

---

# Misspecified $Q$ -Learning with Sparse Linear Function Approximation: Tight Bounds on Approximation Error

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       The recent work by Dong and Yang [2023] showed for misspecified sparse linear  
2       bandits, one can obtain an  $O(\epsilon)$ -optimal policy using a polynomial number of  
3       samples when the sparsity is a constant, where  $\epsilon$  is the misspecification error.  
4       This result is in sharp contrast to misspecified linear bandits without sparsity,  
5       which require an exponential number of samples to get the same guarantee. In  
6       order to study whether the analog result is possible in the reinforcement learning  
7       setting, we consider the following problem: assuming the optimal  $Q$ -function is a  
8        $d$ -dimensional linear function with sparsity  $s$  and misspecification error  $\epsilon$ , whether  
9       we can obtain an  $O(\epsilon)$ -optimal policy using number of samples polynomially  
10      in the feature dimension  $d$ . We first demonstrate why the standard approach  
11      based on Bellman backup or the existing optimistic hypothesis class elimination  
12      approach such as OLIVE Jiang et al. [2017] achieves suboptimal guarantees for  
13      this problem. We then design a novel elimination-based algorithm to show one  
14      can obtain an  $O(H\epsilon)$ -optimal policy with sample complexity polynomially in the  
15      feature dimension  $d$  and planning horizon  $H$ . Lastly, we complement our upper  
16      bound with an  $\tilde{\Omega}(H\epsilon)$  suboptimality lower bound, giving a complete picture of  
17      this problem.

## 18   1 Introduction

19   Bandit and reinforcement learning (RL) problems in real-world applications, such as autonomous  
20   driving [Kiran et al., 2021], healthcare [Esteva et al., 2019], recommendation systems [Bouneffouf  
21   et al., 2012], and advertising [Schwartz et al., 2017], face challenges due to the vast state-action  
22   space. To tackle this, function approximation frameworks, such as using linear functions or neural  
23   networks, have been introduced to approximate the value functions or policies. However, real-world  
24   complexities often mean that function approximation is agnostic; the function class captures only an  
25   approximate version of the optimal value function, and the misspecification error remains unknown.  
26   A fundamental problem is understanding the impact of agnostic misspecification errors in RL.

27   Prior works show even minor misspecifications can lead to exponential (in dimension) sample  
28   complexity in the linear bandit settings [Du et al., 2020, Lattimore et al., 2020] if the goal is to  
29   learn a policy within the misspecification error. That is, finding an  $O(\epsilon)$ -optimal action necessitates  
30   at least  $\Omega(\exp(d))$  queries (or samples) to the environment, where  $\epsilon$  is the misspecification error.  
31   Recently, Dong and Yang [2023] demonstrated that by leveraging the sparsity structure of ground-  
32   truth parameters, one can overcome the exponential sample barrier in the linear bandit setting. They  
33   showed that with sparsity  $k$  in the ground-truth parameters, it is possible to learn an  $O(\epsilon)$ -optimal

34 action with only  $O\left((d/\epsilon)^k\right)$  samples. In particular, when  $k$  is a constant, their algorithm achieves a  
 35 polynomial sample complexity guarantee.

36 A natural question is whether we can obtain a similar sample complexity guarantee in the RL setting.  
 37 This question motivates us to consider a more general question:

38 *Given a that the  $Q^*$  function is a  $d$ -dimensional linear function with sparsity  $k$  and misspecification  
 39 error  $\epsilon$ , can we learn an  $O(\epsilon)$ -optimal policy using  $\text{poly}(d, 1/\epsilon)$  samples, when  $k$  is a constant?*

40 It turns out that by studying this question, we obtain a series of surprising results which cannot be  
 41 explained by existing RL theories.

## 42 1.1 Our Contributions

43 In this paper, we propose an RL algorithm that can handle linear function approximation with sparsity  
 44 structures and misspecification errors. We also show that the suboptimality achieved by our algorithm  
 45 is near-optimal, by proving information-theoretic hardness results. Here we give a more detailed  
 46 description of our technical contributions.

47 **Our Assumption.** Throughout this paper, we assume the RL algorithm has access to a feature map  
 48 where, for each state-action pair  $(s, a)$ , we have the feature  $\phi(s, a)$  with  $\|\phi(s, a)\| \leq 1$ . We make the  
 49 following assumption, which states that there exists a sequence of parameters  $\theta^* = (\theta_0^*, \dots, \theta_{H-1}^*)$   
 50 where each  $\theta_h^* \in \mathbb{S}^{d-1}$  is  $k$ -sparse, that approximates the optimal  $Q$ -function up to an error of  $\epsilon$ .

51 **Assumption 1.** *There exists  $\theta^* = (\theta_0^*, \dots, \theta_{H-1}^*)$  where each  $\theta_h^* \in \mathbb{S}^{d-1}$  is  $k$ -sparse, such that*

$$|\langle \phi(s, a), \theta_h^* \rangle - Q^*(s, a)| \leq \epsilon$$

52 *for all  $h \in [H]$ , all states  $s$  in level  $h$ , and all actions  $a$  in the action space.*

53 When  $H = 1$ , Assumption 1 is equivalent to the bandit setting in Dong and Yang [2023]. Note that  
 54 we have a different optimal parameter  $\theta_h$  for each level  $h \in [H]$ .

55 We can approximate  $\theta^*$  using an  $\epsilon$ -net of the sphere  $\mathbb{S}^{k-1}$  and the set of all  $k$ -sized subset of  $[d]$ .  
 56 Therefore, when  $k$  is a constant, we may assume that each  $\theta_h^*$  lies in a set with size polynomial in  $d$ .  
 57 Then, a natural idea is to enumerate all possible policies induced by the parameters in that finite set,  
 58 and choose the one with the highest cumulative reward. However, although the number of parameter  
 59 candidates in each individual level has polynomial size, the total number of induced policies would  
 60 be exponential in  $H$ , and the sample complexity of such approach would also be exponential in  $H$ .

61 **The Level-by-level Approach.** Note that when the horizon length  $H = 1$ , the problem under  
 62 consideration is equivalent to a bandit problem, which can be solved by previous approaches [Dong  
 63 and Yang, 2023]. For the RL setting, a natural idea is to first apply the bandit algorithm in Dong  
 64 and Yang [2023] on the last level, and then apply the same bandit algorithm on the second last level  
 65 based on previous results and Bellman-backups, and so on. However, we note that to employ such an  
 66 approach, the bandit algorithm needs to provide a “for-all” guarantee, i.e., finding a parameter that  
 67 approximates the rewards of all arms, instead of just finding a near-optimal arm. On the other hand,  
 68 existing bandit algorithms will amplify the approximation error of the input parameters by a constant  
 69 factor, in order to provide a for-all guarantee. Concretely, existing bandit algorithms can only find a  
 70 parameter  $\theta$  so that  $\theta$  approximates the rewards of all arms by an error of  $2\epsilon$ . As we have  $H$  levels  
 71 in the RL setting, the final error would be exponential in  $H$ , and therefore, such a level-by-level  
 72 approach would result in a suboptimality that is exponential in  $H$ .

73 One may ask if we can further improve existing bandit algorithms, so that we can find a parameter  $\theta$   
 74 that approximates the rewards of all arms by an error of  $\epsilon$  plus a statistical error that can be made  
 75 arbitrarily small, instead of  $2\epsilon$ . The following theorem shows that this is information-theoretically  
 76 impossible unless one pays a sample complexity proportional to the size of the action space.

77 **Theorem 1.1.** *Under Assumption 1 with  $d = k = 1$ , any bandit algorithm that returns an estimate  
 78  $\hat{r}$  such that  $|\hat{r}(a) - r(a)| < 2\epsilon$  for all arms  $a$  with probability at least 0.95 requires at least  $0.9n$   
 79 samples, where  $n$  is the total number of arms.*

80 Therefore, amplifying the approximation error by a factor of 2 is not an artifact of existing bandit  
 81 algorithms. Instead, it is information-theoretically impossible.

82 Geometric error amplification is a common issue in the design of RL algorithm with linear function  
 83 approximation [Zanette et al., 2019, Weisz et al., 2021, Wang et al., 2020a, 2021a]. It is interesting  
 84 (and also surprising) that such an issue arises even when the function class has sparsity structures.

85 **Optimistic Value Function Elimination.** Another approach for the design of RL algorithm is based  
 86 on optimistic value function elimination. Such an approach was proposed by Jiang et al. [2017] and  
 87 was then generalized to broader settings [Sun et al., 2019, Du et al., 2021, Jin et al., 2021, Chen et al.,  
 88 2022b]. At each iteration of the algorithm, we pick the value functions in the hypothesis class with  
 89 maximized value. We then use the induced policy to collect a dataset, based on which we eliminate a  
 90 bunch of value functions from the hypothesis class and proceed to the next iteration.

91 When applied to our setting, existing algorithms and analysis achieve a suboptimality that depends  
 92 on the size of the parameter class, which could be prohibitively large. Here, we use the result in Jiang  
 93 et al. [2017] as an example. The suboptimality of their algorithm is  $H\sqrt{M}\epsilon$ , where  $M$  is *Bellman*  
 94 *rank* of the problem. For our setting, we can show that there exists an MDP instance and a feature  
 95 map that satisfies Assumption 1, whose induced Bellman rank is large.

96 **Proposition 1.2.** *There exists an MDP instance  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, P, r)$  with  $|\mathcal{A}| = 2$ ,  $H = \log d$ ,  
 97  $|\mathcal{S}| = d - 1$ , with  $d$ -dimensional feature map  $\phi$  satisfying Assumption 1 with  $k = 1$ , such that its  
 98 Bellman rank is  $d$ .*

99 Given Proposition 1.2, if one naïvely applies the algorithm in Jiang et al. [2017], the suboptimality  
 100 would be  $O(H\sqrt{d}\epsilon)$  in our setting, which necessitates new algorithm and analysis. In Section 4, we  
 101 design a new RL algorithm whose performance is summarized in the following theorem.

102 **Theorem 1.3.** *Under Assumption 1, with probability at least  $1 - \delta$ , Algorithm 1 returns a policy  
 103 with suboptimality at most  $(4\epsilon_{\text{stat}} + 2\epsilon_{\text{net}} + 2\epsilon)H$  by taking  $O(kd^k H^3 \cdot \ln(dH/\epsilon_{\text{net}}\delta) \cdot \epsilon_{\text{net}}^{-k} \epsilon_{\text{stat}}^{-2})$   
 104 samples.*

105 Here  $\epsilon_{\text{stat}}$  is the statistical error. Compared to the existing approaches, Theorem 1.3 achieves a much  
 106 stronger suboptimality guarantee. Later, we will also show that such a guarantee is near-optimal.

107 Although based on the same idea of optimistic value function elimination, our proposed algorithm  
 108 differs significantly from existing approaches [Jiang et al., 2017, Sun et al., 2019, Du et al., 2021,  
 109 Jin et al., 2021, Chen et al., 2022b] to exploit the sparsity structure. While existing approaches  
 110 based on optimistic value function elimination try to find a sequence of parameters that maximize the  
 111 value of the initial states, our new algorithm selects a parameter that maximizes the empirical roll-in  
 112 distribution at all levels. Also, existing algorithms eliminate a large set of parameters in each iteration,  
 113 while we only eliminate the parameters selected during the current iteration in our algorithm.

114 These two modifications are crucial for obtaining a smaller suboptimality guarantee, smaller sample  
 115 complexity, and shorter running time. In existing algorithms, parameters at different levels are  
 116 interdependent, i.e. the choice of parameter at level  $h$  affects the choice of parameter at level  $h + 1$ .  
 117 Our new algorithm simplifies this by maintaining a parameter set for each level, so each level operates  
 118 independently. Further, we can falsify and eliminate any parameter showing large Bellman error  
 119 at any level  $h$ , since otherwise we would have found another parameter with larger induced value  
 120 function at level  $h + 1$  to make the error small. Consequently, since we eliminate at least one  
 121 parameter at each iteration, we obtain fewer iterations and enhanced sample complexity.

122 **The Hardness Result.** One may wonder if the suboptimality guarantee can be further improved. In  
 123 Section 3, we show that the suboptimality guarantee by Theorem 1.3 is near-optimal.

124 We first consider a weaker setting where the algorithm is not allowed to take samples, and the function  
 125 class contains a single sequence of functions. I.e., we are given a function  $\hat{Q} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , such that  
 126  $|\hat{Q}(s, a) - Q^*(s, a)| \leq \epsilon$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .

127 We show that for this weaker setting, simply choosing the greedy policy with respect to  $\hat{Q}$ , which  
 128 achieves a suboptimality guarantee of  $O(H\epsilon)$ , is actually optimal. To prove this, we construct a hard  
 129 instance based on binary tree. Roughly speaking, the optimal action for each level is chosen uniformly  
 130 random from two actions  $a_1$  and  $a_2$ . At all levels, the reward is  $\epsilon$  if the optimal action is chosen,  
 131 and 0 otherwise. For this instance, there exists a fixed  $\hat{Q}$  that provides a good approximation to  
 132 the optimal  $Q$ -function, regardless of the choice of the optimal actions. Therefore,  $\hat{Q}$  reveals no  
 133 information about the optimal actions, and the suboptimality of the returned policy would be at least  
 134  $\Omega(H\epsilon)$ . The formal construction and analysis and construction will be given in Section 3.1.

135 When the algorithm is allowed to take samples, we show that in order to achieve a suboptimality  
136 guarantee of  $H/C\epsilon$ , any algorithm requires  $\exp(\Omega(C))$  samples, even when Assumption 1 is satisfied  
137 with  $d = k = 1$ . Therefore, for RL algorithms with polynomial sample complexity, the suboptimality  
138 guarantee of Theorem 1.3 is tight up to log factors.

139 To prove the above claim, we still consider the setting where  $d = k = 1$ , i.e., a good approximation to  
140 the  $Q$ -function is given to the algorithm. We also use a more complicated binary tree instance, where  
141 we divide all the  $H$  levels into  $H/C$  blocks, each containing  $C$  levels. For each block, only one  
142 state-action pair at the last level has a reward of  $\epsilon$ , and all other state-action pairs in the block has a  
143 reward of 0. Therefore, the value of the optimal policy would be  $H/C \cdot \epsilon$  since there are  $H/C$  blocks  
144 in total. We further show that there is a fixed function  $\hat{Q}$ , which provides a good approximation to the  
145 optimal  $Q$ -function universally for all instances under consideration.

146 Since  $\hat{Q}$  reveals no information about the state-action pair with  $\epsilon$  reward for all blocks, for an RL  
147 algorithm to return a policy with a non-zero value, it must search for a state-action pair with non-zero  
148 reward in a brute force manner, which inevitably incurs a sample complexity of  $\exp(\Omega(C))$  since each  
149 block contains  $C$  levels and  $\exp(\Omega(C))$  state-action pairs at the last level. The formal construction  
150 and analysis and construction will be given in Section 3.2.

## 151 1.2 Related Work

152 A series of studies have delved into MDPs that can be represented by linear functions of predetermined  
153 feature mappings, achieving sample complexity or regret that depends on the feature mapping’s  
154 dimension. This includes linear MDPs, studied in Jin et al. [2020], Wang et al. [2019], Neu and  
155 Pike-Burke [2020], where both transition probabilities and rewards are linear functions of feature  
156 mappings on state-action pairs. Zanette et al. [2020a,b] examines MDPs with low inherent Bellman  
157 error, indicating value functions that are almost linear with respect to these mappings. Another focus  
158 is on linear mixture MDPs [Modi et al., 2020, Jia et al., 2020, Ayoub et al., 2020, Zhou et al., 2021,  
159 Cai et al., 2020], characterized by transition probabilities that combine several basis kernels linearly.  
160 While these studies often assume known feature vectors, Agarwal et al. [2020] investigates a more  
161 challenging scenario where both features and parameters of the linear model are unknown.

162 The literature has also witnessed a substantial surge of research in understanding how function general  
163 approximations can be applied efficiently in the reinforcement learning setting [Osband and Van Roy,  
164 2014, Sun et al., 2019, Ayoub et al., 2020, Wang et al., 2020b, Foster et al., 2021, Chen et al., 2022b,a,  
165 Zhong et al., 2022, Foster et al., 2023, Wagenmaker and Foster, 2023, Zhou and Gu, 2022, Jiang et al.,  
166 2017, Wang et al., 2020b, Du et al., 2021, Jin et al., 2021, Kong et al., 2021, Dann et al., 2021, Zhong  
167 et al., 2022, Liu et al., 2023, Agarwal et al., 2023]. To obtain good sample, error, or regret bounds,  
168 these approaches typically impose benign structures on values, models, or policies, along with benign  
169 misspecification. Amongst these works, Jiang et al. [2017] is particularly related to our work as their  
170 elimination-based algorithm, OLIVE, can be directly applied to our setting. However, as mentioned  
171 in Section 1.1, the suboptimality guarantee of their algorithm is significantly worse than our result.

172 In another line of works, Du et al. [2020], Dong and Yang [2023], Lattimore et al. [2020] specifically  
173 focuses on understanding misspecification in bandit and RL scenarios. Du et al. [2020] illustrated  
174 that to find an  $O(\epsilon)$ -optimal policy in reinforcement learning with  $\epsilon$ -misspecified linear features,  
175 an agent must sample an exponential (in  $d$ ) number of trajectories, applicable to both value-based  
176 and model-based learning. Relaxing this goal, Lattimore et al. [2020] indicated that  $\text{poly}(d/\epsilon)$   
177 samples could suffice to secure an  $O(\epsilon\sqrt{d})$ -optimal policy in a simulator model setting of RL, though  
178 achieving a policy with an error better than  $O(\epsilon\sqrt{d})$  would still require an exponential sample  
179 size. Recently, Dong and Yang [2023] introduced a solution, showing that incorporating structural  
180 information like sparsity in the bandit instance could address this issue, making it feasible to attain  
181  $O(\epsilon)$  with  $O((d/\epsilon)^k)$  sample complexity, which is acceptable when the sparsity  $k$  is small. Another  
182 recent independent work [Amortila et al., 2024] also obtains a suboptimality guarantee of  $O(H\epsilon)$ .  
183 However, their result depends on a coverability assumption and uses a different technique called  
184 disagreement-based regression (DBR), which is distinct from our assumption and techniques.

185 **2 Preliminaries**

186 Throughout the paper, for a given positive integer  $n$ , we use  $[n]$  to denote the set  $\{0, 1, 2, \dots, n-1\}$ .  
 187 In addition,  $f(n) = O(g(n))$  denotes that there exists a constant  $c > 0$  such that  $|f(n)| \leq c|g(n)|$ .  
 188  $f(n) = \Omega(g(n))$  denotes that there exists a constant  $c > 0$  such that  $|f(n)| \geq c|g(n)|$ .

189 **2.1 Reinforcement Learning**

190 Let  $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, H, P, r\}$  be a Markov Decision Process (MDP) where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the  
 191 action space,  $H \in \mathbb{Z}_+$  is the planning horizon,  $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  is the transition kernel which  
 192 takes a state-action pair as input and returns a distribution over states,  $r : \mathcal{S} \times \mathcal{A} \rightarrow \Delta([0, 1])$  is the  
 193 reward distribution. We assume  $\sum_{h \in [H]} r_h \in [0, 1]$  almost surely. For simplicity, throughout this  
 194 paper, we assume the initial state  $s_0$  is deterministic. To streamline our analysis, for each  $h \in [H]$ ,  
 195 we use  $\mathcal{S}_h \subseteq \mathcal{S}$  to denote the set of states at level  $h$ , and assume  $\mathcal{S}_h$  do not intersect with each other.

196 A policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  chooses an action for each state, and may induce a trajectory denoted by  
 197  $(s_0, a_0, r_0, \dots, s_{H-1}, a_{H-1}, r_{H-1})$ , where  $s_{h+1} \sim P(s_h, a_h)$ ,  $a_h = \pi(s_h)$ , and  $r_h \sim r(s_h, a_h)$  for  
 198 all  $h \in [H]$ . Given a policy  $\pi$  and  $h \in [H]$ , for a state-action pair  $(s, a) \in \mathcal{S}_h \times \mathcal{A}$ , the  $Q$ -function  
 199 and value function is defined as

$$Q^\pi(s, a) = \mathbb{E} \left[ \sum_{h'=h}^{H-1} r(s_{h'}, a_{h'}) \mid s_h = s, a_h = a, \pi \right], V^\pi(s) = \mathbb{E} \left[ \sum_{h'=h}^{H-1} r(s_{h'}, a_{h'}) \mid s_h = s, \pi \right].$$

200 We use  $V^\pi$  to denote the value of the policy  $\pi$ , i.e.,  $V^\pi = V^\pi(s_0)$ . We use  $\pi^*$  to denote the optimal  
 201 policy. For simplicity, for a state  $s \in \mathcal{S}$ , we define  $V^*(s) = V^{\pi^*}(s)$ , and for a state-action pair  
 202  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , we define  $Q^*(s, a) = Q^{\pi^*}(s, a)$ . The suboptimality of a policy  $\pi$  is defined as the  
 203 difference between the value of  $\pi$  and that of  $\pi^*$ , i.e.  $V^* - V^\pi$ .

204 For any sequence of  $k$ -sparse parameter  $\theta = (\theta_0, \dots, \theta_{H-1})$ , we define  $\pi_\theta$  to be the greedy strategy  
 205 based on  $\theta$ . In other words, for each  $h \in [H]$ , for a state  $s \in \mathcal{S}_h$ ,  $\pi_\theta(s) = \arg \max_{a \in \mathcal{A}} \langle \phi(s, a), \theta_h \rangle$ .  
 206 For each  $h \in [H]$ , a parameter  $\theta_h$ , and a state  $s \in \mathcal{S}_h$ , we also write  $V_{\theta_h}(s) = \max_{a \in \mathcal{A}} \langle \phi(s, a), \theta_h \rangle$ .

207 We will prove lower bounds for deterministic systems, i.e., MDPs with deterministic transition  $P$  and  
 208 deterministic reward  $r$ . In this setting,  $P$  and  $r$  can be regarded as functions rather than distributions.  
 209 Since deterministic systems can be considered as a special case for general stochastic MDPs, our  
 210 lower bounds still hold for general MDPs.

211 **Interacting with an MDP.** An RL algorithm takes the feature function  $\phi$  and sparsity  $k$  as the  
 212 input, and interacts with the underlying MDP by taking samples in the form of a trajectory. To  
 213 be more specific, at each round, the RL algorithm decides a policy  $\pi$  and receives a trajectory  
 214  $(s_0, a_0, r_0, \dots, s_{H-1}, a_{H-1}, r_{H-1})$  as feedback. Here one trajectory corresponds to  $H$  samples. We  
 215 define the total number of samples required by an RL algorithm as its *sample complexity*. Our goal is  
 216 to design an algorithm that returns a near-optimal policy while minimizing its sample complexity.

217 **The Bandits Setting.** In this paper, we also consider the bandit setting, which is equivalent to an  
 218 MDP with  $H = 1$ . Let  $\mathcal{A}$  be the action space, and  $r : \mathcal{A} \rightarrow \Delta([0, 1])$  be the reward distribution.  
 219 At round  $t$ , the algorithm chooses an action  $a_t \in \mathcal{A}$  and receives a reward  $r_t \sim r(a_t)$ . In this case,  
 220 Assumption 1 asserts that there exists  $\theta^*$ , such that  $|\langle \phi(a), \theta^* \rangle - \mathbb{E}[r(a)]| \leq \epsilon$  for all  $a \in \mathcal{A}$ .

221 **3 Hardness Results**

222 We prove our hardness results. In Section 3.1, we prove that the suboptimality of any RL algorithm  
 223 is  $\Omega(H\epsilon)$  if the algorithm is not allowed to take samples. This serves as a warmup for the more  
 224 complicated construction in Section 3.2, where we show that for any  $C$  satisfying  $1 \leq 2C \leq H$ , any  
 225 RL algorithm requires  $\exp(\Omega(C))$  samples in order to achieve a suboptimality of  $\Omega(H/C \cdot \epsilon)$ .

226 **3.1 Warmup: Hardness Result for RL without Samples**

227 We prove that the suboptimality of any RL algorithm without sample is  $\Omega(H\epsilon)$ . Specifically, we  
 228 consider a setting where the feature  $\phi$  is 1-dimensional and equal to the optimal  $Q$ -value with an

229 error of  $\epsilon$ . This provides a simplified scenario where an approximate optimal Q-function is readily  
 230 available to the algorithm. Theorem 3.1 suggests that even in such a simplified context, the best  
 231 achievable suboptimality is  $O(H\epsilon)$ .

232 **Theorem 3.1.** *Given a MDP instance satisfying Assumption 1, the suboptimality of the policy*  
 233 *returned by any RL algorithm is  $\Omega(H\epsilon)$  with a probability of 0.99 if the algorithm is not allowed*  
 234 *to take samples. This holds even when the dimension and sparsity satisfies  $d = k = 1$  and the*  
 235 *underlying MDP is a deterministic system.*

236 The formal proof of Theorem 3.1 is given in Section A.1 of the Supplementary Material. Below we  
 237 give the construction of the hard instance (illustration of the hard instance is given in Appendix A.1),  
 238 together with an overview of the hardness proof.

239 Our hardness result is based on a binary tree instance. There are  $H$  levels of states, and level  
 240  $h \in [H]$  contains  $2^h$  distinct states. Thus we have  $2^H - 1$  states in total. We use  $s_0, \dots, s_{2^H-2}$   
 241 to denote all the states, where  $s_0$  is the unique state at level 0, and  $s_1, s_2$  are the states at level 1,  
 242 etc. Equivalently,  $\mathcal{S}_h = \{s_{2^h-1}, \dots, s_{2^{h+1}-2}\}$ . The action space  $\mathcal{A}$  contains two actions,  $a_1$  and  
 243  $a_2$ . For each  $h \in [H-1]$ , a state  $s_i \in \mathcal{S}_h$ , we have  $P(s_i, a_1) = s_{2i+1}$  and  $P(s_i, a_2) = s_{2i+2}$ . For  
 244 each  $h \in [H]$ , there exists an action  $a_h^* \in \{a_1, a_2\}$ , such that  $\pi^*(s) = a_h^*$  for all  $s \in \mathcal{S}_h$ . Based  
 245 on  $a_0^*, a_1^*, \dots, a_{H-1}^*$ , for a state  $s \in \mathcal{S}_h$ , we define the reward function as  $r(s, a) = \epsilon$  if  $a = a_h^*$   
 246 and  $r(s, a) = 0$  otherwise. The corresponding Q-function is  $Q^*(s, a) = (H-h)\epsilon$  if  $a = a_h^*$  and  
 247  $Q^*(s, a) = (H-h-1)\epsilon$  otherwise.

248 Now we define the 1-dimensional feature function  $\phi$ . For each  $h \in [H]$ , for all  $(s, a) \in \mathcal{S}_h \times \mathcal{A}$ ,  
 249  $\phi(s, a) = (H-h-1)\epsilon$ . Clearly, by taking  $\theta^* = 1$ , Assumption 1 is satisfied for our  $\phi$ . This finishes  
 250 the construction of our hard instance.

251 Since the RL algorithm is not allowed to take samples, the only information that the algorithm receives  
 252 is the feature function  $\phi$ . However,  $\phi$  is always the same no matter how we set  $a_0^*, a_1^*, \dots, a_{H-1}^*$ ,  
 253 which means the RL algorithm can only output a fixed policy. On the other hand, if  $a_h^*$  is drawn  
 254 uniformly at random from  $\{a_1, a_2\}$ , for any fixed policy  $\pi$ , its expected suboptimality will be  $H\epsilon/2$ ,  
 255 which proves Theorem 3.1. Our formal proof in Section A.1 of the Supplementary Material is based  
 256 on Yao’s minimax principle in order to cope with randomized algorithms.

### 257 3.2 Hardness Result for RL with Samples

258 In this section, we show that for any  $1 \leq 2C \leq H$ , any RL algorithm requires  $\exp(\Omega(C))$  samples  
 259 in order to achieve a suboptimality of  $\Omega(H/C \cdot \epsilon)$ .

260 **Theorem 3.2.** *Given a RL problem instance satisfying Assumption 1 with misspecification  $\epsilon < 1/H$*   
 261 *and let  $C \in \mathbb{R}$  such that  $1 \leq 2C \leq H$ . Any algorithm that returns a policy with suboptimality less*  
 262 *than  $H/(2C) \cdot \epsilon$  with probability at least 0.9 needs least  $0.1 \cdot C \cdot 2^C$  samples.*

263 In the remaining part of this section, we give an overview of the proof of Theorem 3.2. We first define  
 264 the MULTI-INDEX-QUERY problem.

265 **Definition 1.** (MULTI-INDEX-QUERY) *In the  $m$ -INDQ $_n$  problem, we have a sequence of  $m$  indices*  
 266  *$(i_0^*, i_1^*, \dots, i_{m-1}^*) \in [n]^m$ . In each round, the algorithm guesses a pair  $(j, i) \in [m] \times [n]$  and queries*  
 267 *whether  $i = i_j^*$ . The goal is to output  $(j, i_j^*)$  for any  $j \in [m]$ , using as few queries as possible.*

268 **Definition 2.** ( $\delta$ -correct algorithm) *For  $\delta \in (0, 1)$ , we say a randomized algorithm  $A$  is  $\delta$ -correct for*  
 269  *$m$ -INDQ $_n$  if for any  $i^* = \{i_j^*\}_{j \in [m]}$ , with probability at least  $1 - \delta$ ,  $A$  outputs  $(j, i_j^*)$  for some  $j$ .*

270 We first prove a query complexity lower bound for solving  $m$ -INDQ $_n$ .

271 **Lemma 3.3.** *Any  $0.1$ -correct algorithm that solves  $m$ -INDQ $_n$  requires at least  $0.9n$  queries.*

272 Our proof is based on Yao’s minimax principle [Yao, 1977]. See Section A.2 for the full proof.

273 Now we give the construction of our hard instance, together with the high-level intuition of our  
 274 hardness proof. For simplicity, here we assume  $C$  is an integer that divides  $H$ .

275 **The Hard Instance.** Again, our hardness result is based on a binary tree instance. The state space,  
 276 action space, and the transition kernel of our hard instance are exactly the same as the instance in  
 277 Section 3.1. Moreover, similar to the instance in Section 3.1, for each  $h \in [H]$ , there exists an action  
 278  $a_h^* \in \{a_1, a_2\}$ , such that  $\pi^*(s) = a_h^*$  for all  $s \in \mathcal{S}_h$ .

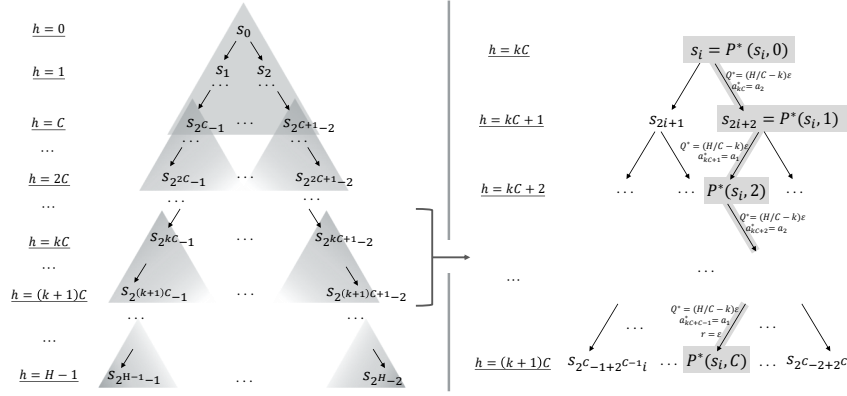


Figure 1: Illustration of the hard instance for Theorem 3.2.

279 To define the reward function  $r$ , we first define an operator  $P^*$ , which can be seen as applying the  
 280 transition kernel for multiple steps by following the optimal policy. For some  $q \in [H/C]$ , a state  $s \in$   
 281  $S_{kC}$ , and an integer  $c \in [C]$ , define  $P^*(s, c) = s$  if  $c = 0$  and  $P^*(s, c) = P(P^*(s, c-1), a_{qC+c-1}^*)$   
 282 otherwise. The reward function  $r(s, a)$  is then defined to be  $\epsilon$  if  $s = P^*(s', C-1)$  for some  
 283  $s' \in S_{qC}$  where  $q \in [H/C]$  and  $a = a_{qC+C-1}^*$ . For all other  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , we define  $r(s, a) = 0$ .  
 284 Accordingly, for each  $(q, c) \in [H/C] \times [C]$ ,  $s \in S_{qC+c}$ , and  $a \in \mathcal{A}$ , we have  $Q^*(s, a) = (H/C - q)\epsilon$   
 285 if  $s = P^*(s', c)$  for some  $s' \in S_{qC}$  and  $a = a_{qC+c}^*$ . For all other  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , we have  
 286  $Q^*(s, a) = (H/C - q - 1)\epsilon$ . This also implies that the value of the optimal policy is  $H/C \cdot \epsilon$ .

287 We define the 1-dimensional feature function  $\phi$  such that, for each  $(q, c) \in [H/C] \times [C]$ ,  $s \in S_{qC+c}$   
 288 and  $a \in \mathcal{A}$ ,  $\phi(s, a) = (H/C - q)\epsilon$ . Clearly, Assumption 1 is satisfied when taking  $\theta^* = 1$ . This  
 289 finishes the construction of our hard instance. An illustration is given in Figure 1.

290 **The Lower Bound.** Now we show that for our hard instance, if there is an RL algorithm that returns  
 291 a policy with suboptimality less than  $H/C \cdot \epsilon$ , then there is an algorithm that solves  $m$ -INDQ $_n$  with  
 292  $n = 2^C$  and  $m = H/C$ . Therefore, the correctness of Theorem 3.2 is implied by Lemma 3.3.

293 We first note that there exists a bijection between  $\{a_1, a_2\}^C$  and  $[2^C]$ . We use  $g : [2^C] \rightarrow \{a_1, a_2\}^T$   
 294 to denote such a bijection. Given an instance of  $m$ -INDQ $_n$  with  $n = 2^C$  and  $m = H/C$ , for  
 295 each  $q \in [H/C]$ , we set  $(a_{qC}^*, a_{qC+1}^*, \dots, a_{(q+1)C-1}^*) = g(i_q^*)$ , where  $(i_0^*, i_1^*, i_2^*, \dots, i_{H/C-1}^*)$  are  
 296 the target indices in the instance of  $m$ -INDQ $_n$ . Each time the RL algorithm samples a trajectory  
 297  $(s_0, a_0, r_0, \dots, s_{H-1}, a_{H-1}, r_{H-1})$ , we make  $H/C$  sequential queries  $(0, i_0), (1, i_1), \dots, (H/C -$   
 298  $1, i_{H/C-1})$  to  $m$ -INDQ $_n$ , where for each  $q \in [H/C]$ ,  $i_q$  is the unique integer in  $[2^C]$  with  $g(i_q) =$   
 299  $(a_{qC}, a_{qC+1}, \dots, a_{(q+1)C-1})$ . For each  $h \in [H]$ , we have  $r_h = \epsilon$  if  $h = (q+1)C - 1$  and  $i_k = i_q^*$   
 300 for some  $k \in [H/C]$ . Otherwise, we have  $r_h = 0$ .

301 Suppose there is an RL algorithm that returns a policy  $\pi$  with suboptimality less than  $H/C \cdot \epsilon$ , and  
 302 since the value of the optimal policy is  $H/C \cdot \epsilon$ , we must have  $r_h = \epsilon$  for some  $h \in [H]$  where  
 303  $(s_0, a_0, r_0, \dots, s_{H-1}, a_{H-1}, r_{H-1})$  is the trajectory obtained by following the policy  $\pi$ . This implies  
 304 the existence of  $q \in [H/C]$  with  $g(i_q^*) = (a_{qC}, a_{qC+1}, \dots, a_{(q+1)C-1})$ . Therefore, if there is an RL  
 305 algorithm that returns a policy with suboptimality less than  $H/C \cdot \epsilon$  for our hard instance, then there  
 306 is an algorithm for solving  $m$ -INDQ $_n$  with  $n = 2^C$  and  $m = H/C$ .

307 **Remark 1.** Our construction is significantly different from Du et al. [2020]. Specifically, we split  
 308 a binary tree with  $H$  levels into  $H/C$  blocks. For each block, we show that any algorithm must  
 309 incur a sample complexity of  $2^C$  in order to find a policy with suboptimality less than  $\epsilon$ . In order to  
 310 show that the overall suboptimality of the RL algorithm is  $H/C \cdot \epsilon$ , we develop a reduction from an  
 311 intermediate problem called MULTI-INDEX-QUERY to RL, which is different from the one used in  
 312 Du et al. [2020].

## 313 4 Main Algorithm

In this section, we present our main algorithm that achieves the guarantee in Theorem 1.3.

---

### Algorithm 1 Elimination Algorithm for Finding the Optimal Hypotheses

---

- 1: **Input:** feature map  $\phi$ , sparsity  $k$ , approximation error  $\epsilon$ , statistical error  $\epsilon_{\text{stat}}$ ,  $\epsilon_{\text{net}}$ , failure rate  $\delta$
- 2: For each  $h \in [H]$ , initialize  $\mathcal{P}_h = \mathcal{P}_h^0 = \{\theta : \theta_{\mathcal{M}} \in \mathcal{N}^k, |\mathcal{M}| = k, \mathcal{M} \subseteq [d]\}$ , where  $\mathcal{N}^k$  is the maximal  $\epsilon_{\text{net}}/2$ -separated subset of the Euclidean sphere  $\mathbb{S}^k$ .
- 3: Calculate  $m = \frac{16k \ln((1+4/\epsilon_{\text{net}})d) + 16 \ln(H/\delta)}{\epsilon_{\text{stat}}}$ .
- 4: **for** iteration  $t = 0, 1, 2, \dots$  **do**
- 5:   Choose  $\theta_0^t = \arg \max_{\theta \in \mathcal{P}_0} V_\theta(s_0)$ .
- 6:   **for**  $h = 1, 2, \dots, H - 1$  **do**
- 7:     Define a policy  $\pi_h^t$ , where  $\pi_h^t(s) = \pi_{\theta_{h'}}^t(s)$  if  $s \in \mathcal{S}_{h'}$  with  $h' < h$ , and arbitrary otherwise. Collect  $m$  trajectories following  $\pi_h^t$  as a dataset

$$\mathcal{D}_h^t = \{(s_0^i, a_0^i, r_0^i, \dots, s_{H-1}^i, a_{H-1}^i, r_{H-1}^i)\}_{i \in [m]}.$$

- 8:     Choose  $\theta_h^t = \arg \max_{\theta \in \mathcal{P}_h} \sum_{i \in [m]} V_\theta(s_h^i)$ , where  $s_h^i$  are from dataset  $\mathcal{D}_h^t$ .
- 9:   **end for**
- 10:   Collect  $m$  trajectories following a policy  $\pi^t = \pi_{\theta^t}$  as a dataset

$$\mathcal{D}_H^t = \{(s_0^i, a_0^i, r_0^i, \dots, s_{H-1}^i, a_{H-1}^i, r_{H-1}^i)\}_{i \in [m]}.$$

- 11:   For each  $h \in [H]$ , calculate using dataset  $\mathcal{D}_H^t$ :

$$\hat{\mathcal{E}}_h^t = \begin{cases} \frac{1}{m} \sum_{i=1}^m (\langle \phi(s_h^i, a_h^i), \theta_h^t \rangle - r_h^i - V_{\theta_{h+1}^t}(s_{h+1}^i)), & \text{if } h \in [H-1] \\ \frac{1}{m} \sum_{i=1}^m (\langle \phi(s_{H-1}^i, a_{H-1}^i), \theta_{H-1}^t \rangle - r_{H-1}^i), & \text{if } h = H-1. \end{cases}$$

- 12:   **if**  $\hat{\mathcal{E}}_h^t \leq 2\epsilon + 2\epsilon_{\text{net}} + 3\epsilon_{\text{stat}}$  for each  $h \in [H-1]$ , and  $\hat{\mathcal{E}}_{H-1}^t \leq \epsilon + \epsilon_{\text{net}} + \epsilon_{\text{stat}}$  **then**
  - 13:     Terminate and output  $\pi_{\theta^t}$ .
  - 14:   **else**
  - 15:     Update  $\mathcal{P}_h = \mathcal{P}_h \setminus \{\theta_h^t\}$ , for all  $h \in [H-1]$  satisfying  $\hat{\mathcal{E}}_h^t \leq 2\epsilon + 2\epsilon_{\text{net}} + 3\epsilon_{\text{stat}}$ , or  $h = H-1$  satisfying  $\hat{\mathcal{E}}_{H-1}^t \leq \epsilon + \epsilon_{\text{net}} + \epsilon_{\text{stat}}$ .
  - 16:   **end if**
  - 17: **end for**
- 

314

315 **Overview.** Here we give an overview of the design of Algorithm 1. We remark that, by Assumption 1,  
 316 each horizon  $h$  has a different optimal  $\theta_h$ . Therefore, a brute-force algorithm would have a sample  
 317 complexity with exponential dependency on  $H$ , while our algorithm has a polynomial dependency  
 318 on  $H$ .

319 First, we approximate all candidate parameter  $\theta$  with a finite set by creating a maximal  $\epsilon_{\text{net}}/2$ -  
 320 separated subset of the euclidean sphere  $\mathbb{S}^{k-1}$ , denoted by  $\mathcal{N}^k$ , and a set of all  $k$ -sized subset of  $[d]$ .  
 321 Then, for each  $h \in [H]$ , we maintain a set of parameter candidates  $\mathcal{P}_h$ . Initially,  $\mathcal{P}_h$  is set to be all  
 322 parameters approximated by  $\mathcal{N}^k$  and  $k$ -sized subset of  $[d]$ , i.e.  $\mathcal{P}_h^0 = \{\theta : \theta_{\mathcal{M}} \in \mathbb{S}^k, |\mathcal{M}| = k, \mathcal{M} \subseteq$   
 323  $[d]\}$  where  $\theta_{\mathcal{M}}$  is the  $k$ -dimension sub-vector of  $\theta$  with indices corresponding to  $\mathcal{M}$ . The set  $\mathcal{P}_h^0$  is  
 324 then finite for all  $h \in [H]$ :  $|\mathcal{P}_h^0| \leq (1 + 4/\epsilon_{\text{net}})^k \cdot \binom{d}{k}$  [Dong and Yang, 2023].

325 During the execution of Algorithm 1, for all  $h \in [H]$ , we eliminate parameter candidates  $\theta$  from  $\mathcal{P}_h$   
 326 if we are certain that  $\theta \neq \theta_h^*$ , where  $\theta^* = (\hat{\theta}_0^*, \hat{\theta}_1^*, \dots, \hat{\theta}_{H-1}^*)$  is a sequence of parameters that is in  
 327  $\mathcal{P}_h^0$  and is closest to the  $\theta^*$  that satisfies Assumption 1, i.e.  $\hat{\theta}_h^* = \arg \min_{\theta \in \mathcal{P}_h^0} \|\theta_h^* - \theta\|$ . Therefore,  
 328 in Algorithm 1, we only consider  $\theta = (\theta_0, \theta_1, \dots, \theta_{H-1})$  if  $\theta_h \in \mathcal{P}_h$  for all  $h \in [H]$ . In the  $t$ -th  
 329 iteration, we choose a parameter  $\theta^t = (\theta_0^t, \theta_1^t, \dots, \theta_{H-1}^t)$  so that  $\theta_h^t$  maximizes  $\mathbb{E}[V_{\theta_h^t}(s_h)]$  and  
 330  $\theta_h^t \in \mathcal{P}_h$  for all  $h \in [H]$ . We then collect  $m$  trajectories to form a dataset  $\mathcal{D}_H^t$  by following the policy



331 induced by  $\theta^t$ . Based on  $\mathcal{D}_H^t$ , we calculate the empirical Bellman error  $\hat{\mathcal{E}}_h^t$  for each  $h \in [H]$ , which  
 332 is the empirical estimate of the average Bellman error defined as follows.

333 **Definition 3** (Average Bellman error). *For a sequence of parameters  $\theta^t = (\theta_0^t, \theta_1^t, \dots, \theta_{H-1}^t)$ , the*  
 334 *average Bellman error of  $\theta^t$  is defined as  $\mathcal{E}_h^t = \mathbb{E}[\langle \phi(s_h, a_h), \theta_h^t \rangle - r(s_h, a_h) - V_{f_{h+1}^t}(s_{h+1})]$  when*  
 335  *$h \in [H - 1]$  and  $\mathcal{E}_{H-1}^t = \mathbb{E}[\langle \phi(s_{H-1}, a_{H-1}), \theta_{H-1}^t \rangle - r(s_{H-1}, a_{H-1})]$  for level  $H - 1$ . Here,*  
 336  *$(s_0, a_0, r_0, \dots, s_{H-1}, a_{H-1}, r_{H-1})$  is a trajectory obtained by following  $\pi_{\theta^t}$ .*

337 Intuitively, the Bellman error at level  $h$  measures the consistency of  $\theta_h^t$  and  $\theta_{h+1}^t$  for the state-action  
 338 distribution induced by  $\pi_{\theta^t}$ . In each iteration of Algorithm 1, we check if  $\hat{\mathcal{E}}_h^t$  is small for all  $h \in [H]$ .  
 339 If so, the algorithm terminates and returns the policy  $\pi_{\theta^t}$ . Otherwise, for all levels  $h \in [H]$  where  $\hat{\mathcal{E}}_h^t$   
 340 is large, we eliminate  $\theta_h^t$  from  $\mathcal{P}_h$  and proceed to the next iteration.

341 Now we give the analysis of Algorithm 1.

342 **Sample Complexity.** To bound the sample complexity of Algorithm 1, it suffices to give an upper  
 343 bound on the number of iterations, since in each iteration, the number of trajectories sampled by the  
 344 algorithm is simply  $H^2 \cdot m = 16H^2(k \ln((1 + 4/\epsilon_{\text{net}})d) + \ln(H/\delta))/(\epsilon_{\text{stat}}^2)$ . The following lemma  
 345 gives an upper bound on the number of iterations of Algorithm 1. The proof is given by counting the  
 346 number of parameters in the parameter space. The detailed proof is given in Appendix B.1.

347 **Lemma 4.1.** *For any MDP instance with horizon  $H$  and satisfying Assumption 1 with sparsity  $k$ ,*  
 348 *Algorithm 1 runs for at most  $(1 + 4/\epsilon_{\text{net}})^k \binom{d}{k} H$  iterations.*

349 **Remark 2.** *Previous works [Weisz et al., 2022, Wang et al., 2021b] show that even when the*  
 350 *optimal Q-function is well-specified, any RL algorithm would require a sample size with exponential*  
 351 *dependency on  $d$  or  $H$ . Note that this is equivalent to the case where the sparsity  $k = d$  and the*  
 352 *approximation error  $\epsilon = 0$  in our setting. Therefore, in our misspecified setting, which is strictly*  
 353 *harder, exponential dependency on  $k$  is unavoidable, unless we can accept an exponential dependency*  
 354 *on  $H$ .*

355 **Suboptimality of the Returned Policy.** We now show that with probability at least  $1 - \delta$ , the  
 356 suboptimality of the returned policy is at most  $(2\epsilon + 2\epsilon_{\text{net}} + 4\epsilon_{\text{stat}})H$ . First, we define a high  
 357 probability event  $E$ , which we will condition on in the remaining part of the analysis.

358 **Definition 4.** *Define  $E$  as the event that  $|\mathcal{E}_h^t - \hat{\mathcal{E}}_h^t| \leq \epsilon_{\text{stat}}$  and  $|\mathbb{E}_{s_h \sim \pi_h^t} V_\theta(s_h) - \sum_{i \in [m]} V_\theta(s_h^i)| \leq$   
 359  $\epsilon_{\text{stat}}$  (where  $s_h^i$  is from  $\mathcal{D}_h^t$ ) for all iterations  $t$ , horizon  $h \in [H]$ , and parameter  $\theta \in \mathcal{P}_h^0$ .*

360 **Lemma 4.2.** *Event  $E$  holds with probability at least  $1 - \delta$ .*

361 To prove Lemma 4.2, we first consider a fixed level  $h$  and iteration  $t$ . Since the empirical Bellman  
 362 error  $\hat{\mathcal{E}}_h^t$  is simply the empirical estimate of  $\mathcal{E}_h^t$ , and  $\sum_{i \in [m]} V_\theta(s_h^i)$  is simply an empirical estimate  
 363 of  $\mathbb{E}_{s_h \sim \pi_h^t} V_\theta(s_h)$ , applying the Chernoff-Hoeffding inequality respectively would suffice. Moreover,  
 364 the number of iterations has an upper bound given by Lemma 4.1. Therefore, Lemma 4.2 follows by  
 365 applying a union bound over all  $h \in [H]$ ,  $t \in [(1 + 4/\epsilon_{\text{net}})^k \binom{d}{k} H]$  and parameter  $\theta \in \mathcal{P}_h^0$ .

366 We next show that, conditioned on event  $E$  defined above, for the sequence of parameters  $\theta^* =$   
 367  $(\theta_0^*, \theta_1^*, \dots, \theta_{H-1}^*)$  that satisfies Assumption 1, we never eliminate  $\hat{\theta}_h^*$  from  $\mathcal{P}_h$ , for all  $h \in [H]$ .

368 **Lemma 4.3.** *Conditioned on event  $E$  defined in Definition 4, for a sequence of parameters*  
 369  *$(\theta_0^*, \theta_1^*, \dots, \theta_{H-1}^*)$  that satisfies Assumption 1, and their approximations  $\hat{\theta}_h^* = \arg \min_{\theta \in \mathcal{P}_h^0} \|\theta_h^* - \theta\|$ ,*  
 370 *during the execution of Algorithm 1,  $\hat{\theta}_h^*$  is never eliminated from  $\mathcal{P}_h$  for all  $h \in [H]$ .*

371 To prove Lemma 4.3, the main observation is that, for  $h \in [H - 1]$  the average Bellman error  
 372 induced by  $\hat{\theta}_h^*$  and  $\theta_{h+1}^t = \arg \max_{\theta \in \mathcal{P}_{h+1}} \mathbb{E}_{s_{h+1}} [V_\theta(s_{h+1})]$  is always upper bounded by  $2(\epsilon + \epsilon_{\text{net}})$ ,  
 373 regardless of the distribution of  $(s_h, a_h)$  (cf. Definition 3). Conditioned on event  $E$ , the empirical  
 374 Bellman error induced by  $\hat{\theta}_h^*$  and  $\theta_{h+1}^t$  is at most  $2\epsilon + 2\epsilon_{\text{net}} + 3\epsilon_{\text{stat}}$ . Similarly, the empirical Bellman  
 375 error induced by  $\hat{\theta}_{H-1}^*$  is at most  $\epsilon + \epsilon_{\text{net}} + \epsilon_{\text{stat}}$ . In Algorithm 1, we eliminate function  $\theta_h^t$  only  
 376 when the empirical Bellman error is larger than these (Line 15). Thus,  $\hat{\theta}_h^*$  is never eliminated.

377 We now show the suboptimality of the policy returned by Algorithm 1 is at most  $(2\epsilon + 2\epsilon_{\text{net}} + 4\epsilon_{\text{stat}})H$ .

378 **Lemma 4.4.** *For any MDP instance satisfying Assumption 1, conditioned on event  $E$  defined in*  
 379 *Definition 4, Algorithm 1 returns a policy  $\pi$  satisfying  $V^* - V^\pi \leq (2\epsilon + 2\epsilon_{\text{net}} + 4\epsilon_{\text{stat}})H$ .*

380 To prove Lemma 4.4, we first recall the policy loss decomposition lemma (Lemma 1 in Jiang et al.  
 381 [2017]), which states that for a policy induced by a sequence of parameters  $\theta = (\theta_0, \theta_1, \dots, \theta_{H-1})$ ,  
 382  $V_{\theta_0}(s_0) - V^{\pi_\theta}$  is upper bounded by the summation of average Bellman error over all levels  $h \in [H]$ .  
 383 When Algorithm 1 terminates, the empirical Bellman error must be small for all  $h \in [H]$ , and  
 384 therefore, the average Bellman error is small by definition of the event  $E$ . Moreover, in Line 5 of  
 385 Algorithm 1, we always choose a parameter  $\theta$  that maximizes  $V_\theta(s_0)$ . Since the sequence of functions  
 386  $\hat{\theta}^* = (\hat{\theta}_0^*, \hat{\theta}_1^*, \dots, \hat{\theta}_{H-1}^*)$  is never eliminated by Lemma 4.3, we must have  $V_{\theta_0}(s_0) \geq V_{\hat{\theta}_0^*}(s_0) \geq$   
 387  $V^* - \epsilon - \epsilon_{\text{net}}$ , which gives an upper bound on the suboptimality of the policy returned by Algorithm 1.  
 388 Combining Lemma 4.1, Lemma 4.2 and Lemma 4.4, we can prove Theorem 1.3.

389 **Remark 3.** *While Algorithm 1 assumes the sparsity constant  $k$  is known, it can be easily adapted*  
 390 *to the setting where  $k$  is not known beforehand. For such a setting, we could enumerate  $k$  starting*  
 391 *from  $k = 1$ , and use the following observations: 1) If the true  $k$ , say  $k^*$ , is larger than  $k$ , then*  
 392 *running our algorithm with sparsity  $k$  will eliminate all the parameters in the parameter space  $\mathcal{P}_h$*   
 393 *for some horizon  $h$ . 2) If for all  $h$ , there exists one parameter in  $\mathcal{P}_h$  that has not been deleted, the*  
 394 *we have identified  $k^*$ . The sample complexity of this process is asymptotically the same as running*  
 395 *Algorithm 1 with known  $k$ , since the true  $k^*$  dominates the sample complexity.*

396 **Implications.** We can think of the bandit setting as an MDP with  $H = 1$  and derive the following.

397 **Corollary 4.5.** *For the bandit setting satisfying Assumption 1, Algorithm 1 returns an action  $\hat{a}$  such*  
 398 *that  $r(a^*) - r(\hat{a}) \leq 2\epsilon + 2\epsilon_{\text{net}} + 4\epsilon_{\text{stat}}$ .*

399 **Remark 4.** *Here we compare Corollary 4.5 with the result in Dong and Yang [2023]. Scrutinizing*  
 400 *the analysis in Dong and Yang [2023], the suboptimality achieved by their algorithm is  $4\epsilon + \epsilon_{\text{stat}}$ ,*  
 401 *which is worse than our suboptimality guarantee. On the other hand, the algorithm in Dong and Yang*  
 402 *[2023] also returns a parameter  $\theta$  such that  $|\langle \phi(a), \theta \rangle - r(a)| \leq 2\epsilon + \epsilon_{\text{stat}}$  for all  $a \in \mathcal{A}$  (which is*  
 403 *the best possible according to Theorem 1.1), where our algorithm only returns a near-optimal action.*

## 404 5 Conclusion

405 We studied RL problem where the optimal  $Q$ -functions can be approximated by linear function with  
 406 constant sparsity  $k$ , up to an error of  $\epsilon$ . We design a new algorithm with polynomial sample complexity,  
 407 while the suboptimality of the returned policy is  $O(H\epsilon)$ , which is shown to be near-optimal by a  
 408 information-theoretic hardness result.

409 Although the suboptimality guarantee achieved by our algorithm is near-optimal, the sample com-  
 410 plexity can be further improved. As an interesting future direction, it would be interesting to design  
 411 an RL algorithm with the same suboptimality guarantee, while obtaining tighter dependence on the  
 412 horizon length  $H$ .

## 413 References

- 414 Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. Flambe: Structural complexity  
 415 and representation learning of low rank MDPs. *Advances in neural information processing systems*,  
 416 33:20095–20107, 2020.
- 417 Alekh Agarwal, Yujia Jin, and Tong Zhang. VOQL: Towards optimal regret in model-free rl with  
 418 nonlinear function approximation. In *The Thirty Sixth Annual Conference on Learning Theory*,  
 419 pages 987–1063. PMLR, 2023.
- 420 Philip Amortila, Tongyi Cao, and Akshay Krishnamurthy. Mitigating covariate shift in misspecified  
 421 regression with applications to reinforcement learning, 2024.
- 422 Alex Ayoub, Zeyu Jia, Csaba Szepesvari, Mengdi Wang, and Lin Yang. Model-based reinforcement  
 423 learning with value-targeted regression. In *International Conference on Machine Learning*, pages  
 424 463–474. PMLR, 2020.

- 425 Djallel Bouneffouf, Amel Bouzeghoub, and Alda Lopes Gançarski. A contextual-bandit algorithm  
426 for mobile context-aware recommender system. In *International conference on neural information*  
427 *processing*, pages 324–331. Springer, 2012.
- 428 Qi Cai, Zhuoran Yang, Chi Jin, and Zhaoran Wang. Provably efficient exploration in policy optimiza-  
429 tion. In *International Conference on Machine Learning*, pages 1283–1294. PMLR, 2020.
- 430 Fan Chen, Song Mei, and Yu Bai. Unified algorithms for rl with decision-estimation coefficients:  
431 No-regret, pac, and reward-free learning. *arXiv preprint arXiv:2209.11745*, 2022a.
- 432 Zixiang Chen, Chris Junchi Li, Angela Yuan, Quanquan Gu, and Michael I Jordan. A general  
433 framework for sample-efficient function approximation in reinforcement learning. *arXiv preprint*  
434 *arXiv:2209.15634*, 2022b.
- 435 Christoph Dann, Mehryar Mohri, Tong Zhang, and Julian Zimmert. A provably efficient model-free  
436 posterior sampling method for episodic reinforcement learning. *Advances in Neural Information*  
437 *Processing Systems*, 34:12040–12051, 2021.
- 438 Jialin Dong and Lin Yang. Does sparsity help in learning misspecified linear bandits? In *International*  
439 *Conference on Machine Learning*, pages 8317–8333. PMLR, 2023.
- 440 Simon Du, Sham Kakade, Jason Lee, Shachar Lovett, Gaurav Mahajan, Wen Sun, and Ruosong  
441 Wang. Bilinear classes: A structural framework for provable generalization in rl. In *International*  
442 *Conference on Machine Learning*, pages 2826–2836. PMLR, 2021.
- 443 Simon S Du, Sham M Kakade, Ruosong Wang, and Lin F Yang. Is a good representation sufficient for  
444 sample efficient reinforcement learning? In *International Conference on Learning Representations*,  
445 2020.
- 446 Andre Esteva, Alexandre Robicquet, Bharath Ramsundar, Volodymyr Kuleshov, Mark DePristo,  
447 Katherine Chou, Claire Cui, Greg Corrado, Sebastian Thrun, and Jeff Dean. A guide to deep  
448 learning in healthcare. *Nature medicine*, 25(1):24–29, 2019.
- 449 Dylan J Foster, Sham M Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity of  
450 interactive decision making. *arXiv preprint arXiv:2112.13487*, 2021.
- 451 Dylan J Foster, Noah Golowich, and Yanjun Han. Tight guarantees for interactive decision making  
452 with the decision-estimation coefficient. *arXiv preprint arXiv:2301.08215*, 2023.
- 453 Zeyu Jia, Lin Yang, Csaba Szepesvari, and Mengdi Wang. Model-based reinforcement learning with  
454 value-targeted regression. In *Learning for Dynamics and Control*, pages 666–686. PMLR, 2020.
- 455 Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Context-  
456 tual decision processes with low bellman rank are pac-learnable. In *International Conference on*  
457 *Machine Learning*, pages 1704–1713. PMLR, 2017.
- 458 Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement  
459 learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143.  
460 PMLR, 2020.
- 461 Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman eluder dimension: New rich classes of rl  
462 problems, and sample-efficient algorithms. *Advances in neural information processing systems*,  
463 34:13406–13418, 2021.
- 464 B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani,  
465 and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE*  
466 *Transactions on Intelligent Transportation Systems*, 2021.
- 467 Dingwen Kong, Ruslan Salakhutdinov, Ruosong Wang, and Lin F Yang. Online sub-sampling for  
468 reinforcement learning with general function approximation. *arXiv preprint arXiv:2106.07203*,  
469 2021.
- 470 Tor Lattimore, Csaba Szepesvari, and Gellert Weisz. Learning with good feature representations in  
471 bandits and in RL with a generative model. In *International Conference on Machine Learning*,  
472 pages 5662–5670. PMLR, 2020.

- 473 Zhihan Liu, Miao Lu, Wei Xiong, Han Zhong, Hao Hu, Shenao Zhang, Sirui Zheng, Zhuoran Yang,  
474 and Zhaoran Wang. One objective to rule them all: A maximization objective fusing estimation  
475 and planning for exploration. *arXiv preprint arXiv:2305.18258*, 2023.
- 476 Aditya Modi, Nan Jiang, Ambuj Tewari, and Satinder Singh. Sample complexity of reinforcement  
477 learning using linearly combined model ensembles. In *International Conference on Artificial*  
478 *Intelligence and Statistics*, pages 2010–2020. PMLR, 2020.
- 479 Gergely Neu and Ciara Pike-Burke. A unifying view of optimism in episodic reinforcement learning.  
480 *Advances in Neural Information Processing Systems*, 33:1392–1403, 2020.
- 481 Ian Osband and Benjamin Van Roy. Model-based reinforcement learning and the eluder dimension.  
482 *Advances in Neural Information Processing Systems*, 27, 2014.
- 483 Eric M Schwartz, Eric T Bradlow, and Peter S Fader. Customer acquisition via display advertising  
484 using multi-armed bandit experiments. *Marketing Science*, 36(4):500–522, 2017.
- 485 Wen Sun, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Model-based  
486 rl in contextual decision processes: Pac bounds and exponential improvements over model-free  
487 approaches. In *Conference on learning theory*, pages 2898–2933. PMLR, 2019.
- 488 Andrew J Wagenmaker and Dylan J Foster. Instance-optimality in interactive decision making:  
489 Toward a non-asymptotic theory. In *The Thirty Sixth Annual Conference on Learning Theory*,  
490 pages 1322–1472. PMLR, 2023.
- 491 Ruosong Wang, Dean P Foster, and Sham M Kakade. What are the statistical limits of offline rl with  
492 linear function approximation? *arXiv preprint arXiv:2010.11895*, 2020a.
- 493 Ruosong Wang, Russ R Salakhutdinov, and Lin Yang. Reinforcement learning with general value  
494 function approximation: Provably efficient approach via bounded eluder dimension. *Advances in*  
495 *Neural Information Processing Systems*, 33:6123–6135, 2020b.
- 496 Ruosong Wang, Yifan Wu, Ruslan Salakhutdinov, and Sham Kakade. Instabilities of offline rl  
497 with pre-trained neural representation. In *International Conference on Machine Learning*, pages  
498 10948–10960. PMLR, 2021a.
- 499 Yining Wang, Ruosong Wang, Simon S Du, and Akshay Krishnamurthy. Optimism in reinforcement  
500 learning with generalized linear function approximation. *arXiv preprint arXiv:1912.04136*, 2019.
- 501 Yuanhao Wang, Ruosong Wang, and Sham M. Kakade. An exponential lower bound for linearly-  
502 realizable mdps with constant suboptimality gap, 2021b. URL [https://arxiv.org/abs/2103.](https://arxiv.org/abs/2103.12690)  
503 12690.
- 504 Gellért Weisz, Philip Amortila, and Csaba Szepesvári. Exponential lower bounds for planning in  
505 mdps with linearly-realizable optimal action-value functions. In *Algorithmic Learning Theory*,  
506 pages 1237–1264. PMLR, 2021.
- 507 Gellért Weisz, Csaba Szepesvári, and András György. Tensorplan and the few actions lower bound  
508 for planning in mdps under linear realizability of optimal value functions, 2022. URL [https:](https://arxiv.org/abs/2110.02195)  
509 [//arxiv.org/abs/2110.02195](https://arxiv.org/abs/2110.02195).
- 510 Andrew Chi-Chin Yao. Probabilistic computations: Toward a unified measure of complexity. In  
511 *18th Annual Symposium on Foundations of Computer Science (sfcs 1977)*, pages 222–227. IEEE  
512 Computer Society, 1977.
- 513 Andrea Zanette, Alessandro Lazaric, Mykel J Kochenderfer, and Emma Brunskill. Limiting extrapo-  
514 lation in linear approximate value iteration. *Advances in Neural Information Processing Systems*,  
515 32, 2019.
- 516 Andrea Zanette, Alessandro Lazaric, Mykel Kochenderfer, and Emma Brunskill. Learning near  
517 optimal policies with low inherent bellman error. In *International Conference on Machine Learning*,  
518 pages 10978–10989. PMLR, 2020a.

- 519 Andrea Zanette, Alessandro Lazaric, Mykel J Kochenderfer, and Emma Brunskill. Provably efficient  
520 reward-agnostic navigation with linear value iteration. *Advances in Neural Information Processing*  
521 *Systems*, 33:11756–11766, 2020b.
- 522 Han Zhong, Wei Xiong, Sirui Zheng, Liwei Wang, Zhaoran Wang, Zhuoran Yang, and Tong Zhang.  
523 Gec: A unified framework for interactive decision making in mdp, pomdp, and beyond. *arXiv*  
524 *preprint arXiv:2211.01962*, 2022.
- 525 Dongruo Zhou and Quanquan Gu. Computationally efficient horizon-free reinforcement learning for  
526 linear mixture mdps. *Advances in neural information processing systems*, 35:36337–36349, 2022.
- 527 Dongruo Zhou, Jiafan He, and Quanquan Gu. Provably efficient reinforcement learning for discounted  
528 mdps with feature mapping. In *International Conference on Machine Learning*, pages 12793–  
529 12802. PMLR, 2021.

530 **A Proofs in Section 3**

531 **A.1 Proof of Theorem 3.1**

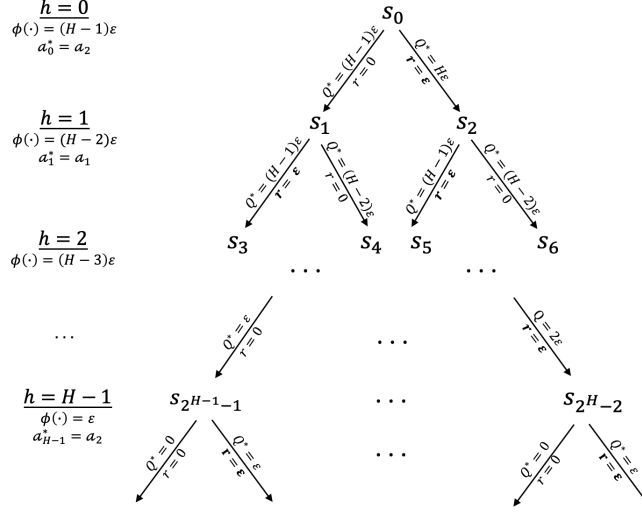


Figure 2: Illustration of the hard instance for Theorem 3.1.

532 *Proof.* Consider an input distribution where  $a_h^*$  is drawn uniformly random from  $\{a_1, a_2\}$ . By Yao's  
533 minimax principle, it suffices to consider the best deterministic algorithm, say  $A$ . Note that, since we  
534 have no sampling ability, a deterministic algorithm in this setting can be seen as a function that takes  
535 in feature function  $\phi$  and returns a policy  $\pi$ . Also, for all instances supported by this distribution, their  
536 inputs  $\phi$  are the same. Thus, the policy returned by  $A$  is fixed. Denote the policy as  $\pi$ , and denote the  
537 trajectory following  $\pi$  as  $(s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{H-1}, a_{H-1}, r_{H-1})$ . The suboptimality of  $\pi$  can  
538 be written as

$$V^* - V^\pi = \sum_{h=0}^{H-1} \epsilon \cdot \mathbb{I}[a_h^* \neq a_h]$$

539 Since  $a_h$  is fixed and  $a_h^*$  is drawn uniformly random from  $\{a_1, a_2\}$ ,  $\mathbb{I}[a_h^* \neq a_h] = 1$  with probability  
540  $1/2$ . Thus,  $(V^* - V^\pi)/\epsilon$  is a binomial random variable, or  $(V^* - V^\pi)/\epsilon \sim B(H, 1/2)$ . The  
541 expectation of  $(V^* - V^\pi)$  is then  $H\epsilon/2$ , and its variance is  $H\epsilon^2/4$ . Using Chebyshev inequality,  
542 with probability 0.99, we have

$$V^* - V^\pi \geq \frac{1}{2}H\epsilon - 5\epsilon\sqrt{H} = \Omega(H\epsilon)$$

543 for sufficiently large  $H \geq 100$ . □

544 **A.2 Proof of Lemma 3.3**

545 *Proof.* Consider an input distribution where  $i^* = (i_0^*, i_1^*, \dots, i_{m-1}^*)$  is drawn uniformly random  
546 from  $[n]^m$ . Let  $c(i^*, a)$  be the query complexity of running algorithm  $a$  to solve the problem with  
547 correct indices  $i^*$ . Assume there exists a 0.1-correct algorithm  $\mathcal{A}$  for  $m$ -INDQ $_n$  that queries less  
548 than  $0.9n$  times in the worst case. Then, using Yao's minimax principle, there exists a deterministic  
549 algorithm  $\mathcal{A}'$  with  $c(i^*, \mathcal{A}') < 0.9n$  for all  $i^* \in [n]^m$ , such that

$$\mathbb{P}[\mathcal{A}' \text{ outputs } (j, i_j^*) \text{ for some } j \in [m]] \geq 0.9.$$

550 We may assume that the sequence of queries made by  $\mathcal{A}'$  is fixed until it correctly guesses one of  $i_j^*$ .  
551 This is because  $\mathcal{A}'$  is deterministic, and the responses  $\mathcal{A}'$  receives are the same (i.e. all guesses are  
552 incorrect) until it correctly queries  $(j, i_j^*)$  for some  $j$ . Let  $S = \{s_1, \dots, s_k\}$  be the sequence of first  
553  $k$  guesses made by  $\mathcal{A}'$ , and let  $I_{BAD} \subset [n]^m$  be a set of all possible  $i^*$ 's such that the guesses in  $S$   
554 are all incorrect. Denote the number of guesses on INDQ $_n^{(j)}$  in  $S$  by  $n_j$ , then  $n_j$ 's are also fixed, and  
555  $\sum_{j \in [m]} n_j = k$ . The size of  $I_{BAD}$  then satisfies

$$|I_{BAD}| = \prod_{j=0}^{m-1} (n - n_j) \geq (n - k)n^{m-1}$$

556 Set  $k$  as the worst-case query complexity of  $\mathcal{A}'$ . Then, for all  $i^* \in I_{BAD}$ , the output of  $\mathcal{A}'$  is incorrect.  
 557 Since  $i^*$  is drawn uniformly random from  $[n]^m$ , the probability of  $\mathcal{A}'$  being incorrect is

$$\mathbb{P}[\mathcal{A}' \text{ is incorrect}] = \frac{|I_{BAD}|}{|[n]^m|} \geq \frac{(n-k)n^{m-1}}{n^m} > \frac{(n-0.9n)n^{m-1}}{n^m} > 0.1,$$

558 where in the second to last inequality we used  $k < 0.9n$ .

559 However, this contradicts with the fact that  $\mathbb{P}[\mathcal{A}' \text{ outputs } (j, i_j^*) \text{ for some } j \in [m]] \geq 0.9$ . Thus, there  
 560 does not exist a 0.1-correct algorithm that solves the problem with less than  $0.9n$  queries in the worst  
 561 case.  $\square$

### 562 A.3 Proof of Theorem 3.2

563 *Proof.* First, we prove our claim based on the assumption that  $C$  is an integer that divides  $H$ . We can  
 564 create the hard instance described in Section 3.2.

565 We reduce the problem to  $H/C$ -INDQ $_{2^C}$ . Assume there exists an algorithm  $\mathcal{A}$  that takes less than  $0.9 \cdot$   
 566  $2^C \cdot C$  samples, such that, with probability at least 0.9, it outputs a policy  $\pi$  with suboptimality  $V^* -$   
 567  $V^\pi < H/C \cdot \epsilon$ . By definition, at round  $i$ ,  $\mathcal{A}$  interacts with the MDP instance by following a trajectory  
 568  $(s_0, a_0^i, r_0^i, \dots, s_{H-1}^i, a_{H-1}^i, r_{H-1}^i)$ . Based on  $\mathcal{A}$ , we create an algorithm  $\mathcal{A}'$  for  $H/C$ -INDQ $_{2^C}$  as  
 569 follows. Consider  $\mathcal{A}$  is querying the trajectory  $(s_0, a_0^i, r_0^i, \dots, s_{H-1}^i, a_{H-1}^i, r_{H-1}^i)$ . For each  $q \in$   
 570  $\{0, \dots, H/C - 1\}$ , we can map  $(a_{qC}^i, \dots, a_{(q+1)C-1}^i)$  to an index in  $[2^C]$  using the bijection  $g$ . Thus,  
 571 we make a sequence of  $H/C$  guesses,  $\{(q, g(a_{qC}^i, \dots, a_{(q+1)C-1}^i))\}_{q=0}^{H/C-1}$ , to the  $H/C$ -INDQ $_{2^C}$ .  
 572 If the guess  $(q, g(a_{qC}^i, \dots, a_{(q+1)C-1}^i))$  is correct for some  $q$ ,  $\mathcal{A}$  receives a reward of  $\epsilon$  at level  
 573  $(q+1)C - 1$ , i.e.  $r_{(q+1)C-1}^i = r(s_{(q+1)C-1}^i, a_{(q+1)C-1}^i) = \epsilon$ . For all other state-action pairs in the  
 574 trajectory, algorithm  $\mathcal{A}$  receives zero reward. Since  $\mathcal{A}$  takes less than  $0.9 \cdot 2^C \cdot C$  samples, it queries  
 575 less than  $0.9 \cdot 2^C \cdot C/H$  trajectories, corresponding  $0.9 \cdot 2^C$  guesses to  $H/C$ -INDQ $_{2^C}$  in total. Recall  
 576 that  $\mathcal{A}$  outputs a policy  $\pi$  with suboptimality  $V^* - V^\pi < H/C \cdot \epsilon$  with probability at least 0.9. This  
 577 means the sequence of guesses to  $H/C$ -INDQ $_{2^C}$  made by  $\pi$  must have at least one of them being  
 578 correct. Thus,  $\mathcal{A}'$  is a 0.1-correct algorithm that solves  $H/C$ -INDQ $_{2^C}$  with less than  $0.9 \cdot 2^C$  guesses.  
 579 However, by Lemma 3.3, such an algorithm does not exist, so  $\mathcal{A}$  does not exist. We conclude that  
 580 any algorithm that returns a policy with suboptimality less than  $H/C \cdot \epsilon$  with probability at least 0.9  
 581 needs to sample at least  $0.9 \cdot C \cdot 2^C$  times.

582 Now we consider when  $C$  is not an integer that divides  $H$ . There are two cases. First, consider  $C$  as  
 583 an integer that does not divide  $H$ . Let  $H' = \lfloor H/C \rfloor \cdot C$ , then we can make the same construction  
 584 as above for the first  $H'$  horizons, and set the reward as zero for all the state-action pairs in the  
 585 remaining  $H - H'$  levels. Because the rewards are the same for levels  $H'$  through  $H - 1$ , different  
 586 values of  $\{\pi_{H'}, \dots, \pi_{H-1}\}$  do not make a difference to  $V^\pi$ . Therefore, we only care about the  
 587 first  $H'$  levels, so we can conclude from our above analysis that, any algorithm that returns a  
 588 policy with suboptimality less than  $H'/C \cdot \epsilon = \lfloor H/C \rfloor \cdot \epsilon$  with probability at least 0.9 needs to  
 589 sample at least  $0.9 \cdot C \cdot 2^C$  times. For the second case, we consider when  $C$  is not an integer. Let  
 590  $C' = \lfloor C \rfloor$ , we can apply our conclusion from the previous case. That is, any algorithm that returns a  
 591 policy with suboptimality less than  $\lfloor H/C' \rfloor \cdot \epsilon$  with probability at least 0.9 needs to sample at least  
 592  $0.9 \cdot C' \cdot 2^{C'}$  times. Since  $2C \leq H$ , we have  $\lfloor H/C' \rfloor \cdot \epsilon \geq \lfloor H/C \rfloor \cdot \epsilon \geq \frac{H\epsilon}{2C}$ . Also observing that  
 593  $0.9 \cdot C' \cdot 2^{C'} \geq 0.1 \cdot C \cdot 2^C$ , we finish the proof.  $\square$

## 594 B Proofs in Section 4

### 595 B.1 Proof of Lemma 4.1

596 *Proof.* In each iteration, we either output a policy or delete at least one function in  $\mathcal{P}_h$  for some  
 597  $h \in [H - 1]$ . Since there are  $\sum_{h \in [H-1]} (|\mathcal{S}^k| \times \binom{d}{k}) \leq (1 + 4/\epsilon_{\text{net}})^k \binom{d}{k} H$  functions in total  
 598 initially, the algorithm is guaranteed to terminate within  $(1 + 4/\epsilon_{\text{net}})^k \binom{d}{k} H$  iterations.  $\square$

599 **B.2 Proof of Lemma 4.2**

**Lemma B.1** (Deviation bound for  $\mathcal{E}_h$ ). *For fixed iteration  $t$  and horizon  $h \in [H]$ , with probability at least  $1 - \delta'$ , we have*

$$|\mathcal{E}_h^t - \hat{\mathcal{E}}_h^t| \leq 4\sqrt{\frac{\ln 2 - \ln \delta'}{2m}}.$$

600 *Hence, we can set  $m > \frac{16(\ln 2 - \ln(\delta'))}{2\epsilon_{\text{stat}}^2}$  to guarantee that  $|\mathcal{E}_h^t - \hat{\mathcal{E}}_h^t| < \epsilon_{\text{stat}}$*

601 *Proof.* Recall that the batch dataset  $\mathcal{D}_t = \{(a_0^i, r_0^i, s_1^i, \dots, a_{H-1}^i, r_{H-1}^i)\}_{i=1}^m$  is collected by playing  
602 policy  $\pi_{\theta^t}$ . We define  $\hat{\mathcal{E}}_h^{t,i} = \langle \phi(s_{h-1}^i, a_h^i), \theta_h^t \rangle - r(s_h^i, a_h^i) - V_{\theta_{h+1}^t}(s_{h+1}^i)$ , then  $\hat{\mathcal{E}}_h^t = \frac{1}{m} \sum_{i=1}^m \hat{\mathcal{E}}_h^{t,i}$ .  
603 By definition of  $\mathcal{E}_h^t$ , it satisfies

$$\mathcal{E}_h^t = \mathbb{E}[\hat{\mathcal{E}}_h^{t,i}].$$

604 Further, since  $\langle \phi(s, a), \theta^t \rangle \in [-1, 1]$  and  $r(s, a) \in [0, 1]$  for any state-action pair  $(s, a)$ , we have  
605  $\hat{\mathcal{E}}_h^{t,i} \in [-3, 1]$ . Thus, using Chernoff-Hoeffding inequality, we get, with probability  $1 - \delta'$ ,

$$|\mathcal{E}_h^t - \hat{\mathcal{E}}_h^t| = \left| \frac{1}{m} \sum_{i=1}^m \left( \hat{\mathcal{E}}_h^{t,i} - \mathbb{E}[\hat{\mathcal{E}}_h^t] \right) \right| \leq 4\sqrt{\frac{\ln 2 - \ln \delta'}{2m}}.$$

606 □

607 **Lemma B.2** (Deviation bound for  $\mathbb{E}_{s_h \sim \pi_h^t} V_\theta(s_h)$ ). *For fixed iteration  $t$ , horizon  $h \in [H]$ , and  
608 parameter  $\theta \in \mathcal{P}_h$ , with probability at least  $1 - \delta'$ , we have*

$$\left| \mathbb{E}_{s_h \sim \pi_h^t} V_\theta(s_h) - \frac{1}{m} \sum_{i \in [m]} V_\theta(s_h^i) \right| \leq \sqrt{\frac{\ln 2 - \ln \delta'}{2m}}.$$

609 *Hence, we can set  $m > \frac{\ln 2 - \ln(\delta')}{2\epsilon_{\text{stat}}^2}$  to guarantee that  $|\mathcal{E}_h^t - \hat{\mathcal{E}}_h^t| < \epsilon_{\text{stat}}$*

610 *Proof.* Recall that the batch dataset  $\mathcal{D}_h^t = \{(s_0^i, a_0^i, r_0^i, \dots, s_{h-1}^i, a_{h-1}^i, r_{h-1}^i, s_h^i)\}_{i=1}^m$  is collected  
611 by playing policy  $\pi_h^t$ . By definition, it satisfies

$$\mathbb{E}_{s_h \sim \pi_h^t} V_\theta(s_h) = \mathbb{E}_{\mathcal{D}_h^t} \left[ \frac{1}{m} \sum_{i \in [m]} V_\theta(s_h^i) \right].$$

612 Further, since  $V_\theta(s) \in [0, 1]$  for any state  $s$ , using Chernoff-Hoeffding inequality, we have with  
613 probability  $1 - \delta'$ ,

$$\left| \mathbb{E}_{s_h \sim \pi_h^t} V_\theta(s_h) - \frac{1}{m} \sum_{i \in [m]} V_\theta(s_h^i) \right| \leq \sqrt{\frac{\ln 2 - \ln \delta'}{2m}}.$$

614 □

615 *Proof.* Define  $E_{t,h}^{\mathcal{E}}$  to be the event

$$E_{t,h}^{\mathcal{E}} = \{|\mathcal{E}_h^t - \hat{\mathcal{E}}_h^t| \leq \epsilon_{\text{stat}}\},$$

616 then by Lemma B.1,  $\mathbb{P}(E_{t,h}^{\mathcal{E}}) \geq 1 - \delta/(2(1 + 4/\epsilon_{\text{net}})^{2k} \binom{d}{k}^2 H^2)$  for all iterations  $t \in [(1 +$   
617  $4/\epsilon_{\text{net}})^k \binom{d}{k} H]$  and horizon  $h \in [H]$ .

618 Define  $E_{t,h,\theta}^V$  to be the event

$$E_{t,h,\theta}^V = \left\{ \left| \mathbb{E}_{s_h \sim \pi_h^t} V_\theta(s_h) - \frac{1}{m} \sum_{i \in [m]} V_\theta(s_h^i) \right| \leq \epsilon_{\text{stat}} \right\},$$



619 then by Lemma B.2,  $\mathbb{P}(E_{t,h,f}^V) \geq 1 - \delta/(2(1 + 4/\epsilon_{\text{net}})^{2k} \binom{d}{k}^2 H^2)$  for all iterations  $t \in [(1 +$   
620  $4/\epsilon_{\text{net}})^k \binom{d}{k} H]$ , horizon  $h \in [H]$ , and  $\theta \in \mathcal{P}_h$ .

621 We can lower bound the probability of  $E$  by union bound

$$\mathbb{P}(E) \geq 1 - \sum_t \sum_{h \in [H]} \mathbb{P}(\bar{E}_{t,h}^{\mathcal{E}}) - \sum_t \sum_{h \in [H]} \sum_{\theta \in \mathcal{P}_h} \mathbb{P}(\bar{E}_{t,h,\theta}^V) \geq 1 - \delta.$$

622

□

### 623 B.3 Proof of Lemma 4.3

624 *Proof.* Let  $\hat{\theta}^* = (\hat{\theta}_0^*, \dots, \hat{\theta}_{H-1}^*)$  be the sequence of parameters such that, for each  $h \in [H]$ , the  
625 non-zero sub-vector of  $\hat{\theta}_h^*$  is in  $\mathcal{N}^k$  and is closest to the non-zero indices in  $\theta^*$ . Then, since  $\mathcal{N}^k$  is  
626  $\epsilon_{\text{net}}/2$ -maximal, we have by Assumption 1 that

$$|\langle \phi(s, a), \hat{\theta}_h^* \rangle - Q^*(s, a)| \leq |\langle \phi(s, a), \theta_h^* \rangle - Q^*(s, a)| + |\langle \phi(s, a), \hat{\theta}_h^* \rangle - \langle \phi(s, a), \theta_h^* \rangle| \leq \epsilon + \epsilon_{\text{net}},$$

627 for all  $s$  in horizon  $h$  and action  $a \in \mathcal{A}$ .

628 At iteration  $t$ , algorithm 1 deletes  $\hat{\theta}_h^*$  if and only if one of the following two cases happens: (1)  
629  $h < H - 1$ ,  $\theta_h^t = \hat{\theta}_h^*$ , and  $\hat{\mathcal{E}}_h^t > 2\epsilon + 2\epsilon_{\text{net}} + 3\epsilon_{\text{stat}}$ , (2)  $h = H - 1$ ,  $\theta_{H-1}^t = \hat{\theta}_{h+1}^*$  and  
630  $\hat{\mathcal{E}}_{H-1}^t > \epsilon + \epsilon_{\text{net}} + \epsilon_{\text{stat}}$ .

631 For any state-action pair  $(s_h, a_h)$  at level  $h$  where  $h \in [H - 1]$ , we observe by definition that

$$Q^*(s_h, a_h) - \mathbb{E}[r(s_h, a_h)] - \mathbb{E}[V_{h+1}^*(s_{h+1})] = 0.$$

632 Thus, we can upper bound  $\mathcal{E}_h^t$  by

$$\begin{aligned} \mathcal{E}_h^t &= \mathbb{E}[\langle \phi(s_h, a_h), \hat{\theta}_h^* \rangle - r(s_h, a_h) - V_{\theta_{h+1}}(s_{h+1})] \\ &\leq \mathbb{E}[(Q^*(s_h, a_h) + \epsilon + \epsilon_{\text{net}}) - r(s_h, a_h) - V_{\theta_{h+1}}(s_{h+1})] \quad (\text{By Assumption 1}) \end{aligned}$$

633 Here,  $(s_0, a_0, r_0, \dots, s_h, a_h, r_h)$  is a trajectory following  $\pi_{\theta^t}$ , and  $s_{h+1} \sim P(s_h, a_h)$ .

634 Recall that  $\theta_h^t$  is chosen by taking the function that gives maximum empirical value at level  $h$ , so

$$\frac{1}{m} \sum_{i \in [m]} V_{\theta_h^t}(s_h^i) \geq \frac{1}{m} \sum_{i \in [m]} V_{\hat{\theta}_h^*}(s_h^i),$$

635 where  $s_h^i$  are taken from the dataset  $\mathcal{D}_h^t$ . Moreover, we are conditioned under event  $E$ , so we have

$$\mathbb{E}[V_{\theta_h^*}(s_h)] - \mathbb{E}[V_{\theta_h^t}(s_h)] \leq \left( \frac{1}{m} \sum_{i \in [m]} V_{\theta_h^t}(s_h^i) + \epsilon_{\text{stat}} \right) - \left( \frac{1}{m} \sum_{i \in [m]} V_{\theta_h^*}(s_h^i) - \epsilon_{\text{stat}} \right) \leq 2\epsilon_{\text{stat}}$$

636 for all  $h$  and  $t$ .

637 For the first case, we consider  $h \in [H - 1]$ . We have

$$\begin{aligned} \mathcal{E}_h^t &\leq \mathbb{E}[(Q^*(s_h, a_h) + \epsilon + \epsilon_{\text{net}}) - r(s_h, a_h) - V_{\theta_{h+1}^t}(s_{h+1})] \\ &\leq \mathbb{E}[Q^*(s_h, a_h) - r(s_h, a_h) - (V_{f_{h+1}^*}(s_{h+1}) - 2\epsilon_{\text{stat}})] + \epsilon + \epsilon_{\text{net}} \\ &= \mathbb{E}[Q^*(s_h, a_h) - r(s_h, a_h) - f_{h+1}^*(s_{h+1}, \pi_{h+1}^*(s_{h+1}))] + \epsilon + \epsilon_{\text{net}} + 2\epsilon_{\text{stat}} \\ &\quad (\text{since } V_{f_{h+1}^*} = \max_{a \in \mathcal{A}} f_{h+1}^*(s_{h+1}, a)) \\ &\leq \mathbb{E}[Q^*(s_h, a_h) - r(s_h, a_h) - (Q^*(s_{h+1}, \pi_{h+1}^*(s_{h+1})) - \epsilon - \epsilon_{\text{net}})] + \epsilon + \epsilon_{\text{net}} + 2\epsilon_{\text{stat}} \\ &\quad (\text{By Assumption 1}) \\ &= \mathbb{E}[Q^*(s_h, a_h) - r(s_h, a_h) - Q^*(s_{h+1}, \pi_{h+1}^*(s_{h+1}))] + 2\epsilon + 2\epsilon_{\text{net}} + 2\epsilon_{\text{stat}} \\ &= 2\epsilon + 2\epsilon_{\text{net}} + 2\epsilon_{\text{stat}}. \quad (\text{since } Q^*(s_{h+1}, \pi_{h+1}^*(s_{h+1})) = V_{h+1}^*(s_{h+1})) \end{aligned}$$

638 Given that we are conditioned under event  $E$ ,  $\hat{\mathcal{E}}_h^t - \mathcal{E}_h^t \leq \epsilon_{\text{stat}}$  for all iteration  $t$  and all horizon  $h$ .

639 Thus,  $\hat{\mathcal{E}}_h^t < 2\epsilon + 2\epsilon_{\text{net}} + 3\epsilon_{\text{stat}}$ .

640 For the second case, we consider  $h = H - 2$ . We have

$$\mathcal{E}_{H-1}^t \leq \mathbb{E}[(r(s_{H-1}, a_{H-1}) + \epsilon) - r(s_{H-1}, a_{H-1})] = \epsilon + \epsilon_{\text{net}},$$

641 because  $H - 1$  is the last level.

642 Again, given that we are conditioned under event  $E$ , we have  $\hat{\mathcal{E}}_{H-1}^t - \mathcal{E}_{H-1}^t \leq \epsilon_{\text{stat}}$ , so  $\hat{\mathcal{E}}_{H-1}^t <$   
 643  $\mathcal{E}_{H-1}^t + \epsilon_{\text{stat}} \leq \epsilon + \epsilon_{\text{net}} + \epsilon_{\text{stat}}$ .  $\square$

#### 644 B.4 Proof of Lemma 4.4

645 *Proof.* Algorithm 1 terminates and returns a policy at iteration  $t$  only if  $\theta^t$  satisfies the conditions in  
 646 line 6, and by Lemma 4.3, there always exists a nice sequence of functions  $\{\hat{\theta}_h^*\}_{h=0}^{H-1}$  that satisfies  
 647 these conditions. Also, Lemma 4.1 indicates that the algorithm terminates within a finite number of  
 648 iterations. Thus, algorithm 1 is guaranteed to terminate and return a policy.

649 Let the output policy be  $\pi_{\theta^t}$ , i.e.  $\hat{\mathcal{E}}_h^t \leq 2\epsilon + 2\epsilon_{\text{net}} + 3\epsilon_{\text{stat}}$  for all  $h \in [H - 2]$  and  $\hat{\mathcal{E}}_{H-1}^t \leq$   
 650  $\epsilon + \epsilon_{\text{net}} + \epsilon_{\text{stat}}$ . The loss of this policy can be bounded by

$$\begin{aligned} V^*(s_0) - V^{\pi_{\theta^t}}(s_0) &= Q^*(s_0, \pi^*(s_0)) - V^{\pi_{\theta^t}}(s_0) \\ &\leq (\langle \phi(s_0, \pi^*(s_0)), \hat{\theta}_0^* \rangle + \epsilon + \epsilon_{\text{net}}) - V^{\pi_{\theta^t}}(s_0) && \text{(By Assumption 1)} \\ &\leq (\langle \phi(s_0, \pi_{\theta_0^t}(s_0)), \theta_0^t \rangle + \epsilon + \epsilon_{\text{net}}) - \mathbb{E}\left[\sum_{h=0}^{H-1} r(s_h, a_h)\right] \\ &\hspace{10em} \text{(since } \theta_0^t \text{ is chosen by taking the maximum)} \\ &= \epsilon + \epsilon_{\text{net}} + \mathbb{E}\left[\sum_{h=0}^{H-1} \langle \phi(s_h, a_h), \theta_h^t \rangle - r(s_h, a_h) - \langle \phi(s_{h+1}, a_{h+1}), \theta_{h+1}^t \rangle\right] \\ &\hspace{10em} \text{(telescoping sum)} \\ &= \epsilon + \epsilon_{\text{net}} + \sum_{h=0}^{H-1} \mathbb{E}\left[\langle \phi(s_h, a_h), \theta_h^t \rangle - r(s_h, a_h) - \langle \phi(s_{h+1}, a_{h+1}), \theta_{h+1}^t \rangle\right] \\ &\hspace{10em} \text{(linearity of expectation)} \\ &= \epsilon + \epsilon_{\text{net}} + \sum_{h=0}^{H-1} \mathcal{E}_h^t \leq \epsilon + \epsilon_{\text{net}} + \sum_{h=0}^{H-1} (\hat{\mathcal{E}}_h^t + \epsilon_{\text{stat}}) \leq (2\epsilon + 2\epsilon_{\text{net}} + 4\epsilon_{\text{stat}})H \end{aligned}$$

651  $\square$

## 652 C Additional Proofs

### 653 C.1 Proof of Theorem 1.1

654 *Proof.* We construct a hard instance as follows. For each  $a \in [n]$ , define  $\phi(a) = \epsilon$ . Let  $\theta^*$  be randomly  
 655 selected from  $\{-1, 1\}$ , and let  $a^*$  is uniformly chosen from  $\mathcal{A}$ . The reward  $r$  is deterministic and is  
 656 defined as

$$r(a) = \begin{cases} 2\theta^* \epsilon & \text{if } a = a^* \\ 0 & \text{otherwise.} \end{cases}$$

657 Therefore  $|r(a) - \theta^* \cdot \phi(a)| \leq \epsilon$  holds true for all actions  $a \in \mathcal{A}$ .

658 By Yao's minimax principle, it suffices to consider deterministic algorithms. Let  $A$  be a deterministic  
 659 algorithm that, by taking less than  $0.9n$  samples, returns a  $\hat{r}$  with  $|\hat{r}(a) - r(a)| < 2\epsilon$  for all  $a \in \mathcal{A}$   
 660 with probability at least 0.95. We can say the sequence of actions made by  $A$  is fixed until it receives  
 661 a reward  $r(a_t) \neq 0$  at some round  $t$ . This is because  $A$  is deterministic, and the responses  $A$  receives  
 662 are the same (i.e. all actions have reward 0) until it queries  $a^*$ . Let  $S = (a_1, \dots, a_t)$  be the sequence

663 of actions made by  $A$ . Let  $\mathcal{A}_{BAD} \subset \mathcal{A}$  be the set of actions that are not in  $S$ . We have

$$\begin{aligned} \mathbb{P}[|\hat{r}(a) - r(a)| < 2\epsilon, \forall a \in \mathcal{A}] &= \mathbb{P}[|\hat{r}(a) - r(a)| < 2\epsilon, \forall a \in \mathcal{A} | a^* \in \mathcal{A}_{BAD}] \mathbb{P}[a^* \in \mathcal{A}_{BAD}] \\ &\quad + \mathbb{P}[|\hat{r}(a) - r(a)| < 2\epsilon, \forall a \in \mathcal{A} | a^* \notin \mathcal{A}_{BAD}] \mathbb{P}[a^* \notin \mathcal{A}_{BAD}] \\ &< \mathbb{P}[|\hat{r}(a) - r(a)| < 2\epsilon, \forall a \in \mathcal{A} | a^* \in \mathcal{A}_{BAD}] \mathbb{P}[a^* \in \mathcal{A}_{BAD}] \\ &\quad + (1 - \mathbb{P}[a^* \in \mathcal{A}_{BAD}]) \end{aligned}$$

664 Since  $t < 0.9n$  and  $a^*$  is chosen uniformly random from  $\mathcal{A}$ , the probability that  $a^* \in \mathcal{A}_{BAD}$  is

$$\mathbb{P}[a^* \in \mathcal{A}_{BAD}] = \frac{|\mathcal{A}_{BAD}|}{|\mathcal{A}|} > \frac{1 - 0.9n}{n} = 0.1.$$

665 When  $a^* \in \mathcal{A}_{BAD}$ , the output of our deterministic algorithm must be fixed. We denote such output  
666 by  $r'$ . Consider a fixed  $a^* \in \mathcal{A}_{BAD}$ , if we have  $|r'(a^*) - 2\epsilon| < 2\epsilon$ , then  $r'(a^*) \in (0, 4\epsilon)$ , and  
667  $|r'(a^*) - (-2\epsilon)| > 2\epsilon$ . Similarly, if we have  $|r'(a^*) - (-2\epsilon)| < 2\epsilon$ , then  $|r'(a^*) - 2\epsilon| > 2\epsilon$ . Since  
668  $\theta^*$  is chosen uniformly random in  $\{-1, 1\}$ , we know  $r(a^*)$  is chosen uniformly random in  $\{-2\epsilon, 2\epsilon\}$ .  
669 Thus,

$$\mathbb{P}[|\hat{r}(a) - r(a)| < 2\epsilon, \forall a \in \mathcal{A} | a^* \in \mathcal{A}_{BAD}] \leq \mathbb{P}[|\hat{r}(a^*) - r(a^*)| < 2\epsilon | a^* \in \mathcal{A}_{BAD}] = 0.5.$$

670 We have

$$\mathbb{P}[|\hat{r}(a) - r(a)| < 2\epsilon, \forall a \in \mathcal{A}] < 0.5 \cdot \mathbb{P}[a^* \in \mathcal{A}_{BAD}] + (1 - \mathbb{P}[a^* \in \mathcal{A}_{BAD}]) < 0.5 \cdot 0.1 + 0.9 = 0.95.$$

671 However, by our assumption on algorithm  $A$ , we have  $\mathbb{P}[|\hat{r}(a) - r(a)| < 2\epsilon, \forall a \in \mathcal{A}] > 0.95$ .  $\square$

## 672 D Bellman Rank

673 The following definition of the general average Bellman error is helpful for our proofs in this section.

674 **Definition 5.** Given any policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ , feature function  $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$  and a sequence of  
675 parameters  $\theta = (\theta_0, \dots, \theta_{H-1})$ , the average Bellman error of  $\theta$  under roll-in policy  $\pi$  at level  $h$  is  
676 defined as

$$\mathcal{E}_h(\theta, \pi) = \mathbb{E}[\langle \phi(s_h, a_h), \theta_h \rangle - r_h - \max_{a \in \mathcal{A}} \langle \phi(s_{h+1}, a), \theta_{h+1} \rangle].$$

677 Here,  $(s_0, a_0, r_0, \dots, s_h, a_h, r_h)$  is a trajectory by following  $\pi$ , and  $s_{h+1} \sim P(s_h, a_h)$ .

678 **Definition 6 (Bellman Rank).** For a given MDP  $\mathcal{M}$ , we say that our parameter space  $\mathcal{F} = \{\theta \in \mathbb{R}^d :$   
679  $\theta$  is  $k$ -sparse,  $\|\theta\|_2 = 1\}$  has a Bellman rank of dimension  $d$  if, for all  $h \in [H]$ , there exist functions  
680  $X_h : \mathcal{F} \rightarrow \mathbb{R}^d$  and  $Y_h : \mathcal{F} \rightarrow \mathbb{R}^d$  such that for all  $\theta, \theta' \in \mathcal{F}$ ,

$$\mathcal{E}_h(\theta, \pi_{\theta'}) = \langle X_h(\theta), Y_h(\theta') \rangle.$$

681 For each  $h \in [H]$ , define  $W_h \in \mathbb{R}^{d \times d}$  as the Bellman error matrix at level  $h$ , where the  $i, j$ -th index  
682 of  $W_h$  is  $\mathcal{E}_h(\theta_i, \pi_{\theta_j})$ . Then, the Bellman rank of  $\mathcal{F}$  is the maximum among the rank of the matrices  
683  $\{W_h\}_{h \in [H]}$ .

684 We prove Proposition 1.2.

685 *Proof.* We again construct a deterministic MDP instance with binary trees. For simplicity, we  
686 assume  $d$  is a power of 2, and we construct the instance with horizon  $H = \log d$ . Thus, we have  
687  $|\mathcal{S}| = 2^H - 1 = d - 1$  states. We also assume the sparsity is  $k = 1$ , so the parameter space  $|\mathcal{F}| = d$ .  
688 The rest details of state space, action space, and the transition kernel are exactly the same as in  
689 Section 3.1.

690 The reward is defined as  $r(s, a_1) = r(s, a_2) = \epsilon$  for  $s \in \mathcal{S}_{H-1}$ , and  $r(s, a) = 0$  for all other  
691 state-action pairs. Correspondingly, the  $Q$ -function satisfies that  $Q^*(s, a) = \epsilon$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .

692 For feature at horizon  $h \in [H]$ , for  $j \geq 2^{h+1}$ , we define the  $j$ -th index of  $\phi(s, a)$  as  $\phi(s, a)[j] = j\epsilon$  for  
693 all  $(s, a) \in \mathcal{S}_h \times \mathcal{A}$ . For  $i \in [2^{h+1}]$  and  $i$  is even,  $\phi(s_{2^{h-1}+i}, a_1)[i] = \epsilon$  and  $\phi(s_{2^{h-1}+i}, a_2)[i] = 0$ .  
694 If  $i \in [2^{h+1}]$  and  $i$  is odd, then  $\phi(s_{2^{h-1}+i}, a_1)[i] = 0$  and  $\phi(s_{2^{h-1}+i}, a_2)[i] = \epsilon$ . We also define  
695  $\phi(s, a)[i] = 0$  for all other state-action pairs.

696 Notice that, for any  $h \in [H]$  and  $i \in [2^{h+1}]$ , we have can let  $\theta$  be the one-hot vector with  $i$ -th index  
 697 being 1, then  $|\langle \phi(s, a), \theta \rangle - Q^*(s, a)| \leq \epsilon$  for all  $(s, a) \in \mathcal{S}_h \times \mathcal{A}$ , so our construction satisfies  
 698 assumption 1.

699 Clearly, for each pair  $(s, a) \in \mathcal{S}_{H-1} \times \mathcal{A}$ , we can find  $\theta = (\theta_0, \dots, \theta_{H-1})$  such that the trajectory  
 700 created by following  $\pi_\theta$ , denoted by  $(s_0, a_0, r_0, \dots, s_{H-1}, a_{H-1}, r_{H-1})$ , satisfies  $s_{H-1} = s$  and  
 701  $a_{H-1} = a$ .

702 Consider two parameter candidates  $\theta, \theta'$ . Let  $(s_{H-1}, a_{H-1})$  be the state and action at level  $H - 1$   
 703 when following  $\pi_{\theta'}$ . Since we are considering deterministic MDP, we can calculate the Bellman error  
 704 at level  $H - 1$  as follows

$$\begin{aligned} \mathcal{E}_{H-1}(\theta, \pi_{\theta'}) &= \langle \phi(s_{H-1}, a_{H-1}), \theta_{H-1} \rangle - r(s_{H-1}, a_{H-1}) - \max_{a \in \mathcal{A}} \langle \phi(s_H, a), \theta_H \rangle \\ &= \langle \phi(s_{H-1}, a_{H-1}), \theta_{H-1} \rangle - \epsilon && \text{(since } H - 1 \text{ is the last level)} \\ &= \begin{cases} 0 & \text{, if } \theta_{H-1} = \theta'_{H-1} \\ -\epsilon & \text{, if } \theta_{H-1} \neq \theta'_{H-1}. \end{cases} \end{aligned}$$

705 Here, the last equality holds because, for each  $\theta_{H-1}$ , there is only one unique  $(s, a) \in \mathcal{S}_{H-1} \times \mathcal{A}$   
 706 such that  $\langle \phi(s, a), \theta_{H-1} \rangle = \epsilon$ .

707 Thus, at level  $H - 1$ , a submatrix of the Bellman error matrix,  $W \in \mathbb{R}^{d \times d}$ , satisfies

$$W_{ij} = \mathcal{E}_{H-1}(\theta_i, \pi_{\theta_j}) = \begin{cases} 0 & \text{, if } i = j \\ -\epsilon & \text{, otherwise.} \end{cases}$$

708 In other words,  $W = \epsilon(I - J)$  where  $I$  is the identity matrix and  $J$  is a  $d \times d$  matrix with all 1s.  
 709 Define matrix  $W' = \frac{1}{\epsilon}(I - 1/(n-1)J)$ , then we have

$$\begin{aligned} WW' &= (I - J)(I - \frac{1}{n-1}J) \\ &= I - \frac{1}{n-1}J - J + \frac{n}{n-1}J = I. \end{aligned}$$

710 This means  $W'$  is the inverse matrix of  $W$ , and  $W$  is full rank. Thus, the Bellman rank is at least  
 711  $d$ . □