}{\varepsilon}, \frac{1}{\delta} \right). \quad (6.4)$$

347 *Then, for any $\varepsilon > 0$, if we take*

$$K := \frac{CH^2 |\mathcal{S}| |\mathcal{A}|}{\varepsilon^2} \left(\iota + |\mathcal{S}| \right) \text{polylog} \left(|\mathcal{S}|, |\mathcal{A}|, H, \frac{1}{\varepsilon}, \frac{1}{\delta} \right),$$

348 *then with probability at least $1 - \delta$, it holds that*

$$\frac{1}{K} \sum_{k=1}^K U_1^k(s_1, a) \leq \varepsilon^2.$$

349 After combining lemmas 6.1 to 6.3, we see that the estimation error can be bounded as

$$\begin{aligned} \left| V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi) - V_1(s_1; \widetilde{\mathbb{P}}^\dagger, r^\dagger, \pi) \right| &\leq \max_a X_1(s_1, a) + C \sqrt{\max_a X_1(s_1, a)} \\ &\leq C \max_a \left[\sqrt{\frac{1}{K} \sum_{k=1}^K U_1^k(s_1, a)} + \frac{1}{K} \sum_{k=1}^K U_1^k(s_1, a) \right] \\ &\leq C (\varepsilon + \varepsilon^2) \\ &\leq C \varepsilon \quad (\text{for } 0 < \varepsilon < \text{const}) \end{aligned}$$

350 Here, the constant C may vary between lines. Rescaling the universal constant C and the failure
 351 probability δ , we complete the upper bound in equation (6.1) and hence the proof for the main result.

352 References

- 353 M Mehdi Afsar, Trafford Crump, and Behrouz Far. Reinforcement learning based recommender
354 systems: A survey. *ACM Computing Surveys*, 55(7):1–38, 2022.
- 355 Arpit Agarwal, Rohan Ghuge, and Viswanath Nagarajan. Batched dueling bandits. In *International*
356 *Conference on Machine Learning*, pages 89–110. PMLR, 2022.
- 357 Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit
358 problem. *Machine learning*, 47:235–256, 2002.
- 359 Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for rein-
360 forcement learning. In *International Conference on Machine Learning*, pages 263–272. PMLR,
361 2017.
- 362 Yu Bai, Tengyang Xie, Nan Jiang, and Yu-Xiang Wang. Provably efficient q-learning with low
363 switching cost. *Advances in Neural Information Processing Systems*, 32, 2019.
- 364 Nicolo Cesa-Bianchi, Ofer Dekel, and Ohad Shamir. Online learning with switching costs and other
365 adaptive adversaries. *Advances in Neural Information Processing Systems*, 26, 2013.
- 366 Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement learning. In
367 *International Conference on Machine Learning*, pages 1042–1051. PMLR, 2019.
- 368 Jinglin Chen, Aditya Modi, Akshay Krishnamurthy, Nan Jiang, and Alekh Agarwal. On the statistical
369 efficiency of reward-free exploration in non-linear rl. *arXiv preprint arXiv:2206.10770*, 2022.
- 370 Christoph Dann and Emma Brunskill. Sample complexity of episodic fixed-horizon reinforcement
371 learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2015.
- 372 Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying pac and regret: Uniform pac bounds
373 for episodic reinforcement learning. *Advances in Neural Information Processing Systems*, 30,
374 2017.
- 375 Christoph Dann, Lihong Li, Wei Wei, and Emma Brunskill. Policy certificates: Towards accountable
376 reinforcement learning. In *International Conference on Machine Learning*, pages 1507–1516,
377 2019.
- 378 Simon Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudik, and John Langford.
379 Provably efficient rl with rich observations via latent state decoding. In *International Conference*
380 *on Machine Learning*, pages 1665–1674. PMLR, 2019.
- 381 Yaqi Duan, Chi Jin, and Zhiyuan Li. Risk bounds and rademacher complexity in batch reinforcement
382 learning. In *International Conference on Machine Learning*, pages 2892–2902. PMLR, 2021.
- 383 John Duchi, Feng Ruan, and Chulhee Yun. Minimax bounds on stochastic batched convex optimiza-
384 tion. In *Conference On Learning Theory*, pages 3065–3162. PMLR, 2018.
- 385 Rick Durrett. *Probability: theory and examples*, volume 49. Cambridge university press, 2019.
- 386 Dylan J Foster, Akshay Krishnamurthy, David Simchi-Levi, and Yunzong Xu. Offline reinforcement
387 learning: Fundamental barriers for value function approximation. *arXiv preprint arXiv:2111.10919*,
388 2021.
- 389 Minbo Gao, Tianle Xie, Simon S Du, and Lin F Yang. A provably efficient algorithm for linear
390 markov decision process with low switching cost. *arXiv preprint arXiv:2101.00494*, 2021.
- 391 Zijun Gao, Yanjun Han, Zhimei Ren, and Zhengqing Zhou. Batched multi-armed bandits problem.
392 *Advances in Neural Information Processing Systems*, 32, 2019.
- 393 Omer Gottesman, Fredrik Johansson, Joshua Meier, Jack Dent, Donghun Lee, Srivatsan Srinivasan,
394 Linying Zhang, Yi Ding, David Wihl, Xuefeng Peng, Jiayu Yao, Isaac Lage, Christopher Mosch,
395 Li wei H. Lehman, Matthieu Komorowski, Matthieu Komorowski, Aldo Faisal, Leo Anthony
396 Celi, David Sontag, and Finale Doshi-Velez. Evaluating reinforcement learning algorithms in
397 observational health settings, 2018.

- 398 Elad Hazan, Sham Kakade, Karan Singh, and Abby Van Soest. Provably efficient maximum entropy
399 exploration. In *International Conference on Machine Learning*, pages 2681–2691. PMLR, 2019.
- 400 Jiawei Huang, Jinglin Chen, Li Zhao, Tao Qin, Nan Jiang, and Tie-Yan Liu. Towards deployment-
401 efficient reinforcement learning: Lower bound and optimality. *arXiv preprint arXiv:2202.06450*,
402 2022.
- 403 Nan Jiang and Jiawei Huang. Minimax value interval for off-policy evaluation and policy optimization.
404 *Advances in Neural Information Processing Systems*, 33:2747–2758, 2020.
- 405 Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is q-learning provably efficient?
406 In *Advances in Neural Information Processing Systems*, pages 4863–4873, 2018.
- 407 Chi Jin, Akshay Krishnamurthy, Max Simchowitz, and Tiancheng Yu. Reward-free exploration for
408 reinforcement learning. In *International Conference on Machine Learning*, pages 4870–4879.
409 PMLR, 2020a.
- 410 Chi Jin, Akshay Krishnamurthy, Max Simchowitz, and Tiancheng Yu. Reward-free exploration for
411 reinforcement learning. In *International Conference on Machine Learning (ICML)*, 2020b.
- 412 Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? *arXiv*
413 *preprint arXiv:2012.15085*, 2020c.
- 414 Anders Jonsson, Emilie Kaufmann, Pierre Ménard, Omar Darwiche Domingues, Edouard Leurent,
415 and Michal Valko. Planning in markov decision processes with gap-dependent sample complexity.
416 *Advances in Neural Information Processing Systems*, 33:1253–1263, 2020.
- 417 Nathan Kallus and Masatoshi Uehara. Efficiently breaking the curse of horizon in off-policy evaluation
418 with double reinforcement learning. *Operations Research*, 2022.
- 419 Emilie Kaufmann, Pierre Ménard, Omar Darwiche Domingues, Anders Jonsson, Edouard Leurent,
420 and Michal Valko. Adaptive reward-free exploration. In *Algorithmic Learning Theory*, pages
421 865–891. PMLR, 2021.
- 422 Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-
423 based offline reinforcement learning. *arXiv preprint arXiv:2005.05951*, 2020.
- 424 B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani,
425 and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE*
426 *Transactions on Intelligent Transportation Systems*, 23(6):4909–4926, 2021.
- 427 Allison Koenecke, Andrew Nam, Emily Lake, Joe Nudell, Minnie Quartey, Zion Mengesha, Connor
428 Touns, John R Rickford, Dan Jurafsky, and Sharad Goel. Racial disparities in automated speech
429 recognition. *Proceedings of the National Academy of Sciences*, 117(14):7684–7689, 2020.
- 430 Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- 431 Gen Li, Yuling Yan, Yuxin Chen, and Jianqing Fan. Optimal reward-agnostic exploration in rein-
432 forcement learning. 2023a.
- 433 Gen Li, Wenhao Zhan, Jason D Lee, Yuejie Chi, and Yuxin Chen. Reward-agnostic fine-tuning:
434 Provable statistical benefits of hybrid reinforcement learning. *arXiv preprint arXiv:2305.10282*,
435 2023b.
- 436 Yao Liu, Adith Swaminathan, Alekh Agarwal, and Emma Brunskill. Provably good batch reinforce-
437 ment learning without great exploration. *arXiv preprint arXiv:2007.08202*, 2020.
- 438 Jihao Long, Jiequn Han, and E Weinan. An l2 analysis of reinforcement learning in high dimensions
439 with kernel and neural network approximation. *arXiv preprint arXiv:2104.07794*, 3, 2021.
- 440 Andreas Maurer and Massimiliano Pontil. Empirical bernstein bounds and sample variance penaliza-
441 tion. *arXiv preprint arXiv:0907.3740*, 2009.

- 442 Pierre Ménard, Omar Darwiche Domingues, Anders Jonsson, Emilie Kaufmann, Edouard Leurent,
443 and Michal Valko. Fast active learning for pure exploration in reinforcement learning. In *International Conference on Machine Learning*, pages 7599–7608. PMLR, 2021.
- 445 Dipendra Misra, Mikael Henaff, Akshay Krishnamurthy, and John Langford. Kinematic state abstraction
446 and provably efficient rich-observation reinforcement learning. In *International conference on machine learning*, pages 6961–6971. PMLR, 2020.
- 448 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare,
449 Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control
450 through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 451 Rémi Munos and Csaba Szepesvári. Finite-time bounds for fitted value iteration. *Journal of Machine
452 Learning Research*, 9(5), 2008.
- 453 Ofir Nachum, Bo Dai, Ilya Kostrikov, Yinlam Chow, Lihong Li, and Dale Schuurmans. Algaedice:
454 Policy gradient from arbitrary experience. *arXiv preprint arXiv:1912.02074*, 2019.
- 455 Thanh Nguyen-Tang, Ming Yin, Sunil Gupta, Svetha Venkatesh, and Raman Arora. On instance-
456 dependent bounds for offline reinforcement learning with linear function approximation. *arXiv
457 preprint arXiv:2211.13208*, 2022.
- 458 Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John
459 Wiley & Sons, Inc., New York, NY, USA, 1994. ISBN 0471619779.
- 460 Dan Qiao and Yu-Xiang Wang. Near-optimal deployment efficiency in reward-free reinforcement
461 learning with linear function approximation. *arXiv preprint arXiv:2210.00701*, 2022.
- 462 Dan Qiao, Ming Yin, Ming Min, and Yu-Xiang Wang. Sample-efficient reinforcement learning with
463 loglog (t) switching cost. In *International Conference on Machine Learning*, pages 18031–18061.
464 PMLR, 2022.
- 465 Dan Qiao, Ming Yin, and Yu-Xiang Wang. Logarithmic switching cost in reinforcement learning
466 beyond linear mdps. *arXiv preprint arXiv:2302.12456*, 2023.
- 467 Shuang Qiu, Jieping Ye, Zhaoran Wang, and Zhuoran Yang. On reward-free rl with kernel and neural
468 function approximations: Single-agent mdp and markov game. In *International Conference on
469 Machine Learning*, pages 8737–8747. PMLR, 2021.
- 470 Paul Raccuglia, Katherine C Elbert, Philip DF Adler, Casey Falk, Malia B Wenny, Aurelio Mollo,
471 Matthias Zeller, Sorelle A Friedler, Joshua Schrier, and Alexander J Norquist. Machine-learning-
472 assisted materials discovery using failed experiments. *Nature*, 533(7601):73–76, 2016.
- 473 Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging offline rein-
474 forcement learning and imitation learning: A tale of pessimism. *arXiv preprint arXiv:2103.12021*,
475 2021.
- 476 Paria Rashidinejad, Hanlin Zhu, Kunhe Yang, Stuart Russell, and Jiantao Jiao. Optimal conser-
477 vative offline rl with general function approximation via augmented lagrangian. *arXiv preprint
478 arXiv:2211.00716*, 2022.
- 479 Yufei Ruan, Jiaqi Yang, and Yuan Zhou. Linear bandits with limited adaptivity and learning
480 distributional optimal design. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on
481 Theory of Computing*, pages 74–87, 2021.
- 482 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,
483 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering
484 the game of go with deep neural networks and tree search. *Nature*, 529(7587):484, 2016.
- 485 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT Press, 2018.
- 486 Mohammad Sadegh Talebi and Odalric-Ambrym Maillard. Variance-aware regret bounds for undis-
487 counted reinforcement learning in mdps. In *Algorithmic Learning Theory*, pages 770–805. PMLR,
488 2018.

- 489 Masatoshi Uehara and Wen Sun. Pessimistic model-based offline reinforcement learning under partial
490 coverage, 2021.
- 491 Andrew J Wagenmaker, Yifang Chen, Max Simchowitz, Simon Du, and Kevin Jamieson. Reward-
492 free rl is no harder than reward-aware rl in linear markov decision processes. In *International*
493 *Conference on Machine Learning*, pages 22430–22456. PMLR, 2022.
- 494 Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cam-
495 bridge university press, 2019.
- 496 Ruosong Wang, Simon S Du, Lin Yang, and Russ R Salakhutdinov. On reward-free reinforcement
497 learning with linear function approximation. *Advances in neural information processing systems*,
498 33:17816–17826, 2020a.
- 499 Ruosong Wang, Dean P Foster, and Sham M Kakade. What are the statistical limits of offline rl with
500 linear function approximation? *arXiv preprint arXiv:2010.11895*, 2020b.
- 501 Tianhao Wang, Dongruo Zhou, and Quanquan Gu. Provably efficient reinforcement learning with
502 linear function approximation under adaptivity constraints. *Advances in Neural Information*
503 *Processing Systems*, 34:13524–13536, 2021.
- 504 Chenjun Xiao, Ilbin Lee, Bo Dai, Dale Schuurmans, and Csaba Szepesvari. The curse of passive data
505 collection in batch reinforcement learning. In *International Conference on Artificial Intelligence*
506 *and Statistics*, pages 8413–8438. PMLR, 2022.
- 507 Tengyang Xie and Nan Jiang. Batch value-function approximation with only realizability. *arXiv*
508 *preprint arXiv:2008.04990*, 2020a.
- 509 Tengyang Xie and Nan Jiang. Q^* approximation schemes for batch reinforcement learning: A
510 theoretical comparison. In *Conference on Uncertainty in Artificial Intelligence*, pages 550–559.
511 PMLR, 2020b.
- 512 Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent
513 pessimism for offline reinforcement learning. *arXiv preprint arXiv:2106.06926*, 2021a.
- 514 Tengyang Xie, Nan Jiang, Huan Wang, Caiming Xiong, and Yu Bai. Policy finetuning: Bridg-
515 ing sample-efficient offline and online reinforcement learning. *Advances in neural information*
516 *processing systems*, 34:27395–27407, 2021b.
- 517 Wei Xiong, Han Zhong, Chengshuai Shi, Cong Shen, Liwei Wang, and Tong Zhang. Nearly minimax
518 optimal offline reinforcement learning with linear function approximation: Single-agent mdp and
519 markov game. *arXiv preprint arXiv:2205.15512*, 2022.
- 520 Zhiyuan Xu, Jian Tang, Jingsong Meng, Weiye Zhang, Yanzhi Wang, Chi Harold Liu, and Dejun
521 Yang. Experience-driven networking: A deep reinforcement learning based approach. In *IEEE*
522 *INFOCOM 2018-IEEE conference on computer communications*, pages 1871–1879. IEEE, 2018.
- 523 Tsung-Yen Yang, Justinian Rosca, Karthik Narasimhan, and Peter J. Ramadge. Accelerating safe
524 reinforcement learning with constraint-mismatched policies, 2021.
- 525 Ming Yin and Yu-Xiang Wang. Asymptotically efficient off-policy evaluation for tabular rein-
526 forcement learning. In *International Conference on Artificial Intelligence and Statistics*, pages
527 3948–3958. PMLR, 2020.
- 528 Ming Yin, Yu-Xiang Wang, Yaqi Duan, and Mengdi Wang. Near-optimal offline reinforcement
529 learning with linear representation: Leveraging variance information with pessimism.
- 530 Ming Yin, Yu Bai, and Yu-Xiang Wang. Near optimal provable uniform convergence in off-policy
531 evaluation for reinforcement learning. *arXiv preprint arXiv:2007.03760*, 2020.
- 532 Ming Yin, Mengdi Wang, and Yu-Xiang Wang. Offline reinforcement learning with differentiable
533 function approximation is provably efficient. *arXiv preprint arXiv:2210.00750*, 2022.
- 534 Chao Yu, Jiming Liu, Shamim Nemati, and Guosheng Yin. Reinforcement learning in healthcare: A
535 survey. *ACM Computing Surveys (CSUR)*, 55(1):1–36, 2021.

- 536 Andrea Zanette. When is realizability sufficient for off-policy reinforcement learning? *ICML 2023*,
537 2023.
- 538 Andrea Zanette and Emma Brunskill. Tighter problem-dependent regret bounds in reinforcement
539 learning without domain knowledge using value function bounds. In *International Conference on*
540 *Machine Learning*, pages 7304–7312. PMLR, 2019.
- 541 Andrea Zanette and Martin J. Wainwright. Bellman residual orthogonalization for offline reinforce-
542 ment learning, 2022. URL <https://arxiv.org/abs/2203.12786>.
- 543 Andrea Zanette, Alessandro Lazaric, Mykel J Kochenderfer, and Emma Brunskill. Provably efficient
544 reward-agnostic navigation with linear value iteration. *Advances in Neural Information Processing*
545 *Systems*, 33:11756–11766, 2020.
- 546 Andrea Zanette, Kefan Dong, Jonathan N Lee, and Emma Brunskill. Design of experiments for
547 stochastic contextual linear bandits. *Advances in Neural Information Processing Systems*, 34:
548 22720–22731, 2021a.
- 549 Andrea Zanette, Martin J Wainwright, and Emma Brunskill. Provable benefits of actor-critic methods
550 for offline reinforcement learning. *Advances in neural information processing systems*, 34:13626–
551 13640, 2021b.
- 552 Wenhao Zhan, Baihe Huang, Audrey Huang, Nan Jiang, and Jason D Lee. Offline reinforcement
553 learning with realizability and single-policy concentrability. *arXiv preprint arXiv:2202.04634*,
554 2022.
- 555 Ruiqi Zhang, Xuezhou Zhang, Chengzhuo Ni, and Mengdi Wang. Off-policy fitted q-evaluation
556 with differentiable function approximators: Z-estimation and inference theory. In *International*
557 *Conference on Machine Learning*, pages 26713–26749. PMLR, 2022.
- 558 Xuezhou Zhang, Yuzhe Ma, and Adish Singla. Task-agnostic exploration in reinforcement learning.
559 *Advances in Neural Information Processing Systems*, 33:11734–11743, 2020a.
- 560 Zihan Zhang, Yuan Zhou, and Xiangyang Ji. Almost optimal model-free reinforcement learning via
561 reference-advantage decomposition. *Advances in Neural Information Processing Systems*, 33:
562 15198–15207, 2020b.
- 563 Zihan Zhang, Simon Du, and Xiangyang Ji. Near optimal reward-free reinforcement learning. In
564 *International Conference on Machine Learning*, pages 12402–12412. PMLR, 2021.

565	Contents	
566	1 Introduction	1
567	2 Related Work	2
568	3 Setup	3
569	3.1 Interaction protocol	3
570	4 Algorithm: balancing optimism and pessimism for experimental design	4
571	4.1 Intuition	4
572	4.2 Sparsified MDP	4
573	4.3 Offline design of experiments	5
574	4.4 Online and planning phase	6
575	5 Main Result	7
576	6 Proof	8
577	A Proof of the main result	18
578	A.1 Definitions	18
579	A.1.1 Sparsified MDP	18
580	A.1.2 High probability events	19
581	A.1.3 Uncertainty function and empirical uncertainty function	20
582	A.2 Proof of the key theorems and lemmas	21
583	A.2.1 Uncertainty Functions upper bounds the estimation error	21
584	A.2.2 Proof of lemma 6.2	22
585	A.2.3 Upper bounding the empirical uncertainty function (lemma 6.3)	23
586	A.3 Omitted proofs	25
587	A.3.1 Proof for lemma A.4 (high probability event)	25
588	A.3.2 Proof for lemma A.14 (property of intermediate uncertainty function)	26
589	A.3.3 Proof for lemma A.15 and lemma A.16 (properties of uncertainty function)	28
590	A.3.4 Proof for lemma A.17, lemma A.18, lemma A.19, lemma A.20 (properties of bonus function)	29
591	A.3.5 Proof for lemma A.21 and lemma A.22 (properties of empirical sparsified MDP)	31
592		
593		
594	B Additional comparisons	33
595	B.1 Comparison with other comparator policy	33
596	B.2 Comparison with offline reinforcement learning	33
597	B.3 Proof of corollary 5.2	35
598	B.4 Proof for lemma B.1 and lemma B.2	38
599	C Lower bound	39

600	D	Technical lemmas and proofs	40
601	E	Details of the planning phase	42
602	F	More related works	42

603 A Proof of the main result

604 A.1 Definitions

605 In this section, we define some crucial concepts that will be used in the proof of the main result.

606 A.1.1 Sparsified MDP

607 First, we restate the definition 4.1 in the main text.

608 **Definition A.1** (Sparsified MDP). *Let s^\dagger be an absorbing state, i.e., such that $\mathbb{P}(s^\dagger | s^\dagger, a) = 1$*
 609 *and $r(s^\dagger, a) = 0$ for all $a \in \mathcal{A}$. The state space in the sparsified MDP \mathcal{M}^\dagger is defined as that of the*
 610 *original MDP with the addition of s^\dagger . The dynamics \mathbb{P}^\dagger of the sparsified MDP are defined as*

$$611 \mathbb{P}^\dagger(s' | s, a) = \begin{cases} \mathbb{P}(s' | s, a) & \text{if } N(s, a, s') \geq \Phi \\ 0 & \text{if } N(s, a, s') < \Phi, \end{cases} \quad \mathbb{P}^\dagger(s^\dagger | s, a) = \sum_{\substack{s' \neq s^\dagger \\ N(s, a, s') < \Phi}} \mathbb{P}(s' | s, a). \quad (\text{A.1})$$

611 *For any deterministic reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, the reward function on the sparsified MDP*
 612 *is defined as $r^\dagger(s, a) = r(s, a)$; for simplicity we only consider deterministic reward functions.*

613

614 In the offline phase of our algorithm, we simulate the virtual episodes from the empirical version of
 615 sparsified MDP (See Algorithm 2). Now we formally this MDP.

616 **Definition A.2** (Empirical sparsified MDP). *Let s^\dagger be the absorbing state defined in the sparsified*
 617 *MDP. The state space in the empirical sparsified MDP $\widehat{\mathcal{M}}^\dagger$ is defined as that of the original MDP*
 618 *with the addition of s^\dagger . The dynamics $\widehat{\mathbb{P}}^\dagger$ of the sparsified MDP are defined as*

$$619 \widehat{\mathbb{P}}^\dagger(s' | s, a) = \begin{cases} \frac{N(s, a, s')}{N(s, a)} & \text{if } N(s, a, s') \geq \Phi \\ 0 & \text{if } N(s, a, s') < \Phi, \end{cases} \quad \widehat{\mathbb{P}}^\dagger(s^\dagger | s, a) = \sum_{\substack{s' \neq s^\dagger \\ N(s, a, s') < \Phi}} \frac{N(s, a, s')}{N(s, a)}. \quad (\text{A.2})$$

619 *For any deterministic reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, the reward function on the empirical*
 620 *sparsified MDP is defined as $r^\dagger(s, a) = r(s, a)$; for simplicity we only consider deterministic reward*
 621 *functions. Here, the counters $N(s, a)$ and $N(s, a, s')$ are the number of (s, a) and (s, a, s') in the*
 622 *offline data, respectively.*

623

624 Finally, in the planning phase, after we interact with the true environment for many online episodes,
 625 construct a fine-estimated sparsified MDP, which is used to extract the optimal policy of given reward
 626 functions. We formally define it below.

627 **Definition A.3** (Fine-estimated sparsified MDP). *Let s^\dagger be the absorbing state defined in the sparsified*
 628 *MDP. The state space in the fine-estimated sparsified MDP $\widetilde{\mathcal{M}}^\dagger$ is defined as that of the original*
 629 *MDP with the addition of s^\dagger . The dynamics $\widetilde{\mathbb{P}}^\dagger$ of the sparsified MDP are defined as*

$$630 \widetilde{\mathbb{P}}^\dagger(s' | s, a) = \begin{cases} \frac{m(s, a, s')}{m(s, a)} & \text{if } N(s, a, s') \geq \Phi \\ 0 & \text{if } N(s, a, s') < \Phi, \end{cases} \quad \widetilde{\mathbb{P}}^\dagger(s^\dagger | s, a) = \sum_{\substack{s' \neq s^\dagger \\ N(s, a, s') < \Phi}} \frac{m(s, a, s')}{m(s, a)}. \quad (\text{A.3})$$

630 *For any deterministic reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, the reward function on the fine-estimated*
 631 *sparsified MDP is defined as $r^\dagger(s, a) = r(s, a)$; for simplicity we only consider deterministic reward*
 632 *functions. Here, the counters $N(s, a)$ and $N(s, a, s')$ are the number of (s, a) and (s, a, s') in the*
 633 *offline data, respectively, while $m(s, a)$ and $m(s, a, s')$ are the counters of (s, a) and (s, a, s') in*
 634 *online episodes.*

635 **A.1.2 High probability events**

636 In this section, we define all high-probability events we need in order to make our theorem to hold.
 637 Specifically, we define

$$\begin{aligned}
 \mathcal{E}^P &:= \left\{ \forall (s, a, s') \text{ s.t. } N(s, a, s') \geq \Phi, \left| \widehat{\mathbb{P}}^\dagger(s' | s, a) - \mathbb{P}^\dagger(s' | s, a) \right| \right. \\
 &\quad \left. \leq \sqrt{\frac{2\widehat{\mathbb{P}}^\dagger(s' | s, a)}{N(s, a)} \log\left(\frac{12|\mathcal{S}|^2|\mathcal{A}|}{\delta}\right)} + \frac{14}{3N(s, a)} \log\left(\frac{12|\mathcal{S}|^2|\mathcal{A}|}{\delta}\right) \right\}; \\
 \mathcal{E}^2 &:= \left\{ \forall k \in \mathbb{N}, (s, a) \in \mathcal{S} \times \mathcal{A}, n^k(s, a) \geq \frac{1}{2} \sum_{i < k} \sum_{h \in [H]} \widehat{d}_{\pi^i, h}^\dagger(s, a) - H \ln\left(\frac{6H|\mathcal{S}||\mathcal{A}|}{\delta}\right) \right\}; \\
 \mathcal{E}^3 &:= \left\{ \forall k \in \mathbb{N}, (s, a) \in \mathcal{S} \times \mathcal{A}, n^k(s, a) \leq 2 \sum_{i < k} \sum_{h \in [H]} \widehat{d}_{\pi^i, h}^\dagger(s, a) + H \ln\left(\frac{6H|\mathcal{S}||\mathcal{A}|}{\delta}\right) \right\}; \\
 \mathcal{E}^4 &:= \left\{ \forall (s, a) \in \mathcal{S} \times \mathcal{A}, \text{KL}\left(\widetilde{\mathbb{P}}^\dagger(s, a); \mathbb{P}^\dagger(s, a)\right) \right. \\
 &\quad \left. \leq \frac{1}{m(s, a)} \left[\log\left(\frac{6|\mathcal{S}||\mathcal{A}|}{\delta}\right) + |\mathcal{S}| \log\left(e \left(1 + \frac{m(s, a)}{|\mathcal{S}|}\right)\right) \right] \right\}; \\
 \mathcal{E}^5 &:= \left\{ \forall (s, a) \in [H] \times \mathcal{S} \times \mathcal{A}, m(s, a) \geq \frac{1}{2} K_{de} \sum_{h \in [H]} w_h^{mix}(s, a) - H \ln\left(\frac{6H|\mathcal{S}||\mathcal{A}|}{\delta}\right) \right\},
 \end{aligned}$$

638 where

$$w_h^{mix}(s, a) := \frac{1}{K_{ucb}} \sum_{k=1}^{K_{ucb}} d_{\pi^k, h}^\dagger(s, a).$$

639 Here, $\text{KL}(\cdot; \cdot)$ denotes the Kullback–Leibler divergence between two distributions. K_{ucb} is the
 640 number of episodes of the sub-routine RF-UCB and π^i is the policy executed at the i -th episode of
 641 RF-UCB (See algorithm 2). K_{de} is the number of episodes executed in the online phase. $n^k(s, a)$ is
 642 the counter of (s, a) before the beginning of the k -th episode in the offline phase and $m(s, a)$ is the
 643 number of (s, a) samples in the online data (See definitions in section A.1). We denote $d_{\pi, h}(s, a)$,
 644 $d_{\pi, h}^\dagger(s, a)$ and $\widehat{d}_{\pi, h}^\dagger(s, a)$ as the occupancy measure of (s, a) at stage h under policy π , on \mathbb{P} (the
 645 true transition dynamics), \mathbb{P}^\dagger (the transition dynamics in the sparsified MDP) and $\widehat{\mathbb{P}}^\dagger$ (the transition
 646 dynamics in the empirical sparsified MDP) respectively.

647 For the event defined above, we have the following guarantee, which is proved in appendix A.3.1.

Lemma A.4. *For $\delta > 0$, we have*

$$\mathcal{E} := \{\mathcal{E}^P \cup \mathcal{E}^2 \cup \mathcal{E}^3 \cup \mathcal{E}^4 \cup \mathcal{E}^5\}$$

648 happens with probability at least $1 - \delta$,

649 **A.1.3 Uncertainty function and empirical uncertainty function**

650 In this section, we restate the definition of uncertainty function and empirical uncertainty functions,
 651 as well as some intermediate uncertainty functions which will be often used in the proof. First, we
 652 restate the definition of bonus function.

653 **Definition A.5** (Bonus Function). *We define*

$$\phi(x) = \frac{H}{x} \left[\log \left(\frac{6H|\mathcal{S}||\mathcal{A}|}{\delta} \right) + |\mathcal{S}| \log \left(e \left(1 + \frac{x}{|\mathcal{S}|} \right) \right) \right]. \quad (\text{A.4})$$

654 *Further, we define*

$$\bar{\phi}(x) = \min \{1, H\phi(x)\}. \quad (\text{A.5})$$

655 Next, we restate the definition of uncertainty function and empirical uncertainty function. Notice that
 656 the uncertainty function does not depend on reward function or specific policy π .

657 **Definition A.6** (Uncertainty function). *We define for any $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$ and any deterministic*
 658 *policy π ,*

$$\begin{aligned} X_{H+1}(s, a) &:= 0, \\ X_h(s, a) &:= \min \left\{ H - h + 1, 9H\phi(m(s, a)) + \left(1 + \frac{1}{H} \right) \tilde{\mathbb{P}}^\dagger(\cdot | s, a)^\top \left(\max_{a'} X_{h+1}(\cdot, a') \right) \right\}, \end{aligned}$$

659 *where we specify*

$$\tilde{\mathbb{P}}^\dagger(\cdot | s, a)^\top \left(\max_{a'} \{X_{h+1}(\cdot, a')\} \right) := \sum_{s'} \tilde{\mathbb{P}}^\dagger(s' | s, a) \left(\max_{a'} \{X_{h+1}(s', a')\} \right).$$

660 *In addition, we define $X_h(s^\dagger, a) = 0$ for any $h \in [H], a \in \mathcal{A}$. Here, the notation $\sum_{s'}$ above means*
 661 *summation over $s' \in \mathcal{S} \cup \{s^\dagger\}$. But since $X_h(s^\dagger, a) = 0$ for any $h \in [H], a \in \mathcal{A}$, this is equivalent*
 662 *to the summation over $s' \in \mathcal{S}$. $\tilde{\mathbb{P}}^\dagger$ is the transition probability of fine-estimated sparsified MDP*
 663 *defined in section A.1.1.*

664 Then, for the reader to better understand the proof, we define some intermediate quantity, which also
 665 measure the uncertainty of value estimation, but they may be reward or policy dependent.

666 **Definition A.7** (Intermediate uncertainty function). *We define for any $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$ and*
 667 *any deterministic policy π ,*

$$\begin{aligned} W_{H+1}^\pi(s, a, r) &:= 0, \\ W_h^\pi(s, a, r) &:= \min \left\{ H - h + 1, \sqrt{\frac{8}{H^2} \bar{\phi}(m(s, a)) \cdot \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s, a)} \left(V_{h+1}(s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi) \right)} \right. \\ &\quad \left. + 9H\phi(m(s, a)) + \left(1 + \frac{1}{H} \right) \left(\tilde{\mathbb{P}}^\dagger \pi_{h+1} \right) W_{h+1}^\pi(s, a, r) \right\}, \end{aligned}$$

668 *where*

$$\left(\tilde{\mathbb{P}}^\dagger \pi_{h+1} \right) W_{h+1}^\pi(s, a, r) := \sum_{s'} \sum_{a'} \tilde{\mathbb{P}}^\dagger(s' | s, a) \pi_{h+1}(a' | s') W_{h+1}^\pi(s', a', r).$$

669 *In addition, we define $W_h^\pi(s^\dagger, a, r) = 0$ for any $h \in [H], a \in \mathcal{A}$ and any reward function r . Here,*
 670 *$\tilde{\mathbb{P}}^\dagger$ is the transition probability of fine-estimated sparsified MDP defined in section A.1.1.*

671

672 The intermediate function W is reward and policy dependent and can be used to upper bound the
 673 value estimation error. Further, we need to define some policy-dependent version of the uncertainty
 674 function, denoted as $X_h^\pi(s, a)$, as well as another quantity $Y_h^\pi(s, a, r)$.

675 **Definition A.8.** *We define $X_{H+1}^\pi(s, a) = Y_{H+1}^\pi(s, a, r) = 0$ for any π, r, s, a and*

$$X_h^\pi(s, a) := \min \left\{ H - h + 1, 9H\phi(m(s, a)) + \left(1 + \frac{1}{H} \right) \left(\tilde{\mathbb{P}}^\dagger \pi_{h+1} \right) X_{h+1}^\pi(s, a) \right\};$$

$$Y_h^\pi(s, a, r) := \sqrt{\frac{8}{H^2} \bar{\phi}(m(s, a)) \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s, a)} \left(V_{h+1}(s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi) \right)} + \left(1 + \frac{1}{H} \right) \left(\tilde{\mathbb{P}}^\dagger \pi_{h+1} \right) Y_{h+1}^\pi(s, a, r)$$

676 where

$$\begin{aligned} \left(\tilde{\mathbb{P}}^\dagger \pi_{h+1} \right) X_{h+1}^\pi(s, a) &= \sum_{s'} \sum_{a'} \tilde{\mathbb{P}}^\dagger(s' | s, a) \pi_{h+1}(a' | s') X_{h+1}^\pi(s', a'); \\ \left(\tilde{\mathbb{P}}^\dagger \pi_{h+1} \right) Y_{h+1}^\pi(s, a, r) &= \sum_{s'} \sum_{a'} \tilde{\mathbb{P}}^\dagger(s' | s, a) \pi_{h+1}(a' | s') Y_{h+1}^\pi(s', a', r). \end{aligned}$$

677 In addition, we define $X_h^\pi(s^\dagger, a) = Y_h^\pi(s^\dagger, a, r) = 0$ for any $a \in \mathcal{A}$, $h \in [H]$ and any reward r , any
678 policy π . Here, $\tilde{\mathbb{P}}^\dagger$ is the transition probability of fine-estimated sparsified MDP defined in section
679 A.1.1.

680

681 Finally, we define the empirical uncertainty functions.

682 **Definition A.9.** We define $U_{H+1}^k(s, a) = 0$ for any $k \in [K_{ucb}]$ and $s \in \mathcal{S}$, $a \in \mathcal{A}$. Further, we define

$$U_h^k(s, a) = H \min \{1, \phi(n^k(s, a))\} + \hat{\mathbb{P}}^\dagger(s, a)^\top \left(\max_{a'} U_{h+1}^k(\cdot, a') \right), \quad (\text{A.6})$$

683 where ϕ is the bonus function defined in eq. (A.4) and $n^k(s, a)$ is the counter of the times we encounter
684 (s, a) until the beginning of the k -th episode when running the sub-routine RF-UCB in the offline
685 phase.

686 A.2 Proof of the key theorems and lemmas

687 A.2.1 Uncertainty Functions upper bounds the estimation error

688 *Proof.* This proof follows closely from the techniques in [Ménard et al., 2021], but here we consider
689 the homogeneous MDP. We let $\tilde{d}_{\pi, h}^\dagger(s, a)$ be the probability of encountering (s, a) at stage h when
690 running policy π on $\tilde{\mathbb{P}}^\dagger$, starting from s_1 . Now we assume \mathcal{E} happens and fix a reward function r and
691 policy π . From lemma A.14, we know

$$\begin{aligned} \left| V_1(s_1; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi) - V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi) \right| &\leq \left| Q_h(s_1, \pi_1(a_1); \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi) - Q_h(s_1, \pi_1(s_1); \mathbb{P}^\dagger, r^\dagger, \pi) \right| \\ &\leq W_h^\pi(s_1, \pi(s_1), r), \end{aligned}$$

692 where $W_h^\pi(s, a, r)$ is the intermediate uncertainty function defined in definition A.7. Here, we use
693 the policy-dependent version of uncertainty function $W^\pi(s, a, r)$, which will then upper bounded
694 using another two policy-dependent quantities $X_h^\pi(s, a)$ and $Y_h^\pi(s, a, r)$. We define these policy-
695 dependent uncertainty functions to upper bound the estimation error of specific policy, but since we
696 are considering a reward-free setting, which entails a low estimation error for any policy and reward,
697 these quantities cannot be directly used in the algorithm to update the policy. Next, we claim

$$W_h^\pi(s, a, r) \leq X_h^\pi(s, a) + Y_h^\pi(s, a, r),$$

698 where $X_h^\pi(s, a)$ and $Y_h^\pi(s, a, r)$ are defined in definition A.8. Actually, this is easy from the definition
699 of $X_h^\pi(s, a)$, $Y_h^\pi(s, a, r)$ and $W_h^\pi(s, a, r)$. For $h = H + 1$, this is trivial. Assume this is true for
700 $h + 1$, then the case of h is given by the definition and the fact that $\min\{x, y + z\} \leq \min\{x, y\} + z$
701 for any $x, y, z \geq 0$. Therefore, we have

$$\left| V_1(s_1; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi) - V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi) \right| \leq X_1^\pi(s_1, \pi_1(s_1)) + Y_1^\pi(s_1, \pi_1(s_1), r). \quad (\text{A.7})$$

702 Next, we eliminate the dependency to policy π and obtain an upper bound of estimation error using
703 the policy-independent uncertainty function $X_h(s, a)$. This is done by bounding $Y_1^\pi(s_1, \pi_1(s_1), r)$ by
704 $X_h^\pi(s, a)$ and then upper bounding $X_h^\pi(s, a)$ by $X_h(s, a)$. From definition A.8, the Cauchy-Schwarz
705 inequality and the fact that $r \in [0, 1]$, we have for any reward function and any deterministic policy
706 π ,

$$Y_1^\pi(s_1, \pi_1(s_1), r)$$

$$\begin{aligned}
&\leq \sum_{s,a} \sum_{h=1}^H \tilde{d}_{\pi,h}^\dagger(s,a) \left(1 + \frac{1}{H}\right)^{h-1} \sqrt{\frac{8}{H^2} \bar{\phi}(m(s,a)) \cdot \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s,a)} \left(V_{h+1}(s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi)\right)} \\
&\hspace{15em} \text{(Induction by successively using the definition of } Y) \\
&\leq \frac{e}{H} \sqrt{8 \sum_{s,a} \sum_{h=1}^H \tilde{d}_{\pi,h}^\dagger(s,a) \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s,a)} \left(V_{h+1}(s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi)\right)} \cdot \sqrt{\sum_{s,a} \sum_{h=1}^H \tilde{d}_{\pi,h}^\dagger(s,a) \bar{\phi}(m(s,a))} \\
&\hspace{15em} \text{(Cauchy-Schwarz Inequality)} \\
&\leq \frac{e}{H} \sqrt{8H^2 \sum_{s,a} \sum_{h=1}^H \tilde{d}_{\pi,h}^\dagger(s,a)} \cdot \sqrt{\sum_{s,a} \sum_{h=1}^H \tilde{d}_{\pi,h}^\dagger(s,a) \bar{\phi}(m(s,a))} \\
&\leq e \sqrt{8 \sum_{s,a} \sum_{h=1}^H \tilde{d}_{\pi,h}^\dagger(s,a) \bar{\phi}(m(s,a))}.
\end{aligned}$$

707 We notice that the right hand side of the last inequality above is the policy value of some specific
708 reward function when running π on $\tilde{\mathbb{P}}$. Concretely, if the transition probability is $\tilde{\mathbb{P}}$ and
709 the reward function at (s, a) is $\bar{\phi}(m(s, a))$, then the state value function at the initial state s_1
710 is $\sum_{s,a} \sum_{h=1}^H \tilde{d}_{\pi,h}^\dagger(s, a) \bar{\phi}(m(s, a))$. This specific reward function is non-negative and uniformly
711 bounded by one, so it holds that $\sum_{s,a} \sum_{i=h}^H \tilde{d}_{\pi,i}^\dagger(s, a) \bar{\phi}(m(s, a)) \leq H - h + 1$. Moreover, from the
712 definition of ϕ and $\bar{\phi}$ (eq. (A.4)), we know $\bar{\phi}(m(s, a)) \leq H\phi(m(s, a))$. Then, from the definition
713 of $X_h(s, a)$ in definition A.6, we apply an inductive argument to obtain

$$\sum_{s,a} \sum_{h=1}^H \tilde{d}_{\pi,h}^\dagger(s, a) \bar{\phi}(m(s, a)) \leq X_1^\pi(s_1, \pi(s_1))$$

714 for any deterministic policy π . So we have

$$Y_1^\pi(s_1, \pi_1(s_1), r) \leq 2\sqrt{2}e \sqrt{X_1^\pi(s_1, \pi(s_1))}.$$

715 Therefore, combining the last inequality with (A.7), we have

$$\left| V_1(s_1; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi) - V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi) \right| \leq X_1^\pi(s_1, \pi_1(s_1)) + 2\sqrt{2}e \sqrt{X_1^\pi(s_1, \pi_1(s_1))}.$$

716 From the definition of $X_h(s, a)$ (definition A.6) and $X_h^\pi(s, a)$ (definition A.8), we can see

$$X_h^\pi(s, a) \leq X_h(s, a)$$

717 for any (h, s, a) and any deterministic policy π . Therefore, we conclude that

$$\begin{aligned}
\left| V_1(s_1; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi) - V_1(s_1; \mathbb{P}^\dagger, r^\dagger, \pi) \right| &\leq X_1(s_1, \pi_1(s_1)) + 2\sqrt{2}e \sqrt{X_1(s_1, \pi_1(s_1))} \\
&\leq \max_a X_1(s_1, a) + 2\sqrt{2}e \sqrt{\max_a X_1(s_1, a)}.
\end{aligned}$$

718

□

719 A.2.2 Proof of lemma 6.2

720 *Proof.* It suffices to prove that for any $k \in [K]$, $h \in [H]$, $s \in \mathcal{S}$, $a \in \mathcal{A}$, it holds that

$$X_h(s, a) \leq C \left(1 + \frac{1}{H}\right)^{3(H-h)} U_h^k(s, a)$$

721 since the theorem comes from an average of the inequalities above and the fact that $(1 + 1/H)^H \leq e$.
722 For any fixed k , we prove it by induction on h . When $h = H + 1$, both sides are zero by definition.
723 Suppose the claim holds for $h + 1$; we will prove that it also holds for h . We denote K_{ucb} as the

724 number of virtual episodes in the offline phase. From lemma A.16, the decreasing property of
 725 $\phi(\cdot)$ (lemma A.17) and lemma A.21, we have

$$\begin{aligned} X_h(s, a) &\leq CH\phi(n^{K_{ucb}}(s, a)) + \left(1 + \frac{2}{H}\right) \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} \{X_{h+1}(\cdot, a')\}\right) \quad (\text{Lemma A.16}) \\ &\leq CH\phi(n^k(s, a)) + \left(1 + \frac{2}{H}\right) \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} \{X_{h+1}(\cdot, a')\}\right) \quad (\text{Lemma A.17}) \\ &\leq CH\phi(n^k(s, a)) + \left(1 + \frac{1}{H}\right)^3 \widehat{\mathbb{P}}^\dagger(\cdot | s, a)^\top \left(\max_{a'} \{X_{h+1}(\cdot, a')\}\right) \\ &\quad (\text{Lemma A.21 and } 1 + 2/H \leq (1 + 1/H)^2) \end{aligned}$$

726 The inductive hypothesis gives us for any s' ,

$$\left(\max_{a'} \{X_{h+1}(s', a')\}\right) \leq C \left(1 + \frac{1}{H}\right)^{3(H-h-1)} \left(\max_{a'} U_{h+1}^k(s', a')\right),$$

727 which implies

$$\begin{aligned} X_h(s, a) &\leq CH\phi(n^k(s, a)) + \left(1 + \frac{1}{H}\right)^{3(H-h)} \widehat{\mathbb{P}}^\dagger(\cdot | s, a)^\top \left(\max_{a'} \{U_{h+1}^k(\cdot, a')\}\right) \\ &\leq C \left(1 + \frac{1}{H}\right)^{3(H-h)} \left[H\phi(n^k(s, a)) + \widehat{\mathbb{P}}^\dagger(\cdot | s, a)^\top \left(\max_{a'} \{U_{h+1}^k(\cdot, a')\}\right) \right]. \quad (\text{A.8}) \end{aligned}$$

728 Again, by definition of $X_h(s, a)$, we have

$$X_h(s, a) \leq H - h + 1 \leq C \left(1 + \frac{1}{H}\right)^{3(H-h)} H. \quad (\text{A.9})$$

729 Combining (A.8) and (A.9), as well as the definition of U_h^k (definition A.9), we prove the case of
 730 stage h and we conclude the theorem by induction. \square

731 A.2.3 Upper bounding the empirical uncertainty function (lemma 6.3)

732 *Proof.* First, we are going to prove

$$\frac{1}{K} \sum_{k=1}^K U_1^k(s, a) \leq \frac{H^2 |\mathcal{S}| |\mathcal{A}|}{K} \left[\log \left(\frac{6H |\mathcal{S}| |\mathcal{A}|}{\delta} \right) + |\mathcal{S}| \log \left(e \left(1 + \frac{KH}{|\mathcal{S}|} \right) \right) \right] \log \left(1 + \frac{HK}{|\mathcal{S}| |\mathcal{A}|} \right). \quad (\text{A.10})$$

733 From the definition (see algorithm 2), we know an important property of uncertainty function is that
 734 π_h^k is greedy with respect to U_h^k . So we have

$$\begin{aligned} U_h^k(s, a) &= H \min \{1, \phi(n^k(s, a))\} + \widehat{\mathbb{P}}^\dagger(s, a)^\top \left(\max_{a'} U_{h+1}^k(\cdot, a')\right) \\ &= H \min \{1, \phi(n^k(s, a))\} + \sum_{s', a'} \widehat{\mathbb{P}}^\dagger(s' | s, a) \pi_{h+1}^k(a' | s') U_{h+1}^k(s', a'). \end{aligned}$$

735 Therefore, we have

$$\begin{aligned} &\frac{1}{K_{ucb}} \sum_{k=1}^{K_{ucb}} U_1^k(s_1, a) \\ &\leq \frac{1}{K_{ucb}} \sum_{k=1}^{K_{ucb}} \max_a U_1^k(s_1, a) \\ &\leq \frac{H}{K_{ucb}} \sum_{k=1}^{K_{ucb}} \sum_{h=1}^H \sum_{(s, a)} \widehat{d}_{\pi^k, h}^\dagger(s, a) \min \{1, \phi(n^k(s, a))\} \quad (\text{definition A.9}) \end{aligned}$$

$$\leq \frac{4H^2}{K_{ucb}} \left[\log \left(\frac{6H|\mathcal{S}||\mathcal{A}|}{\delta} \right) + |\mathcal{S}| \log \left(e \left(1 + \frac{K_{ucb}H}{|\mathcal{S}|} \right) \right) \right] \sum_{h=1}^H \sum_{k=1}^{K_{ucb}} \sum_{(s,a)} \frac{\widehat{d}_{\pi^k, h}^\dagger(s, a)}{\max \left\{ 1, \sum_{i=1}^{k-1} \sum_{h \in [H]} \widehat{d}_{\pi^i, h}^\dagger(s, a) \right\}} \quad (\text{lemma A.20})$$

736 We apply Lemma D.7 and Jensen's Inequality to obtain

$$\begin{aligned} & \sum_{h \in [H]} \sum_{k=1}^{K_{ucb}} \sum_{(s,a)} \widehat{d}_{\pi^k, h}^\dagger(s, a) \frac{1}{\max \left\{ 1, \sum_{i=1}^{k-1} \sum_{h \in [H]} \widehat{d}_{\pi^i, h}^\dagger(s, a) \right\}} \\ & \leq 4 \sum_{(s,a)} \log \left(1 + \sum_{i=1}^{K_{ucb}} \sum_{h \in [H]} \widehat{d}_{\pi^i, h}^\dagger(s, a) \right) \quad (\text{lemma D.7}) \\ & \leq 4 |\mathcal{S}| |\mathcal{A}| \log \left(1 + \frac{1}{|\mathcal{S}| |\mathcal{A}|} \sum_{i=1}^{K_{ucb}} \sum_{h \in [H]} \sum_{(s,a)} \widehat{d}_{\pi^i, h}^\dagger(s, a) \right) \\ & \leq 4 |\mathcal{S}| |\mathcal{A}| \log \left(1 + \frac{HK_{ucb}}{|\mathcal{S}| |\mathcal{A}|} \right) \end{aligned}$$

737 Inserting the last display to the upper bound, we conclude the proof of the upper bound on U .

738 Then, to prove the sample complexity, we use lemma D.8. We take

$$B = \frac{64H^2 |\mathcal{A}|}{\varepsilon^2} \left[\log \left(\frac{6H|\mathcal{S}||\mathcal{A}|}{\delta} \right) + |\mathcal{S}| \right] \text{ and } x = \frac{K_{ucb}}{|\mathcal{S}|}$$

739 in Lemma D.8. From lemma D.8, we know there exists a universal constant C_1 such that when
740 $x \geq C_1 B \log(B)^2$, it holds that $B \log(1+x) (1 + \log(1+x)) \leq x$, which implies whenever

$$K_{ucb} \geq \frac{64H^2 |\mathcal{S}| |\mathcal{A}|}{\varepsilon^2} \left[\log \left(\frac{6H|\mathcal{S}||\mathcal{A}|}{\delta} \right) + |\mathcal{S}| \right],$$

741 it holds that

$$\frac{64H^2 |\mathcal{S}| |\mathcal{A}|}{K_{ucb}} \left[\log \left(\frac{6H|\mathcal{S}||\mathcal{A}|}{\delta} \right) + |\mathcal{S}| \right] \log \left(1 + \frac{K_{ucb}}{|\mathcal{S}|} \right) \left[1 + \log \left(1 + \frac{K_{ucb}}{|\mathcal{S}|} \right) \right] \leq \varepsilon^2. \quad (\text{A.11})$$

742 For notation simplicity, we define

$$L = 16H^2 |\mathcal{S}| |\mathcal{A}| \left[\log \left(\frac{6H|\mathcal{S}||\mathcal{A}|}{\delta} \right) + |\mathcal{S}| \right].$$

743 Comparing the left hand side of (A.11) and the right hand side of (A.10), we know

$$\begin{aligned} & \text{Right hand side of (A.10)} \\ & = \frac{16H^2 |\mathcal{S}| |\mathcal{A}|}{K_{ucb}} \log \left(1 + \frac{HK_{ucb}}{|\mathcal{S}| |\mathcal{A}|} \right) \left[\log \left(\frac{6H|\mathcal{S}||\mathcal{A}|}{\delta} \right) + |\mathcal{S}| + |\mathcal{S}| \log \left(1 + \frac{K_{ucb}H}{|\mathcal{S}|} \right) \right] \\ & \leq \frac{16H^2 |\mathcal{S}| |\mathcal{A}|}{K_{ucb}} \left[\log(H) + \log \left(1 + \frac{K_{ucb}}{|\mathcal{S}|} \right) \right] \left[\log \left(\frac{6H|\mathcal{S}||\mathcal{A}|}{\delta} \right) + |\mathcal{S}| + |\mathcal{S}| \log \left(1 + \frac{K_{ucb}H}{|\mathcal{S}|} \right) \right] \\ & \leq \frac{16H^2 |\mathcal{S}| |\mathcal{A}|}{K_{ucb}} \left[\log(H) + \log \left(1 + \frac{K_{ucb}}{|\mathcal{S}|} \right) \right] \left[\log \left(\frac{6H|\mathcal{S}||\mathcal{A}|}{\delta} \right) + |\mathcal{S}| \right] \left[1 + \log(H) + \log \left(1 + \frac{K_{ucb}}{|\mathcal{S}|} \right) \right] \\ & = \frac{L}{K_{ucb}} \left[\log(H) + \log \left(1 + \frac{K_{ucb}}{|\mathcal{S}|} \right) \right] \left[1 + \log(H) + \log \left(1 + \frac{K_{ucb}}{|\mathcal{S}|} \right) \right]. \end{aligned}$$

744 It is easy to see that $\log(H) \leq \log(1 + K_{ucb}/|\mathcal{S}|)$, so the right hand side of last inequality is upper
745 bounded by

$$\frac{4L}{K_{ucb}} \log \left(1 + \frac{K_{ucb}}{|\mathcal{S}|} \right) \left[1 + \log \left(1 + \frac{K_{ucb}}{|\mathcal{S}|} \right) \right].$$

746 From (A.11), we know this is upper bounded by ε^2 as long as $K_{ucb} \geq 4L/\varepsilon^2$. From the definition of
 747 L and equations (A.10) (A.11), we have when $K_{ucb} \geq 4L/\varepsilon^2$, it holds that

$$\frac{1}{K_{ucb}} \sum_{k=1}^{K_{ucb}} U_1^k(s_1, a) \leq \text{right hand side of (A.10)} \leq \text{left hand side of (A.11)} \leq \varepsilon^2.$$

748 In conclusion, this admits a sample complexity of

$$K_{ucb} = \frac{CH^2 |\mathcal{S}| |\mathcal{A}|}{\varepsilon^2} \left(\iota + |\mathcal{S}| \right) \text{polylog} \left(|\mathcal{S}|, |\mathcal{A}|, H, \frac{1}{\varepsilon}, \frac{1}{\delta} \right).$$

749

□

750 A.3 Omitted proofs

751 A.3.1 Proof for lemma A.4 (high probability event)

752 The lemma A.4 is proved by combining all the lemmas below via a union bound. Below, we always
 753 denote $N(s, a, s')$ as the number of (s, a, s') in the offline dataset.

754 **Lemma A.10.** \mathcal{E}^P holds with probability at least $1 - \delta/6$.

755 *Proof.* For any fixed (s, a, s') such that $N(s, a, s') \geq \Phi$, we denote $n_i(s, a)$ as the index of the i -th
 756 time when we visit (s, a) . For notation simplicity, we fix the state-action pair (s, a) here and write
 757 $n_i(s, a)$ as n_i simply for $i = 1, 2, \dots, N(s, a)$. Notice that the total visiting time $N(s, a)$ is random,
 758 so our argument is based on conditional probability. We denote $X_i = \mathbb{I}(s'_{n_i} = s')$ as the indicator
 759 of whether the next state is s' when we visited (s, a) at the i -th time. From the data generating
 760 mechanism and the definition for the reference MDP, we know conditional on the total number
 761 of visiting $N(s, a)$ and all the indexes $n_1 \leq \dots \leq n_{N(s, a)}$, it holds that $X_1, X_2, \dots, X_{N(s, a)}$ are
 762 independent Bernoulli random variable with successful probability being $\mathbb{P}^\dagger(s' | s, a)$. We denote \bar{X}
 763 as their arithmetic average. Using Empirical Bernstein Inequality (Lemma D.3), and the fact that
 764 $\frac{1}{N(s, a)} \sum_{i=1}^{N(s, a)} (X_i - \bar{X})^2 \leq \frac{1}{N(s, a)} \sum_{i=1}^{N(s, a)} X_i^2 = \frac{1}{N(s, a)} \sum_{i=1}^{N(s, a)} X_i = \widehat{\mathbb{P}}^\dagger(s' | s, a)$, we have

$$\mathbb{P}(\mathcal{E}^P | n_i (1 \leq i \leq N(s, a)), N(s, a) = n) \geq 1 - \delta/6.$$

765 Taking integral of the conditional probability, we conclude that the unconditional probability of the
 766 event \mathcal{E}^P is at least $1 - \delta/6$. Therefore, we conclude. □

767 **Lemma A.11.** For $\delta > 0$, it holds that $\mathbb{P}(\mathcal{E}^4) \geq 1 - \delta/6$.

768 *Proof.* Consider for a fixed triple $(s, a) \in \mathcal{S} \times \mathcal{A}$. Let $m(s, a)$ denote the number of times the tuple
 769 (s, a) was encountered in total during the online interaction phase. We define X_i as follows. For $i \leq$
 770 $m(s, a)$, we let X_i be the subsequent state s' when (s, a) was encountered the i -th time in the whole
 771 run of the algorithm. Otherwise, we let X_i be an independent sample from $\mathbb{P}^\dagger(s, a)$. By construction,
 772 $\{X_i\}$ is a sequence of i.i.d. categorical random variables from \mathcal{S} with distribution $\mathbb{P}^\dagger(\cdot | s, a)$. We
 773 denote $\widetilde{\mathbb{P}}_X^{\dagger, i} = \frac{1}{i} \sum_{j=1}^i \delta_{X_j}$ as the empirical probability mass and \mathbb{P}_X^\dagger as the probability mass of X_i .
 774 Then, from Lemma D.4, we have with probability at least $1 - \delta/6$, it holds that for any $i \in \mathbb{N}$,

$$\text{KL} \left(\widetilde{\mathbb{P}}_X^{\dagger, i}; \mathbb{P}_X^\dagger \right) \leq \frac{1}{i} \left[\log \left(\frac{6}{\delta} \right) + |\mathcal{S}| \log \left(e \left(1 + \frac{i}{|\mathcal{S}|} \right) \right) \right],$$

775 which implies

$$\text{KL} \left(\widetilde{\mathbb{P}}_X^{\dagger}; \mathbb{P}^\dagger(s, a) \right) \leq \frac{1}{m(s, a)} \left[\log \left(\frac{6}{\delta} \right) + |\mathcal{S}| \log \left(e \left(1 + \frac{m(s, a)}{|\mathcal{S}|} \right) \right) \right].$$

776 Using a union bound for $(s, a) \in \mathcal{S} \times \mathcal{A}$, we conclude. □

777

778 **Lemma A.12** (Lower Bound on Counters). For $\delta > 0$, it holds that $\mathbb{P}(\mathcal{E}^2) \geq 1 - \delta/6$ and
 779 $\mathbb{P}(\mathcal{E}^5) \geq 1 - \delta/6$.

780 *Proof.* We denote n_h^k as the number of times we encounter (s, a) at stage h before the beginning of
 781 episode k . Concretely speaking, we define $n_h^1(s, a) = 0$ for any (h, s, a) . Then, we define

$$n_h^k(s, a) = n_h^k(s, a) + \mathbb{I}[(s, a) = (s_h^k, a_h^k)].$$

782 We fixed an $(s, a, h) \in [H] \times \mathcal{S} \times \mathcal{A}$ and denote \mathcal{F}_k as the sigma field generated by the first $k - 1$
 783 episodes when running RF-UCB and $X_k = \mathbb{I}[(s_h^k, a_h^k) = (s, a)]$. Then, we know X_k is \mathcal{F}_{k+1} -
 784 measurable and $\mathbb{E}[X_k | \mathcal{F}_k] = \hat{d}_{\pi^k, h}^\dagger(s, a)$ is \mathcal{F}_k -measurable. Taking $W = \ln(6/\delta)$ in Lemma D.1
 785 and applying a union bound, we know with probability $1 - \delta/6$, the following event happens:

$$\forall k \in \mathbb{N}, (h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}, \quad n_h^k(s, a) \geq \frac{1}{2} \sum_{i < k} \hat{d}_{\pi^i, h}^\dagger(s, a) - \ln \left(\frac{6H |\mathcal{S}| |\mathcal{A}|}{\delta} \right).$$

786 To finish the proof, it remains to note that the event above implies the event we want by summing
 787 over $h \in [H]$ for each $k \in \mathbb{N}$ and each $(s, a) \in \mathcal{S} \times \mathcal{A}$. For the \mathcal{E}^5 , the proof is almost the same. \square

788 **Lemma A.13** (Upper Bound on Counters). For $\delta > 0$, it holds that $\mathbb{P}(\mathcal{E}^3) \geq 1 - \delta/6$.

789 *Proof.* We denote n_h^k as the number of times we encounter (s, a) at stage h before the beginning of
 790 episode k . Concretely speaking, we define $n_h^1(s, a) = 0$ for any (h, s, a) . Then, we define

$$n_h^k(s, a) = n_h^k(s, a) + \mathbb{I}[(s, a) = (s_h^k, a_h^k)].$$

791 We fixed an $(s, a, h) \in [H] \times \mathcal{S} \times \mathcal{A}$ and denote \mathcal{F}_k as the sigma field generated by the first $k - 1$
 792 episodes when running RF-UCB and $X_k = \mathbb{I}[(s_h^k, a_h^k) = (s, a)]$. Then, we know X_k is \mathcal{F}_{k+1} -
 793 measurable and $\mathbb{E}[X_k | \mathcal{F}_k] = \hat{d}_{\pi^k, h}^\dagger(s, a)$ is \mathcal{F}_k -measurable. Taking $W = \ln(6/\delta)$ in Lemma D.2
 794 and applying a union bound, we know with probability $1 - \delta/6$, the following event happens:

$$\forall k \in \mathbb{N}, (h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}, \quad n_h^k(s, a) \leq 2 \sum_{i < k} \hat{d}_{\pi^i, h}^\dagger(s, a) + \ln \left(\frac{6H |\mathcal{S}| |\mathcal{A}|}{\delta} \right).$$

795 To finish the proof, it remains to note that the event above implies the event we want by summing
 796 over $h \in [H]$ for each $k \in \mathbb{N}$ and each $(s, a) \in \mathcal{S} \times \mathcal{A}$. \square

797 A.3.2 Proof for lemma A.14 (property of intermediate uncertainty function)

798 **Lemma A.14** (Intermediate Uncertainty Function). We define for any $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$ and
 799 any deterministic policy π ,

$$W_{H+1}^\pi(s, a, r) := 0,$$

$$W_h^\pi(s, a, r) := \min \left\{ H - h + 1, \sqrt{\frac{8}{H^2} \bar{\phi}(m(s, a)) \cdot \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s, a)} \left(V_{h+1} \left(s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi \right) \right)} \right. \\ \left. + 9H\phi(m(s, a)) + \left(1 + \frac{1}{H} \right) \left(\tilde{\mathbb{P}}^\dagger \pi_{h+1} \right) W_{h+1}^\pi(s, a, r) \right\},$$

800 where

$$\left(\tilde{\mathbb{P}}^\dagger \pi_{h+1} \right) W_{h+1}^\pi(s, a, r) := \sum_{s'} \sum_{a'} \tilde{\mathbb{P}}^\dagger(s' | s, a) \pi_{h+1}(a' | s') W_{h+1}^\pi(s', a', r).$$

801 In addition, we define $W_h^\pi(s^\dagger, a, r) = 0$ for any $h \in [H]$, $a \in \mathcal{A}$ and any reward function r . Then,
 802 under \mathcal{E} , for any $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$, any deterministic policy π , and any deterministic reward
 803 function r (with its augmentation r^\dagger), it holds that

$$\left| Q_h \left(s, a; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi \right) - Q_h \left(s, a; \mathbb{P}^\dagger, r^\dagger, \pi \right) \right| \leq W_h^\pi(s, a, r).$$

804 *Proof.* We prove by induction. For $h = H + 1$, we know $W_{H+1}(s, a, r) = 0$ by definition and the
805 left hand side of the inequality we want also vanishes. Suppose the claim holds for $h + 1$ and for any
806 s, a . The Bellman equation gives us

$$\begin{aligned} \Delta_h &:= Q_h \left(s, a; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi \right) - Q_h \left(s, a; \mathbb{P}^\dagger, r^\dagger, \pi \right) \leq \underbrace{\sum_{s'} \left(\tilde{\mathbb{P}}^\dagger (s' | s, a) - \mathbb{P}^\dagger (s' | s, a) \right) V_{h+1} (s'; \mathbb{P}^\dagger, r^\dagger, \pi)}_I \\ &+ \underbrace{\sum_{s'} \tilde{\mathbb{P}}^\dagger (s' | s, a) \left| V_{h+1} (s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi) - V_{h+1} (s'; \mathbb{P}^\dagger, r^\dagger, \pi) \right|}_II \end{aligned}$$

807 Under \mathcal{E} , we know that $\text{KL} \left(\tilde{\mathbb{P}}^\dagger; \mathbb{P}^\dagger \right) \leq \frac{1}{H} \phi (m(s, a))$ for any (s, a) , so from Lemma D.5 we have

$$\begin{aligned} |I| &\leq \sqrt{\frac{2}{H} \phi (m(s, a)) \cdot \text{Var}_{s' \sim \mathbb{P}^\dagger(s, a)} (V_{h+1} (s'; \mathbb{P}^\dagger, r^\dagger, \pi))} + \frac{2}{3} (H - h) \frac{\phi (m(s, a))}{H} \\ &\leq \sqrt{\frac{2}{H} \phi (m(s, a)) \cdot \text{Var}_{s' \sim \mathbb{P}^\dagger(s, a)} (V_{h+1} (s'; \mathbb{P}^\dagger, r^\dagger, \pi))} + \frac{2}{3} \phi (m(s, a)), \end{aligned}$$

808 where ϕ is the bonus defined as (A.4). Here, we let f in Lemma D.5 be $V_{h+1} (s'; \mathbb{P}^\dagger, r^\dagger, \pi)$. So the
809 range will be $H - h$ and the upper bound for KL divergence is given by high-probability event \mathcal{E} . We
810 further apply Lemma D.6 to get

$$\text{Var}_{s' \sim \mathbb{P}^\dagger(s, a)} (V_{h+1} (s'; \mathbb{P}^\dagger, r^\dagger, \pi)) \leq 2 \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s, a)} (V_{h+1} (s'; \mathbb{P}^\dagger, r^\dagger, \pi)) + 4H\phi (m(s, a))$$

811 and

$$\begin{aligned} \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s, a)} (V_{h+1} (s'; \mathbb{P}^\dagger, r^\dagger, \pi)) &\leq 2 \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s, a)} (V_{h+1} (s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi)) \\ &+ 2H \sum_{s'} \tilde{\mathbb{P}}^\dagger (s' | s, a) \left| V_{h+1} (s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi) - V_{h+1} (s'; \mathbb{P}^\dagger, r^\dagger, \pi) \right|. \end{aligned}$$

812 Therefore, we have

$$\begin{aligned} |I| &\leq \frac{2}{3} \phi (m(s, a)) + \left[\frac{8}{H} \phi (m(s, a)) \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s, a)} (V_{h+1} (s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi)) + 8\phi (m(s, a))^2 \right. \\ &\left. + 8\phi (m(s, a)) \sum_{s'} \tilde{\mathbb{P}}^\dagger (s' | s, a) \left| V_{h+1} (s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi) - V_{h+1} (s'; \mathbb{P}^\dagger, r^\dagger, \pi) \right| \right]^{1/2}. \end{aligned}$$

813 Using the fact that $\sqrt{x + y} \leq \sqrt{x} + \sqrt{y}$ and $\sqrt{xy} \leq \frac{x+y}{2}$ for positive x, y , and the definition of II ,
814 we obtain

$$|I| \leq \sqrt{\frac{8}{H} \phi (m(s, a)) \cdot \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s, a)} (V_{h+1} (s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi))} + 6H\phi (m(s, a)) + \frac{1}{H} |II|,$$

815 which implies

$$|\Delta_h| \leq \sqrt{\frac{8}{H} \phi (m(s, a)) \cdot \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s, a)} (V_{h+1} (s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi))} + 6H\phi (m(s, a)) + \left(1 + \frac{1}{H} \right) |II|.$$

816 If $H\phi (m(s, a)) \leq 1$, then we have by definition

$$\begin{aligned} |\Delta_h| &\leq \sqrt{\frac{8}{H^2} \cdot H\phi (m(s, a)) \cdot \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s, a)} (V_{h+1} (s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi))} + 6H\phi (m(s, a)) + \left(1 + \frac{1}{H} \right) |II| \\ &\leq \sqrt{\frac{8}{H^2} \bar{\phi} (m(s, a)) \cdot \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s, a)} (V_{h+1} (s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi))} + 6H\phi (m(s, a)) + \left(1 + \frac{1}{H} \right) |II| \end{aligned}$$

817 by definition. Otherwise, if $H\phi (m(s, a)) \geq 1$, we have

$$\sqrt{\frac{8}{H} \phi (m(s, a)) \cdot \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s, a)} (V_{h+1} (s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi))} \leq \sqrt{8H\phi (m(s, a))} \leq 3H\phi (m(s, a)).$$

818 Therefore, we have in either case

$$|\Delta_h| \leq \sqrt{\frac{8}{H^2} \bar{\phi}(m(s, a)) \cdot \text{Var}_{s' \sim \tilde{\mathbb{P}}^\dagger(s, a)} \left(V_{h+1} \left(s'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi \right) \right) + 9H\phi(m(s, a))} + \left(1 + \frac{1}{H} \right) |II|.$$

819 The induction hypothesis gives us that

$$\begin{aligned} |II| &\leq \sum_{s'} \sum_{a'} \tilde{\mathbb{P}}^\dagger(s' | s, a) \pi_{h+1}(a' | s') \left| Q_{h+1} \left(s', a'; \tilde{\mathbb{P}}^\dagger, r^\dagger, \pi \right) - Q_{h+1} \left(s', a'; \mathbb{P}^\dagger, r^\dagger, \pi \right) \right| \\ &\leq \left(\tilde{\mathbb{P}}^\dagger \pi_{h+1} \right) W_{h+1}^\pi(s, a, r). \end{aligned}$$

820 Inserting this into the upper bound of Δ_h , we prove the case of h by the definition of $W_h^\pi(s, a, r)$
821 and the simple fact that $|\Delta_h| \leq H - h + 1$. \square

822 A.3.3 Proof for lemma A.15 and lemma A.16 (properties of uncertainty function)

823 In this section, we prove some lemmas for upper bounding the uncertainty functions $X_h(s, a)$. We
824 first provide a basic upper bound for $X_h(s, a)$. The uncertainty function is defined in definition A.6.

825 **Lemma A.15.** *Under the high probability event \mathcal{E} (defined in appendix A.1.2) for all $(h, s, a) \in$
826 $[H] \times \mathcal{S} \times \mathcal{A}$, it holds that*

$$X_h(s, a) \leq 11H \min \{1, \phi(m(s, a))\} + \left(1 + \frac{2}{H} \right) \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} X_{h+1}(\cdot, a') \right).$$

827 In addition, for $s = s^\dagger$, from definition, we know the above upper bound naturally holds.

828 *Proof.* For any fixed $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$, from the definition of the uncertainty function, we
829 know

$$X_h(s, a) \leq 9H\phi(m(s, a)) + \left(1 + \frac{1}{H} \right) \tilde{\mathbb{P}}^\dagger(\cdot | s, a)^\top \left(\max_{a'} X_{h+1}(\cdot, a') \right).$$

830 To prove the lemma, it suffices to bound the difference between $\tilde{\mathbb{P}}^\dagger(\cdot | s, a)^\top \left(\max_{a'} X_{h+1}(\cdot, a') \right)$
831 and $\mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} X_{h+1}(\cdot, a') \right)$. Under \mathcal{E} (which happens with probability at least $1 - \delta$), it
832 holds that $\text{KL} \left(\tilde{\mathbb{P}}^\dagger(\cdot | s, a); \mathbb{P}^\dagger(\cdot | s, a) \right) \leq \frac{1}{H} \phi(m(s, a))$. Applying Lemma D.5 and a simple
833 bound for variance gives us

$$\begin{aligned} &\left| \tilde{\mathbb{P}}^\dagger(\cdot | s, a)^\top \left(\max_{a'} X_{h+1}(\cdot, a') \right) - \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} X_{h+1}(\cdot, a') \right) \right| \\ &\leq \sqrt{\frac{2}{H} \text{Var}_{s' \sim \mathbb{P}^\dagger(s, a)} \left(\max_{a'} X_{h+1}(s', a') \right) \phi(m(s, a))} + \frac{2}{3} \phi(m(s, a)) \\ &\leq \sqrt{\left[\frac{2}{H} \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} X_{h+1}(\cdot, a')^2 \right) \right] \phi(m(s, a))} + \frac{2}{3} \phi(m(s, a)) \\ &\leq \sqrt{\left[\frac{2}{H} \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} X_{h+1}(\cdot, a') \right) \right] \cdot H\phi(m(s, a))} + \frac{2}{3} \phi(m(s, a)) \\ &\leq \frac{1}{H} \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} X_{h+1}(\cdot, a') \right) + 2H\phi(m(s, a)). \end{aligned}$$

834 The last line comes from $\sqrt{ab} \leq \frac{a+b}{2}$ for any positive a, b , and the fact that $\frac{1}{2}H\phi(m(s, a)) +$
835 $\frac{2}{3}\phi(m(s, a)) \leq 2H\phi(m(s, a))$. Insert this bound into the definition of $X_h(s, a)$ to obtain

$$X_h(s, a) \leq 11H\phi(m(s, a)) + \left(1 + \frac{2}{H} \right) \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} X_{h+1}(\cdot, a') \right).$$

836 Noticing that $X_h(s, a) \leq H - h + 1 \leq 11H$, we conclude. \square

837

838 **Lemma A.16.** *There exists a universal constant $C \geq 1$ such that under the high probability event \mathcal{E} ,*
839 *when $3K_{ucb} \geq K_{de}$, we have for any $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$, it holds that*

$$X_h(s, a) \leq CH \frac{K_{ucb}}{K_{de}} \phi(n^{K_{ucb}}(s, a)) + \left(1 + \frac{2}{H}\right) \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} \{X_{h+1}(\cdot, a')\}\right).$$

840 *In addition, for $s = s^\dagger$, from definition, we know the above upper bound naturally holds.*

841 *Proof.* Here, the universal constant may vary from line to line. Under \mathcal{E} , we have

$$\begin{aligned} & X_h(s, a) \\ & \leq CH \phi(m(s, a)) + \left(1 + \frac{1}{H}\right) \tilde{\mathbb{P}}^\dagger(\cdot | s, a)^\top \left(\max_{a'} X_{h+1}(\cdot, a')\right). \end{aligned} \quad (\text{definition})$$

$$\leq CH \min\{1, \phi(m(s, a))\} + \left(1 + \frac{2}{H}\right) \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} \{X_{h+1}(\cdot, a')\}\right) \quad (\text{Lemma A.15})$$

$$\leq CH \phi\left(K_{de} \sum_{h \in [H]} w_h^{mix}(s, a)\right) + \left(1 + \frac{2}{H}\right) \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} \{X_{h+1}(\cdot, a')\}\right) \quad (\text{Lemma A.18})$$

$$\begin{aligned} & = CH \phi\left(\frac{K_{de}}{K_{ucb}} \sum_{k=1}^{K_{ucb}} \sum_{h \in [H]} d_{\pi^k, h}^\dagger(s, a)\right) + \left(1 + \frac{2}{H}\right) \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} \{X_{h+1}(\cdot, a')\}\right) \\ & \quad (\text{definition}) \end{aligned}$$

$$\leq CH \phi\left(\frac{K_{de}}{3K_{ucb}} \sum_{k=1}^{K_{ucb}} \sum_{h \in [H]} \hat{d}_{\pi^k, h}^\dagger(s, a)\right) + \left(1 + \frac{2}{H}\right) \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} \{X_{h+1}(\cdot, a')\}\right) \quad (\text{Lemma A.22 and A.17})$$

$$\leq CH \frac{K_{ucb}}{K_{de}} \phi\left(\sum_{k=1}^{K_{ucb}} \sum_{h \in [H]} \hat{d}_{\pi^k, h}^\dagger(s, a)\right) + \left(1 + \frac{2}{H}\right) \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} \{X_{h+1}(\cdot, a')\}\right) \quad (\text{Lemma A.17})$$

842 Since $X_h(s, a) \leq (H - h) \leq CH \frac{K_{ucb}}{K_{de}}$, we can modify the last display as

$$\begin{aligned} X_h(s, a) & \leq CH \frac{K_{ucb}}{K_{de}} \min\left\{1, \phi\left(\sum_{k=1}^{K_{ucb}} \sum_{h \in [H]} \hat{d}_{\pi^k, h}^\dagger(s, a)\right)\right\} + \left(1 + \frac{2}{H}\right) \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} \{X_{h+1}(\cdot, a')\}\right) \\ & \leq CH \frac{K_{ucb}}{K_{de}} \phi(n^{K_{ucb}}(s, a)) + \left(1 + \frac{2}{H}\right) \mathbb{P}^\dagger(\cdot | s, a)^\top \left(\max_{a'} \{X_{h+1}(\cdot, a')\}\right) \end{aligned} \quad (\text{Lemma A.19})$$

843 Therefore, we conclude. \square

844 **A.3.4 Proof for lemma A.17, lemma A.18, lemma A.19, lemma A.20 (properties of bonus** 845 **function)**

846 For our bonus function, we have the following basic property.

847 **Lemma A.17.** $\phi(x)$ is non-increasing when $x > 0$. For any $\alpha \leq 1$, we have $\phi(\alpha x) \leq \frac{1}{\alpha} \phi(x)$.

848 *Proof.* We define $f(x) := \frac{1}{x} [C + D \log(e(1 + \frac{x}{D}))]$, where $C, D \geq 1$. Then, for $x > 0$,

$$f'(x) = -\frac{C + D \log(e(1 + \frac{x}{D}))}{x^2} + \frac{D}{x(D+x)} \leq -\frac{C + D \log(1 + \frac{x}{D})}{x^2} \leq 0.$$

849 Taking $C = \log(6H|\mathcal{S}||\mathcal{A}|\delta)$ and $D = |\mathcal{S}|$, we conclude the first claim. For the second claim, it
850 is trivial since the logarithm function is increasing. \square

851

852 **Lemma A.18.** Under \mathcal{E} , we have for any $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\min \{1, \phi(m(s, a))\} \leq 4\phi \left(K_{de} \sum_{h=1}^H w_h^{mix}(s, a) \right).$$

Proof. For fixed $(s, a) \in \mathcal{S} \times \mathcal{A}$, when $H \ln(6H |\mathcal{S}| |\mathcal{A}| / \delta) \leq \frac{1}{4} \left(K_{de} \sum_{h \in [H]} w_h^{mix}(s, a) \right)$, we know that \mathcal{E}^5 implies $m(s, a) \geq \frac{1}{4} \sum_{h \in [H]} K_{de} w_h^{mix}(s, a)$. From Lemma A.17, we know

$$\phi(m(s, a)) \leq \phi \left(\frac{1}{4} \sum_{h \in [H]} K_{de} w_h^{mix}(s, a) \right) \leq 4\phi \left(\sum_{h \in [H]} K_{de} w_h^{mix}(s, a) \right).$$

When $H \ln(6H |\mathcal{S}| |\mathcal{A}| / \delta) > \frac{1}{4} \left(K_{de} \sum_{h \in [H]} w_h^{mix}(s, a) \right)$, simple algebra shows that

$$\min \{1, \phi(m(s, a))\} \leq 1 \leq \frac{4H \ln(6H |\mathcal{S}| |\mathcal{A}| / \delta)}{K_{de} \sum_{h \in [H]} w_h^{mix}(s, a)} \leq 4\phi \left(K_{de} \sum_{h \in [H]} w_h^{mix}(s, a) \right).$$

853 Therefore, we conclude. □

854

855 **Lemma A.19.** Under \mathcal{E} , we have for any $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\min \left\{ 1, \phi \left(\sum_{k=1}^{K_{ucb}} \sum_{h=1}^H \hat{d}_{\pi^k, h}^\dagger(s, a) \right) \right\} \leq 4\phi(n^{K_{ucb}}(s, a)).$$

856 *Proof.* We consider under the event \mathcal{E}^3 , it holds that for any $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\sum_{k=1}^{K_{ucb}} \sum_{h=1}^H \hat{d}_{\pi^k, h}^\dagger(s, a) \geq \frac{1}{2} n^{K_{ucb}}(s, a) - \frac{1}{2} H \ln \left(\frac{6H |\mathcal{S}| |\mathcal{A}|}{\delta} \right).$$

857 If $H \ln(6H |\mathcal{S}| |\mathcal{A}| / \delta) \leq \frac{1}{2} n^{K_{ucb}}(s, a)$, then we have $\sum_{k=1}^{K_{ucb}} \sum_{h=1}^H \hat{d}_{\pi^k, h}^\dagger(s, a) \geq \frac{1}{4} n^{K_{ucb}}(s, a)$,
858 which implies

$$\phi \left(\sum_{k=1}^{K_{ucb}} \sum_{h=1}^H \hat{d}_{\pi^k, h}^\dagger(s, a) \right) \leq \phi \left(\frac{1}{4} n^{K_{ucb}}(s, a) \right) \leq 4\phi(n^{K_{ucb}}(s, a)).$$

859 Otherwise, if $H \ln(6H |\mathcal{S}| |\mathcal{A}| / \delta) > \frac{1}{2} n^{K_{ucb}}(s, a)$, then we have

$$\min \left\{ 1, \phi \left(\sum_{k=1}^{K_{ucb}} \hat{d}_{\pi^k, h}^\dagger(s, a) \right) \right\} \leq 1 \leq \frac{2H \ln(6H |\mathcal{S}| |\mathcal{A}| / \delta)}{n^{K_{ucb}}(s, a)} \leq 4\phi(n^{K_{ucb}}(s, a)).$$

860 Combining the two cases above, we conclude. □

861

862 **Lemma A.20.** Under \mathcal{E} , we have for any $k \in [K_{ucb}]$ and any $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned} & \min \{1, \phi(n^k(s, a))\} \\ & \leq \frac{4H}{\max \left\{ 1, \sum_{i < k} \sum_{h \in [H]} \hat{d}_{\pi^i, h}^\dagger(s, a) \right\}} \left[\log \left(\frac{6H |\mathcal{S}| |\mathcal{A}|}{\delta} \right) + |\mathcal{S}| \log \left(e \left(1 + \frac{\sum_{i < k} \sum_{h \in [H]} \hat{d}_{\pi^i, h}^\dagger(s, a)}{|\mathcal{S}|} \right) \right) \right]. \end{aligned}$$

863 *Proof.* Under \mathcal{E}^2 , it holds that

$$n^k(s, a) \geq \frac{1}{2} \sum_{i < k} \sum_{h \in [H]} \widehat{d}_{\pi^i, h}^\dagger(s, a) - H \ln \left(\frac{6H|\mathcal{S}||\mathcal{A}|}{\delta} \right).$$

864 If $H \ln(6H|\mathcal{S}||\mathcal{A}|/\delta) \leq \frac{1}{4} \sum_{i < k} \sum_{h \in [H]} \widehat{d}_{\pi^i, h}^\dagger(s, a)$, then $n^k(s, a) \geq \frac{1}{4} \sum_{i < k} \sum_{h \in [H]} \widehat{d}_{\pi^i, h}^\dagger(s, a)$
865 and hence,

$$\min \{1, \phi(n^k(s, a))\} \leq \phi(n^k(s, a)) \leq \phi \left(\frac{1}{4} \sum_{i < k} \sum_{h \in [H]} \widehat{d}_{\pi^i, h}^\dagger(s, a) \right) \leq 4\phi \left(\sum_{i < k} \sum_{h \in [H]} \widehat{d}_{\pi^i, h}^\dagger(s, a) \right).$$

866 This result equals to the right hand side in the lemma, because $\sum_{i < k} \sum_{h \in [H]} \widehat{d}_{\pi^i, h}^\dagger(s, a) \geq$
867 $4H \ln(6H|\mathcal{S}||\mathcal{A}|/\delta) \geq 1$ (so taking maximum does not change anything). Otherwise, if
868 $H \ln(6H|\mathcal{S}||\mathcal{A}|/\delta) > \frac{1}{4} \sum_{i < k} \sum_{h \in [H]} \widehat{d}_{\pi^i, h}^\dagger(s, a)$, then

$$\min \{1, \phi(n^k(s, a))\} \leq 1 \leq \frac{4H \ln(H|\mathcal{S}||\mathcal{A}|/\delta')}{\sum_{i < k} \sum_{h \in [H]} \widehat{d}_{\pi^i, h}^\dagger(s, a)}.$$

869 Since $1 \leq 4H \ln(6H|\mathcal{S}||\mathcal{A}|/\delta)$, we have

$$\min \{1, \phi(n^k(s, a))\} \leq 1 \leq \frac{4H \ln(H|\mathcal{S}||\mathcal{A}|/\delta')}{\max \{1, \sum_{i < k} \sum_{h \in [H]} \widehat{d}_{\pi^i, h}^\dagger(s, a)\}} \leq \text{RHS}$$

870 The last inequality comes from simple algebra. Therefore, we conclude. \square

871 A.3.5 Proof for lemma A.21 and lemma A.22 (properties of empirical sparsified MDP)

872 In this section, we state two important properties of the empirical sparsified MDP and prove them.

873 We remark that we do not include $\mathbb{P}(s^\dagger | s, a)$ in these two lemmas, since by definition $s^\dagger \notin \mathcal{S}$.

Lemma A.21 (Multiplicative Accuracy). *We set*

$$\Phi \geq 6H^2 \log \left(\frac{12|\mathcal{S}|^2|\mathcal{A}|}{\delta} \right).$$

874 *Then, when $H \geq 2$, under \mathcal{E}^P , we have for any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$,*

$$\left(1 - \frac{1}{H}\right) \widehat{\mathbb{P}}^\dagger(s' | s, a) \leq \mathbb{P}^\dagger(s' | s, a) \leq \left(1 + \frac{1}{H}\right) \widehat{\mathbb{P}}^\dagger(s' | s, a).$$

875 *Proof.* For $N(s, a, s') < \Phi$, both sides are zero. For $N(s, a, s') \geq \Phi$, recall $\widehat{\mathbb{P}}^\dagger(s' | s, a) =$

876 $\frac{N(s, a, s')}{N(s, a)}$, then from lemma A.10, with probability at least $1 - \delta$,

$$\begin{aligned} & \left| \widehat{\mathbb{P}}^\dagger(s' | s, a) - \mathbb{P}^\dagger(s' | s, a) \right| \\ & \leq \sqrt{\frac{2\widehat{\mathbb{P}}^\dagger(s' | s, a)}{N(s, a)} \log \left(\frac{12|\mathcal{S}|^2|\mathcal{A}|}{\delta} \right)} + \frac{14}{3N(s, a)} \log \left(\frac{12|\mathcal{S}|^2|\mathcal{A}|}{\delta} \right) \\ & \hspace{15em} (\text{lemma A.10 and definition of } \mathcal{E}^P) \\ & = \left[\sqrt{\frac{2}{N(s, a, s')} \log \left(\frac{12|\mathcal{S}|^2|\mathcal{A}|}{\delta} \right)} + \frac{14}{3N(s, a, s')} \log \left(\frac{12|\mathcal{S}|^2|\mathcal{A}|}{\delta} \right) \right] \cdot \widehat{\mathbb{P}}^\dagger(s' | s, a) \\ & \hspace{15em} (\widehat{\mathbb{P}}^\dagger(s' | s, a) = \frac{N(s, a, s')}{N(s, a)}) \\ & \leq \left[\sqrt{\frac{1}{3H^2}} + \frac{7}{9H^2} \right] \cdot \widehat{\mathbb{P}}^\dagger(s' | s, a) \leq \frac{\widehat{\mathbb{P}}^\dagger(s' | s, a)}{H}, \end{aligned}$$

877 where the second line comes from the lower bound on Φ . We conclude. \square

878

879 **Lemma A.22** (Bound on Ratios of Visitation Probability). *For any deterministic policy π and any*
 880 *$(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$, it holds that*

$$\frac{1}{4}d_{\pi,h}^\dagger(s, a) \leq \widehat{d}_{\pi,h}^\dagger(s, a) \leq 3d_{\pi,h}^\dagger(s, a).$$

881 *Here, recall that we denote $d_{\pi,h}^\dagger(s, a)$ and $\widehat{d}_{\pi,h}^\dagger(s, a)$ as the occupancy measure of (s, a) at stage h*
 882 *under policy π , on \mathbb{P}^\dagger (the transition dynamics in the sparsified MDP) and $\widehat{\mathbb{P}}^\dagger$ (the transition dynamics*
 883 *in the empirical sparsified MDP) respectively.*

884 We remark that for $s^\dagger \notin \mathcal{S}$ the inequality does not necessarily hold.

885 *Proof.* We denote $T_{h,s,a}$ as all truncated trajectories $(s_1, a_1, s_2, a_2, \dots, s_h, a_h)$ up to stage h such that
 886 $(s_h, a_h) = (s, a)$. Notice that if $\tau_h = (s_1, a_1, s_2, a_2, \dots, s_h, a_h) \in T_{h,s,a}$, then it holds that $s_i \neq s^\dagger$
 887 for $1 \leq i \leq h-1$. We denote $\mathbb{P}(\cdot; \mathbb{P}', \pi)$ as the probability under a specific transition dynamics \mathbb{P}'
 888 and policy π . For any fixed $(h, s, a) \in [H] \times \mathcal{S} \times \mathcal{A}$ and any fixed $\tau \in T_{h,s,a}$, we apply Lemma A.21
 889 to get

$$\begin{aligned} \mathbb{P}[\tau; \widehat{\mathbb{P}}^\dagger, \pi] &= \prod_{i=1}^h \pi_i(a_i | s_i) \prod_{i=1}^{h-1} \widehat{\mathbb{P}}^\dagger(s_{i+1} | s_i, a_i) \\ &\leq \left(1 + \frac{1}{H}\right)^H \prod_{i=1}^h \pi_i(a_i | s_i) \prod_{i=1}^{h-1} \mathbb{P}^\dagger(s_{i+1} | s_i, a_i) \leq 3\mathbb{P}[\tau; \mathbb{P}^\dagger, \pi] \end{aligned}$$

890 and

$$\begin{aligned} \mathbb{P}[\tau; \widehat{\mathbb{P}}^\dagger, \pi] &= \prod_{i=1}^h \pi_i(a_i | s_i) \prod_{i=1}^{h-1} \widehat{\mathbb{P}}^\dagger(s_{i+1} | s_i, a_i) \\ &\geq \left(1 - \frac{1}{H}\right)^H \prod_{i=1}^h \pi_i(a_i | s_i) \prod_{i=1}^{h-1} \mathbb{P}^\dagger(s_{i+1} | s_i, a_i) \geq \frac{1}{4}\mathbb{P}[\tau; \mathbb{P}^\dagger, \pi]. \end{aligned}$$

891 We conclude by rewriting the visiting probability as

$$d_{\pi,h}^\dagger(s, a) = \sum_{\tau \in T_{h,s,a}} \mathbb{P}[\tau; \mathbb{P}^\dagger, \pi]; \quad \widehat{d}_{\pi,h}^\dagger(s, a) = \sum_{\tau \in T_{h,s,a}} \mathbb{P}[\tau; \widehat{\mathbb{P}}^\dagger, \pi].$$

892

□

893 B Additional comparisons

894 B.1 Comparison with other comparator policy

895 Our main result compares the sub-optimality of the policy π_{final} against the optimal policy on the
 896 sparsified MDP. We can further derive the sub-optimality of our output with respect to any comparator
 897 policy on the original MDP \mathcal{M} . If we denote π_* , π_*^\dagger and π_{final} as the global optimal policy, the
 898 optimal policy on the sparsified MDP and the policy output by our algorithm, respectively, and denote
 899 π as the comparator policy, we have

$$\begin{aligned}
 & V_1(s_1, \mathbb{P}, r, \pi) - V_1(s_1, \mathbb{P}, r, \pi_{final}) \\
 & \leq V_1(s_1, \mathbb{P}, r, \pi) - V_1(s_1, \mathbb{P}^\dagger, r, \pi) + \underbrace{V_1(s_1, \mathbb{P}^\dagger, r, \pi) - V_1(s_1, \mathbb{P}^\dagger, r, \pi_*^\dagger)}_{\leq 0} \\
 & \quad + \underbrace{V_1(s_1, \mathbb{P}^\dagger, r, \pi_*^\dagger) - V_1(s_1, \mathbb{P}^\dagger, r, \pi_{final})}_{\lesssim \varepsilon} + \underbrace{V_1(s_1, \mathbb{P}^\dagger, r, \pi_{final}) - V_1(s_1, \mathbb{P}, r, \pi_{final})}_{\leq 0} \\
 & \lesssim \underbrace{V_1(s_1, \mathbb{P}, r, \pi_*) - V_1(s_1, \mathbb{P}^\dagger, r, \pi_*)}_{\text{Approximation Error}} + \varepsilon. \tag{B.1}
 \end{aligned}$$

900 Here, the second term is non-positive from the definition of π_*^\dagger , the third term is upper bounded by ε
 901 due to our main theorem (Theorem 5.1), and the last term is non-positive from the definition of the
 902 sparsified MDP. Since the connectivity graph of the sparsified MDP is a sub-graph of the original
 903 MDP, for any policy, the policy value on the sparsified MDP must be no higher than that on the
 904 original MDP.

905 At a high level, the ε term in the last line of (B.1) represents the error from the finite online episodes,
 906 while the approximation error term $V_1(s_1, \mathbb{P}, r, \pi) - V_1(s_1, \mathbb{P}^\dagger, r, \pi)$ measures the policy value
 907 difference of π on the original MDP and the sparsified one, representing the *coverage quality of the*
 908 *offline dataset*. If the dataset covers most of what π covers, then this gap should be small. When π is
 909 the global optimal policy π_* , this means the data should cover the state-actions pairs where optimal
 910 policy covers. The approximation error here plays a similar role as the concentrability coefficient in
 911 the offline reinforcement learning.

912 B.2 Comparison with offline reinforcement learning

913 Our algorithm leverages some information from the offline dataset, so it is natural to ask how we
 914 expect that offline dataset to be, compared to traditional offline reinforcement learning does. In
 915 offline RL, we typically require the *concentrability condition*, namely good coverage for the offline
 916 dataset, in order to achieve a polynomial sample complexity. Specifically, if we assume the offline
 917 data are sampled by first sampling (s, a) i.i.d. from μ and then sampling the subsequent state from the
 918 transition dynamics, then the concentrability condition says the following constant C^* is well-defined
 919 and finite.

$$C^* := \sup_{(s,a)} \frac{d^{\pi_*}(s, a)}{\mu(s, a)} < \infty.$$

920 The concentrability coefficient can be defined in several alternative ways, either for a set of policies
 921 or with respect to a single policy [Chen and Jiang, 2019, Zhan et al., 2022, Xie et al., 2021b, Zanette
 922 et al., 2021b]. Here, we follow the definition in [Xie et al., 2021b]. This means, the sampling
 923 distribution must covers the region where the global optimal policy covers, which is a very similar
 924 intuition obtained from our setting.

925 [Xie et al., 2021b] also gave optimal sample complexity (in terms of state-action pairs) for an offline
 926 RL algorithm is

$$N = \tilde{O} \left(\frac{C^* H^3 |\mathcal{S}|}{\varepsilon^2} + \frac{C^* H^{5.5} |\mathcal{S}|}{\varepsilon} \right),$$

927 which is minimax optimal up to logarithm terms and higher order terms. Similar sample complexity
 928 were also given in several literature [Yin and Wang, 2020, Yin et al., 2020, Xie and Jiang, 2020b].

929 **Uniform data distribution** For simplicity, we first assume μ to be uniform on all state-action
 930 pairs and the reward function to be given. Consider we have N state-action pairs in the offline data,
 931 which are sampled i.i.d. from the distribution μ . Notice that here, the global optimal policy π_* still
 932 differs from the optimum on the sparsified MDP π_*^\dagger , since even if we get enough samples from each
 933 (s, a) pairs, we might not get enough samples for every (s, a, s') and hence, not all (s, a, s') will be
 934 included in the set of known tuples.

935 Concretely, if we consider the case when we sample each state-action pair for $N/(|\mathcal{S}||\mathcal{A}|)$ times
 936 and simply treat the transition frequency as the true probability, then for any $N(s, a, s') < \Phi$, it
 937 holds that $\mathbb{P}(s' | s, a) = \frac{N(s, a, s')}{N(s, a)} = \frac{N(s, a, s')|\mathcal{S}||\mathcal{A}|}{N} \leq \frac{\Phi|\mathcal{S}||\mathcal{A}|}{N}$. So for any any $N(s, a, s') \geq \Phi$, we
 938 know $\mathbb{P}(s' | s, a) = \mathbb{P}^\dagger(s' | s, a)$; while for any $N(s, a, s') < \Phi$, we have $\mathbb{P}(s' | s, a) \leq \frac{\Phi|\mathcal{S}||\mathcal{A}|}{N}$ and
 939 $\mathbb{P}^\dagger(s' | s, a) = 0$. Therefore, we have

$$|\mathbb{P}(s' | s, a) - \mathbb{P}^\dagger(s' | s, a)| \leq \frac{\Phi|\mathcal{S}||\mathcal{A}|}{N}$$

940 From the value difference lemma (lemma D.11), we can upper bound the approximation error by

$$\begin{aligned} & V_1(s_1, \mathbb{P}, r, \pi) - V_1(s_1, \mathbb{P}^\dagger, r, \pi) \\ &= \mathbb{E}_{\mathbb{P}, \pi} \left[\sum_{h=1}^H \sum_{s_{h+1}} (\mathbb{P}^\dagger(s_{h+1} | s_h, a_h) - \mathbb{P}(s_{h+1} | s_h, a_h)) \cdot V_h(s_{h+1}, \mathbb{P}, r, \pi) \Big| s_h = s \right] \\ & \hspace{20em} \text{(lemma D.11)} \\ &\leq \mathbb{E}_{\mathbb{P}, \pi} \left[\sum_{h=1}^H \sum_{s_{h+1}} \frac{\Phi|\mathcal{S}||\mathcal{A}|}{N} \cdot H \Big| s_h = s \right] \hspace{2em} \text{(the value function is upper bounded by } H) \\ &\leq \frac{\Phi H^2 |\mathcal{S}|^2 |\mathcal{A}|}{N} \hspace{15em} \text{(summation over } h \in [H] \text{ and } s_{h+1} \in \mathcal{S}) \\ &\asymp \tilde{O} \left(\frac{H^4 |\mathcal{S}|^2 |\mathcal{A}|}{N} \right). \hspace{10em} \text{(definition of } \Phi \text{ in (5.1))} \end{aligned}$$

941 Therefore, to get an ε -optimal policy compared to the global optimal one, we need the number of
 942 state-action pairs in the initial offline dataset \mathcal{D} to be

$$N = \tilde{O} \left(\frac{H^4 |\mathcal{S}|^2 |\mathcal{A}|}{\varepsilon} \right).$$

943 From the theorem 5.1, the offline data size we need here is actually significantly smaller than what
 944 we need for an offline algorithm. As long as $\sup_{(s, a)} d^{\pi_*}(s, a)$ is not too small, for instance, larger
 945 than $H^{-1.5}$, then we shave off the whole $1/\varepsilon^2$ term. The order of offline sample complexity here is
 946 actually $O(1/\varepsilon)$ instead of $O(1/\varepsilon^2)$ typical in offline RL, and this is significantly smaller in small
 947 ε regime. To compensate the smaller offline sample size, actually we need more online sample to
 948 obtain an globally ε -optimal policy, and we summarize the general requirement for offline and online
 949 sample size in corollary 5.2.

950 **Non-uniform data distribution** Assume the data generating distribution μ is not uniform but still
 951 supported on all (s, a) pairs such that $d^{\pi_*}(s, a) > 0$, so that the concentrability coefficient in offline
 952 RL is still well defined. We simply consider the case when each state-action pair (s, a) is sampled by
 953 $N\mu(s, a)$ times and treat the transition frequency as the true underlying probability. Then, following
 954 a very similar argument as in the last paragraph, the number of state-action pairs needed in the initial
 955 offline dataset in order to extract an ε -globally optimal policy is

$$N = \tilde{O} \left(\frac{H^4 |\mathcal{S}|}{\varepsilon} \sum_{s, a} \frac{d^{\pi_*}(s, a)}{\mu(s, a)} \right).$$

Here, the quantity

$$C^\dagger := \sum_{s, a} \frac{d^{\pi_*}(s, a)}{\mu(s, a)}$$

956 plays a similar role of classical concentrability coefficient and also measures the distribution shift
 957 between two policies. In the worst case, this coefficient can be $|\mathcal{S}| |\mathcal{A}| C^*$, resulting in an extra
 958 $|\mathcal{S}| |\mathcal{A}|$ factor compared to the optimal offline sample complexity. However, we still shave off the
 959 entire $1/\varepsilon^2$ term and also shave off $H^{1.5}$ in the $1/\varepsilon$ term.

960 **Partial coverage data** Under partial coverage, we expect the output policy π_{final} to be competitive
 961 with the value of the best policy supported in the region covered by the offline dataset. In such case,
 962 theorem 5.1 provides guarantees with the best comparator policy on the sparsified MDP \mathcal{M}^\dagger . In order
 963 to gain further intuition, it is best to ‘translate’ such guarantees into guarantees on \mathcal{M} .

964 In the worst case, the data distribution μ at a certain (s, a) pair can be zero when $d^\pi(s, a) > 0$,
 965 which implies the concentrability coefficient $C^* = \infty$. Here, π is an arbitrary comparator policy. In
 966 this case, either classical offline RL algorithm or our policy finetuning algorithm cannot guarantee
 967 an ε -optimal policy compared to the global optimal policy. However, we can still output a locally
 968 ε -optimal policy, compared to the optimal policy on the sparsified MDP.

969 In order to compare π_{final} to any policy on the original MDP, we have the corollary 5.2, which will
 970 be proved in appendix B.3.

971 The statement in corollary 5.2 is a quite direct consequence of theorem 5.1, and it expresses the
 972 sub-optimality gap of π_{final} with respect to any comparator policy π on the original MDP \mathcal{M} . It can
 973 also be written in terms of the sub-optimality: If we fix a comparator policy π , then with probability
 974 at least $1 - \delta$, for any reward function r , the policy π_{final} returned by algorithm 2 satisfies:

$$\begin{aligned} V_1(s_1; \mathbb{P}, r, \pi) - V_1(s_1; \mathbb{P}, r, \pi_{final}) &= \underbrace{\tilde{O}\left(\frac{H |\mathcal{S}| \sqrt{|\mathcal{A}|}}{\sqrt{K_{de}}}\right)}_{\text{Online error}} + \underbrace{\frac{H^4 |\mathcal{S}|}{N} \sum_{s,a} \frac{d^\pi(s, a)}{\mu(s, a)}}_{\text{Offline error}} \\ &= \tilde{O}\left(\frac{H |\mathcal{S}| \sqrt{|\mathcal{A}|}}{\sqrt{K_{de}}} + \frac{H^4 |\mathcal{S}|^2 |\mathcal{A}|}{N} \sup_{s,a} \frac{d^\pi(s, a)}{\mu(s, a)}\right), \end{aligned}$$

975 where K_{de} is the number of online episodes and N is the number of state-action pairs in offline data.
 976 Here, the sub-optimality depends on an *online error* as well as on an *offline error*. The online error is
 977 the one that also arises in the statement of theorem 5.1. It is an error that can be reduced by collecting
 978 more online samples, i.e., by increasing K , with the typical inverse square-root dependence $1/\sqrt{K}$.

979 However, the upper bound suggests that even in the limit of infinite online data, the value of π_{final}
 980 will not approach that of π_* because of a residual error due to the *offline* dataset \mathcal{D} . Such residual
 981 error depends on certain concentrability factor expressed as $\sum_{s,a} \frac{d^\pi(s, a)}{\mu(s, a)} \leq |\mathcal{S}| |\mathcal{A}| \sup_{s,a} \frac{d^\pi(s, a)}{\mu(s, a)}$,
 982 whose presence is intuitive: if a comparator policy π is not covered well, our algorithm does not
 983 have enough information to navigate to the area that π tends to visit, and so it is unable to refine
 984 its estimates there. However the dependence on the number of offline samples $N = |\mathcal{D}|$ is through
 985 its *inverse*, i.e., $1/N$ as opposed to the more typical $1/\sqrt{N}$: such gap represents the improvable
 986 performance when additional online data are collected non-reactively.

987 It is useful to compare corollary 5.2 with what is achievable by using a minimax-optimal online
 988 algorithm[Xie et al., 2021b]. In this latter case, one can bound the sub-optimality gap for any
 989 comparator policy π with high probability as

$$V_1(s_1; \mathbb{P}, r^\dagger, \pi) - V_1(s_1; \mathbb{P}, r^\dagger, \pi_{final}) \leq \tilde{O}\left(\sqrt{\frac{H^3 |\mathcal{S}|}{N} \sup_{s,a} \frac{d^\pi(s, a)}{\mu(s, a)}}\right). \quad (\text{B.2})$$

990 B.3 Proof of corollary 5.2

991 Let’s denote $\mathcal{D} = \{(s_i, a_i, s'_i)\}_{i \in [N]}$ as the offline dataset, where N as the total number of tuples. We
 992 keep the notation the same as in the main text. We use $N(s, a)$ and $N(s, a, s')$ to denote the counter
 993 of (s, a) and (s, a, s') in the offline data \mathcal{D} . The state-action pairs are sampled i.i.d. from $\mu(s, a)$
 994 and the subsequent states are sampled from the transition dynamics. We fix a comparator policy π
 995 and assume $\mu(s, a) > 0$ for any (s, a) such that $d^\pi(s, a) > 0$, which implies a finite concentrability
 996 constant C^* . Here, $d^\pi(s, a)$ is the occupancy probability of (s, a) when executing policy π , averaged
 997 over all stages $h \in [H]$.

998 Similar to the Section B.1, we have

$$\begin{aligned}
& V_1(s_1, \mathbb{P}, r, \pi) - V_1(s_1, \mathbb{P}, r, \pi_{final}) \\
& \leq V_1(s_1, \mathbb{P}, r, \pi) - V_1(s_1, \mathbb{P}^\dagger, r, \pi) + \underbrace{V_1(s_1, \mathbb{P}^\dagger, r, \pi) - V_1(s_1, \mathbb{P}^\dagger, r, \pi_*^\dagger)}_{\leq 0} \\
& + V_1(s_1, \mathbb{P}^\dagger, r, \pi_*^\dagger) - V_1(s_1, \mathbb{P}^\dagger, r, \pi_{final}) + \underbrace{V_1(s_1, \mathbb{P}^\dagger, r, \pi_{final}) - V_1(s_1, \mathbb{P}, r, \pi_{final})}_{\leq 0} \\
& \lesssim \underbrace{V_1(s_1, \mathbb{P}, r, \pi) - V_1(s_1, \mathbb{P}^\dagger, r, \pi)}_{\text{Approximation Error}} + \underbrace{V_1(s_1, \mathbb{P}^\dagger, r, \pi_*^\dagger) - V_1(s_1, \mathbb{P}^\dagger, r, \pi_{final})}_{\text{Estimation Error}}. \tag{B.3}
\end{aligned}$$

999 where π_*^\dagger and π_{final} are the optimal policy on the sparsified MDP and the policy output by our
1000 algorithm, respectively, and π is the fixed comparator policy. Here, the second term is non-positive
1001 from the definition of π_*^\dagger , the last term is non-positive from the definition of the sparsified MDP.
1002 This is because, for any state-action pair (s, a) and any fixed policy π , the probability of reaching
1003 (s, a) under \mathbb{P}^\dagger will not exceed that under the true transition probability \mathbb{P} . If we denote the visiting
1004 probability under π and \mathbb{P} (or \mathbb{P}^\dagger resp.) as $d_{\pi, h}(s, a)$ ($d_{\pi, h}^\dagger(s, a)$ resp.), then we have

$$d_{\pi, h}^\dagger(s, a) \leq d_{\pi, h}(s, a)$$

1005 for any $h \in [H]$, $s \in \mathcal{S}$, $a \in \mathcal{A}$. Note that, for $s = s^\dagger$, this does not hold necessarily. Then, for any
1006 policy π , we have

$$\begin{aligned}
V_1(s_1, \mathbb{P}^\dagger, r, \pi) &= \sum_{h=1}^H \sum_{s, a} d_{\pi, h}^\dagger(s, a) r(s, a) && \text{(definition of policy value)} \\
&= \sum_{h=1}^H \sum_{s \neq s^\dagger, a} d_{\pi, h}^\dagger(s, a) r(s, a) && (r(s^\dagger, a) = 0 \text{ for any } a) \\
&\leq \sum_{h=1}^H \sum_{s \neq s^\dagger, a} d_{\pi, h}(s, a) r(s, a) = V_1(s_1, \mathbb{P}, r, \pi).
\end{aligned}$$

1007 Therefore, we get $V_1(s_1, \mathbb{P}^\dagger, r, \pi_{final}) - V_1(s_1, \mathbb{P}, r, \pi_{final}) \leq 0$ when we take $\pi = \pi_{final}$.

1008 From the main result (Theorem 5.1), we know the estimation error is bounded as

$$\underbrace{V_1(s_1, \mathbb{P}^\dagger, r, \pi_*^\dagger) - V_1(s_1, \mathbb{P}^\dagger, r, \pi_{final})}_{\text{Estimation Error}} \lesssim \tilde{O}\left(\frac{H|\mathcal{S}|\sqrt{|\mathcal{A}|}}{\sqrt{K_{de}}}\right), \tag{B.4}$$

1009 where K_{de} is the number of online episodes. Therefore, to make the right hand side of (B.4) less than
1010 $\varepsilon/2$, one needs at least $\tilde{O}\left(\frac{H^2|\mathcal{S}||\mathcal{A}|}{\varepsilon^2}\right)$ online episodes. This is exactly what the main result shows.

1011

1012 So it suffices to bound the approximation error term. From the value difference lemma (Lemma
1013 D.11), we have

$$\begin{aligned}
& |V_1(s_1, \mathbb{P}, r, \pi) - V_1(s_1, \mathbb{P}^\dagger, r, \pi)| \\
&= \left| \mathbb{E}_{\mathbb{P}, \pi} \left[\sum_{i=1}^H (\mathbb{P}^\dagger(\cdot | s_i, a_i) - \mathbb{P}(\cdot | s_i, a_i))^\top V_{i+1}(\cdot; \mathbb{P}^\dagger, r, \pi) \Big| s_1 = s \right] \right| \\
&= \left| \sum_{i=1}^H \sum_{s_i, a_i} d_{\pi, h}(s_i, a_i) (\mathbb{P}^\dagger(\cdot | s_i, a_i) - \mathbb{P}(\cdot | s_i, a_i))^\top V_{i+1}(\cdot; \mathbb{P}^\dagger, r, \pi) \right| \\
& \hspace{15em} \text{(by expanding the expectation)} \\
&\leq H |\mathcal{S}| \cdot \sum_{i=1}^H \sum_{s_i, a_i} d_{\pi, h}(s_i, a_i) \left(\sup_{s'} |\mathbb{P}^\dagger(s' | s_i, a_i) - \mathbb{P}(s' | s_i, a_i)| \right). \\
& \hspace{15em} (V_{i+1} \leq H \text{ and the inner product has } |\mathcal{S}| \text{ terms})
\end{aligned}$$

1014 We define

$$d_\pi(s, a) = \frac{1}{H} \sum_{h=1}^H d_{\pi, h}(s, a) \quad (\text{B.5})$$

1015 as the average visiting probability. Then, we have

$$|V_1(s_1, \mathbb{P}, r, \pi) - V_1(s_1, \mathbb{P}^\dagger, r, \pi)| \leq |\mathcal{S}| H^2 \sum_{s, a} \left[\sup_{s'} |\mathbb{P}^\dagger(s' | s, a) - \mathbb{P}(s' | s, a)| d_\pi(s, a) \right], \quad (\text{B.6})$$

1016 So it suffices to upper bound $|\mathbb{P}^\dagger(s' | s, a) - \mathbb{P}(s' | s, a)|$. Notice here we only consider $s \neq s^\dagger$ and
 1017 $s' \neq s^\dagger$, since the value function starting from s^\dagger is always zero.

1018 For (s, a, s') , if $N(s, a, s') \geq \Phi$, it holds that $\mathbb{P}^\dagger(s' | s, a) = \mathbb{P}(s' | s, a)$. Otherwise, from
 1019 the definition, we know $\mathbb{P}^\dagger(s' | s, a) = 0$, so it suffices to bound $\mathbb{P}(s' | s, a)$ in this case. From
 1020 lemma B.1, we know that with probability at least $1 - \delta/2$, for any (s, a, s') such that $N(s, a, s') < \Phi$,
 1021 it holds that

$$\mathbb{P}(s' | s, a) \leq \frac{2N(s, a, s') + 2 \log(2|\mathcal{S}|^2 |\mathcal{A}| / \delta)}{N(s, a)}.$$

1022 Then we deal with two cases. When $N\mu(s, a) \geq 6\Phi$, from lemma B.2 we have

$$\mathbb{P}(s' | s, a) \leq \frac{4N(s, a, s') + 2 \log(4|\mathcal{S}|^2 |\mathcal{A}| / \delta)}{N\mu(s, a) - 2 \log(2|\mathcal{S}| |\mathcal{A}| / \delta)} \leq \frac{4\Phi + 2 \log(4|\mathcal{S}|^2 |\mathcal{A}| / \delta)}{N\mu(s, a) - 2 \log(2|\mathcal{S}| |\mathcal{A}| / \delta)}.$$

1023 From the definition of Φ , we know that $2 \log(4|\mathcal{S}|^2 |\mathcal{A}| / \delta) \leq \Phi$ and $2 \log(2|\mathcal{S}| |\mathcal{A}| / \delta) \leq \Phi$,
 1024 which implies

$$\mathbb{P}(s' | s, a) \leq \frac{5\Phi}{N\mu(s, a) - \Phi} \leq \frac{6\Phi}{N\mu(s, a)}.$$

1025 The last inequality comes from our assumption for $N\mu(s, a) \geq 6\Phi$.

1026 In the other case, when $N\mu(s, a) < 6\Phi$, it holds that

$$\mathbb{P}(s' | s, a) \leq 1 \leq \frac{6\Phi}{N\mu(s, a)}$$

1027 for any (s, a, s') . Therefore, for any for any (s, a, s') such that $N(s, a, s') < \Phi$, we have

$$\mathbb{P}(s' | s, a) \leq \frac{6\Phi}{N\mu(s, a)}. \quad (\text{B.7})$$

1028 Combining equations (B.7) and (B.6), we know for any comparator policy π , it holds that

$$\underbrace{|V_1(s_1, \mathbb{P}, r, \pi) - V_1(s_1, \mathbb{P}^\dagger, r, \pi)|}_{\text{Approximation Error}} \lesssim \frac{\Phi H^2 |\mathcal{S}|}{N} \sum_{s, a} \frac{d_\pi(s, a)}{\mu(s, a)} \lesssim \tilde{O} \left(\frac{H^4 |\mathcal{S}|}{N} \sum_{s, a} \frac{d_\pi(s, a)}{\mu(s, a)} \right),$$

1029 where N is the total number of transitions in the offline data. In order to make the right hand side of
 1030 last display less than $\varepsilon/2$, one needs at least

$$\tilde{O} \left(\frac{H^4 |\mathcal{S}|}{\varepsilon} \sum_{s, a} \frac{d_\pi(s, a)}{\mu(s, a)} \right) \leq \tilde{O} \left(\frac{H^4 |\mathcal{S}|^2 |\mathcal{A}| C^*}{\varepsilon} \right),$$

1031 offline transitions. Combining the proof for estimation error and approximation error, we conclude.

1032 **B.4 Proof for lemma B.1 and lemma B.2**

1033 **Lemma B.1.** *With probability at least $1 - \delta/2$, for any $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, it holds that*

$$N(s, a, s') \geq \frac{1}{2}N(s, a)\mathbb{P}(s' | s, a) - \log\left(\frac{2|\mathcal{S}|^2|\mathcal{A}|}{\delta}\right).$$

1034 *Proof.* We fixed (s, a, s') and denote I as the index set where $(s_i, a_i) = (s, a)$ for $i \in I$. We range the
 1035 indexes in I as $i_1 < i_2 < \dots < i_{N(s, a)}$. For $j \leq N(s, a)$, we denote $X_j = \mathbb{I}(s'_{i_j} = s')$, which is the
 1036 indicator of whether the next state is s' when we encounter (s, a) the j -th time. When $j \geq N(s, a)$,
 1037 we denote X_j as independent Bernoulli random variables with successful rate $\mathbb{P}(s' | s, a)$. Then, we
 1038 know X_j for all $j \in \mathbb{N}$ are i.i.d. sequence of Bernoulli random variables. From Lemma D.1, we know
 1039 with probability at least $1 - \delta/2$, for any positive integer n , it holds that

$$\sum_{j=1}^n X_j \geq \frac{1}{2} \sum_{j=1}^n \mathbb{P}(s' | s, a) - \log\left(\frac{2}{\delta}\right).$$

1040 We take $n = N(s, a)$ (although $N(s, a)$ is random, we can still take it because for any n the inequality
 1041 above holds) to get

$$N(s, a, s') = \sum_{j=1}^{N(s, a)} X_j \geq \frac{1}{2}N(s, a)\mathbb{P}(s' | s, a) - \log\left(\frac{2}{\delta}\right).$$

1042 Applying a union bound for all (s, a, s') , we conclude. □

1043

1044 **Lemma B.2.** *With probability at least $1 - \delta/2$, for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, it holds that*

$$N(s, a) \geq \frac{1}{2}N\mu(s, a) - \log\left(\frac{2|\mathcal{S}||\mathcal{A}|}{\delta}\right).$$

1045 *Proof.* If $\mu(s, a) = 0$, this is trivial. We fixed an (s, a) such that $\mu(s, a) > 0$. For $j \leq N$, we denote
 1046 $X_j = \mathbb{I}(s_j = s, a_j = a)$, which is the indicator of whether the j -th state-action pair we encounter in
 1047 the offline dataset is (s, a) . When $j \geq N$, we denote X_j as independent Bernoulli random variables
 1048 with successful rate $\mu(s, a)$. Then, we know X_j for all $j \in \mathbb{N}$ are i.i.d. sequence of Bernoulli random
 1049 variables. From Lemma D.1, we know with probability at least $1 - \delta/2$, for any positive integer n , it
 1050 holds that

$$\sum_{j=1}^n X_j \geq \frac{1}{2} \sum_{j=1}^n \mu(s, a) - \log\left(\frac{2}{\delta}\right).$$

1051 We take $n = N$ to get

$$N(s, a) = \sum_{j=1}^N X_j \geq \frac{1}{2}N\mu(s, a) - \log\left(\frac{2}{\delta}\right).$$

1052 Applying a union bound for all (s, a) such that $\mu(s, a) > 0$, we conclude. □

1053 **C Lower bound**

1054 In this section we briefly discuss the optimality of the algorithm. Although the following consid-
 1055 erations are also mentioned in the main text, here we mention how they naturally lead to a lower
 1056 bound.

1057 **Lower bound for reward-free exploration** Consider the MDP class \mathcal{M} defined in the proof of
 1058 the lower bound of Theorem 4.1 in [Jin et al., 2020b]. Assume that the dataset arises from a logging
 1059 policy π_{log} which induces the condition $N(s, a, s') \geq \Phi$ for all $(s, a, s') \in \mathcal{S} \times \mathcal{A}$ for every instance
 1060 of the class. In this case, every MDP instance $\mathcal{M} \in \mathcal{M}$ and its sparsified version \mathcal{M}^\dagger coincide. Then
 1061 the concatenation of the logging policy π_{log} and of the policy π_{final} produced by our algorithm (i.e.,
 1062 algorithm 3) can be interpreted as a reactive policy, which must comply with the reward free lower
 1063 bound established in Theorem 4.1 of [Jin et al., 2020b]. More precisely, the reward-free sample
 1064 complexity lower bound established in Theorem 4.1 in [Jin et al., 2020b] is

$$\Omega\left(\frac{|\mathcal{S}|^2|\mathcal{A}|H^2}{\varepsilon^2}\right) \tag{C.1}$$

1065 trajectories. This matches the sample complexity of theorem 5.1. Notice that the number of samples
 1066 originally present in the dataset can be

$$|\mathcal{S}|^2|\mathcal{A}| \times \tilde{O}(H^2) = \tilde{O}(H^2|\mathcal{S}|^2|\mathcal{A}|), \tag{C.2}$$

1067 a term independent of the accuracy ε . Given that when $\varepsilon \leq 1$ we have

$$\tilde{O}(H^2|\mathcal{S}|^2|\mathcal{A}|) + \Omega\left(\frac{|\mathcal{S}|^2|\mathcal{A}|H^2}{\varepsilon^2}\right) = \Omega\left(\frac{|\mathcal{S}|^2|\mathcal{A}|H^2}{\varepsilon^2}\right),$$

1068 our algorithm is unimprovable beyond constant terms and logarithmic terms in a minimax sense.

1069 **Lower bound for non-reactive exploration** Consider the MDP class \mathcal{M} defined in the proof of
 1070 the lower bound in Theorem 1 of [Xiao et al., 2022]. It establishes an exponential sample complexity
 1071 for non-reactive exploration *when no prior knowledge is available*. In other words, in absence of any
 1072 data about the MDP, non-reactive exploration must suffer an exponential sample complexity. In such
 1073 case, our theorem 5.1 (correctly) provides vacuous guarantees, because the sparsified MDP \mathcal{M} is
 1074 degenerate (all edges lead to the absorbing state).

1075 **Combining the two constructions** It is possible to combine the MDP class \mathcal{M}_1 from the paper
 1076 [Xiao et al., 2022] with the MDP class \mathcal{M}_2 from the paper [Jin et al., 2020b]. In lieu of a formal
 1077 proof, here we provide only a sketch of the construction that would induce a lower bound. More
 1078 precisely, consider a starting state s_1 where only two actions— $a = 1$ and $a = 2$ —are available.
 1079 Taking $a = 1$ leads to the start state of an instance of the class \mathcal{M}_1 , while taking $a = 2$ leads to the
 1080 start state of an instance of the class \mathcal{M}_2 ; in both cases the transition occurs with probability one and
 1081 zero reward is collected.

1082 Furthermore, assume that the reward function given over the MDPs in \mathcal{M}_1 is shifted such that the
 1083 value of a policy that takes $a = 1$ in s_1 and then plays optimally is 1 and that the reward functions on
 1084 \mathcal{M}_2 is shifted such that the value of a policy which takes $a = 2$ initially and then plays optimally is
 1085 2ε .

1086 In addition, assume that the dataset arises from a logging policy π_{log} which takes $a = 2$ initially and
 1087 then visits all (s, a, s') uniformly.

1088 Such construction and dataset identify a sparsified MDP which coincide with \mathcal{M}_2 with the addition
 1089 of s_1 (and its transition to \mathcal{M}_2 with zero reward). Intuitively, a policy with value arbitrarily close to 1
 1090 must take action $a = 1$ which leads to \mathcal{M}_1 , which is the portion of the MDP that is unexplored in the
 1091 dataset. In this case, unless the agent collects exponentially many trajectories in the online phase, the
 1092 lower bound from [Xiao et al., 2022] implies that it is not possible to discover a policy with value
 1093 close to 1 (e.g., larger than $1/2$). On the other hand, our theorem 5.1 guarantees that π_{final} has a
 1094 value at least ε , because π_{final} is ε -optimal on the sparsified MDP—i.e., ε -optimal when restricted
 1095 to an instance on \mathcal{M}_2 —with high probability using at most $\sim H^2|\mathcal{S}|^2|\mathcal{A}|/\varepsilon^2$ trajectories. This value
 1096 is unimprovable given the lower bound of Jin et al. [2020b], which applies to the class \mathcal{M}_2 .

1097 For completeness, in the next sub-section we refine the lower bound of [Xiao et al., 2022] to handle
 1098 mixture policies.

1099 **D Technical lemmas and proofs**

Lemma D.1 (Lemma F.4 in [Dann et al., 2017]). *Let \mathcal{F}_i for $i = 1 \dots$ be a filtration and X_1, \dots, X_n be a sequence of Bernoulli random variables with $\mathbb{P}(X_i = 1 \mid \mathcal{F}_{i-1}) = P_i$ with P_i being \mathcal{F}_{i-1} -measurable and X_i being \mathcal{F}_i measurable. It holds that*

$$\mathbb{P}\left(\exists n : \sum_{t=1}^n X_t < \sum_{t=1}^n P_t/2 - W\right) \leq e^{-W}.$$

Lemma D.2. *Let \mathcal{F}_i for $i = 1 \dots$ be a filtration and X_1, \dots, X_n be a sequence of Bernoulli random variables with $\mathbb{P}(X_i = 1 \mid \mathcal{F}_{i-1}) = P_i$ with P_i being \mathcal{F}_{i-1} -measurable and X_i being \mathcal{F}_i measurable. It holds that*

$$\mathbb{P}\left(\exists n : \sum_{t=1}^n X_t > \sum_{t=1}^n 2P_t + W\right) \leq e^{-W}.$$

1100 *Proof.* Notice that $\frac{1}{u^2} [\exp(u) - u - 1]$ is non-decreasing on \mathbb{R} , where at zero we continuously
 1101 extend this function. For any $t \in \mathbb{N}$, since $X_t - P_t \leq 1$, we have $\exp(X_t - P_t) - (X_t - P_t) - 1 \leq$
 1102 $(X_t - P_t)^2 (e - 2) \leq (X_t - P_t)^2$. Taking expectation conditional on \mathcal{F}_{t-1} and noticing that $P_t - X_t$
 1103 is a Martingale difference sequence w.r.t. the filtration \mathcal{F}_t , we have

$$\mathbb{E}[\exp(X_t - P_t) \mid \mathcal{F}_{t-1}] \leq 1 + \mathbb{E}[(X_t - P_t)^2 \mid \mathcal{F}_{t-1}] \leq \exp\left[\mathbb{E}[(X_t - P_t)^2 \mid \mathcal{F}_{t-1}]\right] \leq \exp(P_t),$$

where the last inequality comes from the fact that conditional on \mathcal{F}_{t-1} , X_t is a Bernoulli random variable. We define $M_n := \exp[\sum_{t=1}^n (X_t - 2P_t)]$, which is a supermartingale from our derivation above. We define now the stopping time $\tau = \min\{t \in \mathbb{N} : M_t > e^W\}$ and the sequence $\tau_n = \min\{t \in \mathbb{N} : M_t > e^W \vee t \geq n\}$. Applying the convergence theorem for nonnegative supermartingales (Theorem 4.2.12 in [Durrett, 2019]), we get that $\lim_{t \rightarrow \infty} M_t$ is well-defined almost surely. Therefore, M_τ is well-defined even when $\tau = \infty$. By the optional stopping theorem for nonnegative supermartingales (Theorem 4.8.4 in [Durrett, 2019]), we have $\mathbb{E}[M_{\tau_n}] \leq \mathbb{E}[M_0] \leq 1$ for all n and applying Fatou's lemma, we obtain $\mathbb{E}[M_\tau] = \mathbb{E}[\lim_{n \rightarrow \infty} M_{\tau_n}] \leq \liminf_{n \rightarrow \infty} \mathbb{E}[M_{\tau_n}] \leq 1$. Using Markov's inequality, we can finally bound

$$\mathbb{P}\left(\exists n : \sum_{t=1}^n X_t > 2 \sum_{t=1}^n P_t + W\right) > \mathbb{P}(\tau < \infty) \leq \mathbb{P}(M_\tau > e^W) \leq e^{-W} \mathbb{E}[M_\tau] \leq e^{-W}.$$

1104

□

Lemma D.3 (Empirical Bernstein Inequality, Theorem 11 in [Maurer and Pontil, 2009]). *Let $n \geq 2$ and x_1, \dots, x_n be i.i.d random variables such that $|x_i| \leq A$ with probability 1. Let $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ and $\widehat{V}_n = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$, then with probability $1 - \delta$ we have*

$$\left| \frac{1}{n} \sum_{i=1}^n x_i - \mathbb{E}[x] \right| \leq \sqrt{\frac{2\widehat{V}_n \log(2/\delta)}{n}} + \frac{14A}{3n} \log(2/\delta)$$

Lemma D.4 (Concentration for KL Divergence, Proposition 1 in [Jonsson et al., 2020]). *Let $X_1, X_2, \dots, X_n, \dots$ be i.i.d. samples from a distribution supported over $\{1, \dots, m\}$, of probabilities given by $\mathbb{P} \in \Sigma_m$, where Σ_m is the probability simplex of dimension $m - 1$. We denote by $\widehat{\mathbb{P}}_n$ the empirical vector of probabilities. Then, for any $\delta \in [0, 1]$, with probability at least $1 - \delta$, it holds that*

$$\forall n \in \mathbb{N}, \text{KL}\left(\widehat{\mathbb{P}}_n, \mathbb{P}\right) \leq \frac{1}{n} \log\left(\frac{1}{\delta}\right) + \frac{m}{n} \log\left(e\left(1 + \frac{n}{m}\right)\right).$$

1105 *We remark that there is a slight difference between it and the original version. In [Jonsson et al.,*
 1106 *2020], they use $m - 1$ instead of m . But since the second term of the right hand side above is*
 1107 *increasing with m , our version also holds.*

Lemma D.5 (Bernstein Transportation, Lemma 11 in [Talebi and Maillard, 2018]). *For any function f and any two probability measure \mathbb{Q}, \mathbb{P} which satisfy $\mathbb{Q} \ll \mathbb{P}$, we denote $\mathbb{V}_P[f] := \text{Var}_{X \sim \mathbb{P}}(f(X))$ and $\mathbb{S}(f) := \sup_x f(x) - \inf_x f(x)$. We assume $\mathbb{V}_P[f]$ and $\mathbb{S}(f)$ are finite, then we have*

$$\begin{aligned}\mathbb{E}_Q[f] - \mathbb{E}_P[f] &\leq \sqrt{2\mathbb{V}_P[f]\text{KL}(Q, P)} + \frac{2}{3}\mathbb{S}(f)\text{KL}(Q, P), \\ \mathbb{E}_P[f] - \mathbb{E}_Q[f] &\leq \sqrt{2\mathbb{V}_P[f]\text{KL}(Q, P)}.\end{aligned}$$

Lemma D.6 (Difference of Variance, Lemma 12 in [Ménard et al., 2021]). *Let \mathbb{P}, \mathbb{Q} be two probability measure on a discrete sample space of cardinality \mathcal{S} . Let f, g be two functions defined on \mathcal{S} such that $0 \leq g(s), f(s) \leq b$ for all $s \in \mathcal{S}$, we have that*

$$\begin{aligned}\text{Var}_{\mathbb{P}}(f) &\leq 2\text{Var}_{\mathbb{P}}(g) + 2b\mathbb{E}_{\mathbb{P}}|f - g| \quad \text{and} \\ \text{Var}_{\mathbb{Q}}(f) &\leq \text{Var}_{\mathbb{Q}}(f) + 3b^2\|\mathbb{P} - \mathbb{Q}\|_1,\end{aligned}$$

1108 Further, if $\text{KL}(\mathbb{P}; \mathbb{Q}) \leq \alpha$, it holds that

$$\text{Var}_{\mathbb{Q}}(f) \leq 2\text{Var}_{\mathbb{P}}(f) + 4b^2\alpha.$$

Lemma D.7. *For any sequence of numbers z_1, \dots, z_n with $0 \leq z_k \leq 1$, we have*

$$\sum_{k=1}^n \frac{z_k}{\max\left[1; \sum_{i=1}^{k-1} z_i\right]} \leq 4 \log\left(\sum_{i=1}^n z_i + 1\right)$$

Proof.

$$\begin{aligned}\sum_{k=1}^n \frac{z_k}{\max\left[1; \sum_{i=1}^{k-1} z_i\right]} &\leq 4 \sum_{k=1}^n \frac{\sum_{i=1}^k z_i - \sum_{i=1}^{k-1} z_i}{2 + 2\sum_{i=1}^{k-1} z_i} \\ &\leq 4 \sum_{k=1}^n \frac{\sum_{i=1}^k z_i - \sum_{i=1}^{k-1} z_i}{1 + \sum_{i=1}^k z_i} \\ &\leq 4 \sum_{k=1}^n \int_{\sum_{i=1}^{k-1} z_i}^{\sum_{i=1}^k z_i} \frac{1}{1+x} dx \leq 4 \log\left(\sum_{i=1}^n z_i + 1\right).\end{aligned}$$

1109

□

1110 **Lemma D.8.** *For $B \geq 16$ and $x \geq 3$, there exists a universal constant $C_1 \geq 4$, such that when*

$$x \geq C_1 B \log(B)^2,$$

1111 *it holds that*

$$B \log(1+x)(1 + \log(1+x)) \leq x.$$

1112 *Proof.* We have

$$B \log(1+x)(1 + \log(1+x)) \leq B(1 + \log(1+x))^2 \leq B(1 + \log(2x))^2 \leq B[\log(6x)]^2.$$

1113 We define $f(x) := x - B[\log(6x)]^2$, then we have $f'(x) = 1 - \frac{2B \log(6x)}{x}$. Since $x \geq 2C_1 B \log(B)$,
1114 we have

$$f'(x) \geq 1 - \frac{\log(12C_1 B \log(B))}{C_1 \log(B)}.$$

1115 We can take $C_1 \geq 1$ such that $C_1 \log(B) - \log(12C_1 B \log(B)) \geq 0$ whenever $B \geq 16$. Therefore,
1116 we know $f(x)$ is increasing when $x \geq C_1 B \log(B)^2$. Then, it suffices to prove

$$[\log(6C_1 B \log(B)^2)]^2 \leq C_1 \log(B)^2.$$

1117 Since $[\log(6C_1 B \log(B)^2)]^2 \leq 2\log(B)^2 + 2[\log(6C_1 \log(B)^2)]^2$, it suffices to prove

$$\log(6C_1 \log(B)^2) \leq \sqrt{\frac{C_1 - 2}{2}} \log(B).$$

1118 When $C_1 \geq 4$, the difference of right hand side and left hand side is always increasing w.r.t. B for
1119 fixed C_1 . Therefore, it suffices to prove the case when $B = 16$. Noticing that we can always take
1120 a sufficiently large uniform constant C_1 such that the inequality above holds when $B = 16$, we
1121 conclude. □

Lemma D.9 (Chain rule of Kullback-Leibler divergence, Exercise 3.2 in [Wainwright, 2019]). *Given two n -variate distributions \mathbb{Q} and \mathbb{P} , show that the Kullback-Leibler divergence can be decomposed as*

$$D(\mathbb{Q}; \mathbb{P}) = D(\mathbb{Q}_1; \mathbb{P}_1) + \sum_{j=2}^n \mathbb{E}_{\mathbb{Q}_1^{j-1}} \left[D \left(\mathbb{Q}_j \left(\cdot \mid X_1^{j-1} \right); \mathbb{P}_j \left(\cdot \mid X_1^{j-1} \right) \right) \right],$$

1122 where $\mathbb{Q}_j \left(\cdot \mid X_1^{j-1} \right)$ denotes the conditional distribution of X_j given (X_1, \dots, X_{j-1}) under \mathbb{Q} ,
 1123 with a similar definition for $\mathbb{P}_j \left(\cdot \mid X_1^{j-1} \right)$.

Lemma D.10 (Bretagnolle-Huber Inequality, Theorem 14.1 in [Lattimore and Szepesvári, 2020]). *Let \mathbb{P} and \mathbb{Q} be probability measures on the same measurable space (Ω, \mathcal{F}) , and let $A \in \mathcal{F}$ be an arbitrary event. Then,*

$$\mathbb{P}(A) + \mathbb{Q}(A^c) \geq \frac{1}{2} \exp(-D(\mathbb{P}; \mathbb{Q}))$$

1124 where $A^c = \Omega \setminus A$ is the complement of A .

1125 **Lemma D.11** (Value Difference Lemma, Lemma E.15 in [Dann et al., 2017]). *For any two MDPs*
 1126 \mathcal{M}' *and* \mathcal{M}'' *with rewards* r' *and* r'' *and transition probabilities* \mathbb{P}' *and* \mathbb{P}'' , *the difference in values*
 1127 *with respect to the same policy* π *can be written as*

$$\begin{aligned} V'_i(s) - V''_i(s) &= \mathbb{E}'' \left[\sum_{t=i}^H (r'(s_t, a_t, t) - r''(s_t, a_t, t)) \mid s_i = s \right] \\ &+ \mathbb{E}'' \left[\sum_{t=i}^H (\mathbb{P}'(s_t, a_t, t) - \mathbb{P}''(s_t, a_t, t))^\top V'_{t+1} \mid s_i = s \right] \end{aligned}$$

1128 where $V'_{H+1}(s) = V''_{H+1}(s) = 0$ for any state s and the expectation \mathbb{E}' is taken with respect to \mathbb{P}'
 1129 and π and \mathbb{E}'' with respect to \mathbb{P}'' and π .

1130 E Details of the planning phase

1131 In this section, we provide some details of the planning phase in algorithm 3. In the planning phase,
 1132 we are given a reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ and we compute an estimate of sparsified transition
 1133 dynamics $\tilde{\mathbb{P}}^\dagger$, which is formally defined appendix A.1.1. The goal of the planning phase is to compute
 1134 the optimal policy π_{final} on the MDP specified by the transition dynamics $\tilde{\mathbb{P}}^\dagger$ and reward function
 1135 r^\dagger , where r^\dagger is the sparsified version of $r : r^\dagger(s, a) = r(s, a)$ and $r^\dagger(s^\dagger, a) = 0$ for any $a \in \mathcal{A}$. To
 1136 compute the optimal policy, we iteratively apply the Bellman optimality equation. First, we define
 1137 $\tilde{Q}_H(s, a) = r^\dagger(s, a)$ for any (s, a) and solve

$$\pi_{final, H}(s) = \arg \max_{a \in \mathcal{A}} r^\dagger(s, a).$$

1138 Then, for $h = H - 1, H - 2, \dots, 2, 1$, we iteratively define

$$\tilde{Q}_h(s, a) := r^\dagger(s, a) + \sum_{s'} \tilde{\mathbb{P}}^\dagger(s' \mid s, a) \tilde{Q}_{h+1}(s', \pi_{final, h+1}(s'))$$

1139 for any (s, a) , and then solve

$$\pi_{final, h}(s) = \arg \max_{a \in \mathcal{A}} \tilde{Q}_h(s, a)$$

1140 for any $s \in \mathcal{S}$. For s^\dagger and any $h \in [H]$, $\pi_{final, h}(s^\dagger)$ can be arbitrary action. Then, from the property
 1141 of Bellman optimality equation, we know π_{final} is the optimal policy on $\tilde{\mathbb{P}}^\dagger$ and r^\dagger .

1142 F More related works

1143 **Other low-switching algorithms** Low-switching learning algorithms were initially studied in the
 1144 context of bandits, with the UCB2 algorithm [Auer et al., 2002] achieving an $O(\mathcal{A} \log K)$ switching

1145 cost. Gao et al. [2019] demonstrated a sufficient and necessary $O(\mathcal{A} \log \log K)$ switching cost for
1146 attaining the minimax optimal regret in multi-armed bandits. In both adversarial and stochastic online
1147 learning, [Cesa-Bianchi et al., 2013] designed an algorithm that achieves an $O(\log \log K)$ switching
1148 cost.

1149 **Reward-free reinforcement learning** In reward-free reinforcement learning (RFRL) the goal is to
1150 find a near-optimal policy for any given reward function. [Jin et al., 2020a] proposed an algorithm
1151 based on EULER [Zanette and Brunskill, 2019] that can find a ε policy with $\tilde{O}(H^5 |\mathcal{S}|^2 |\mathcal{A}| / \varepsilon^2)$
1152 trajectories. Subsequently, [Kaufmann et al., 2021] reduces the sample complexity by a factor H by
1153 using uncertainty functions to upper bound the value estimation error. The sample complexity was
1154 further improved by another H factor by [Ménard et al., 2021].

1155 A lower bound of $\tilde{O}(H^2 |\mathcal{S}|^2 |\mathcal{A}| / \varepsilon^2)$ was established for homogeneous MDPs by [Jin et al., 2020a],
1156 while an additional H factor is conjectured for non-homogeneous cases. [Zhang et al., 2021] proposed
1157 the first algorithm with matching sample complexity in the homogeneous case. Similar results are
1158 available with linear [Wang et al., 2020a, Wagenmaker et al., 2022, Zanette et al., 2020] and general
1159 function approximation [Chen et al., 2022, Qiu et al., 2021].

1160 **Offline reinforcement learning** In offline reinforcement learning the goal is to learn a near-optimal
1161 policy from an existing dataset which is generated from a (possibly very different) logging policy.
1162 Offline RL in tabular domains has been investigated extensively [Yin and Wang, 2020, Jin et al.,
1163 2020c, Nachum et al., 2019, Rashidinejad et al., 2021, Kallus and Uehara, 2022, Xie and Jiang,
1164 2020a]. Similar results were shown using linear [Yin et al., Xiong et al., 2022, Nguyen-Tang et al.,
1165 2022, Zanette et al., 2021b] and general function approximation [Xie et al., 2021a, Long et al., 2021,
1166 Zhang et al., 2022, Duan et al., 2021, Jiang and Huang, 2020, Uehara and Sun, 2021, Zanette
1167 and Wainwright, 2022, Rashidinejad et al., 2022, Yin et al., 2022]. Offline RL is effective when
1168 the dataset ‘covers’ a near optimal policy, as measured by a certain concentrability factor. In the
1169 function approximation setting additional conditions, such as Bellman completeness, may need to be
1170 approximately satisfied [Munos and Szepesvári, 2008, Chen and Jiang, 2019, Zanette, 2023, Wang
1171 et al., 2020b, Foster et al., 2021, Zhan et al., 2022].

1172 **Task-agnostic reinforcement learning** Another related line of work is task-agnostic RL, where
1173 N tasks are considered during the planning phase, and the reward functions is collected from the
1174 environment instead of being provided directly. [Zhang et al., 2020a] presented the first task-agnostic
1175 algorithm, UBEZero, with a sample complexity of $\tilde{O}(H^5 |\mathcal{S}| |\mathcal{A}| \log(N) / \varepsilon^2)$. Recently, [Li et al.,
1176 2023a] proposed an algorithm that leverages offline RL techniques to estimate a well-behaved
1177 behavior policy in the reward-agnostic phase, achieving minimax sample complexity. Other works
1178 exploring effective exploration schemes in RL include [Hazan et al., 2019, Du et al., 2019, Misra
1179 et al., 2020]. [Li et al., 2023b] also considered an offline-to-online reinforcement learning algorithm
1180 which explores the environment using two mixed policies in a reward-free mode.