

# Exploration-Free Reinforcement Learning with Linear Function Approximation

Anonymous authors

Paper under double-blind review

**Keywords:** exploration-free, linear function approximation, no-regret.

## Summary

In the context of Markov Decision Processes (MDPs) with linear Bellman completeness, a generalization of linear MDPs, we reconsider the learning capabilities of a *greedy* algorithm. The motivation is that, when exploration is costly or dangerous, an exploration-free approach may be preferable to optimistic or randomized solutions. We show that, under a condition of sufficient diversity in the feature distribution, Least-Squares Value Iteration (LSVI) can achieve sublinear regret. Specifically, we show that the expected cumulative regret is at most  $\tilde{O}(H^3 \sqrt{dK/\lambda_0})$ , where  $K$  is the number of episodes,  $H$  is the task horizon,  $d$  is the dimension of the feature map and  $\lambda_0$  is a measure of feature diversity. We empirically validate our theoretical findings on synthetic linear MDPs. Our analysis is a first step towards exploration-free reinforcement learning in MDPs with large state spaces.

## Contribution(s)

1. The definition of a new diversity condition for linear MDPs.  
**Context:** Inspired from prior work of [Bastani et al. \(2021\)](#) and [Kannan et al. \(2018\)](#).
2. Proved that a greedy algorithm (LSVI) achieves sublinear cumulative regret with high probability when the here defined diversity condition is satisfied.  
**Context:** Proof built upon the related work on linear contextual bandit of [Bastani et al. \(2021\)](#).
3. Proved that a greedy algorithm (LSVI) achieves sublinear cumulative regret with high probability when the here defined diversity condition is satisfied, under a misspecified setting.  
**Context:** Proof built upon the related work on approximately linear MDPs of [Zanette et al. \(2020\)](#).

# Exploration-Free Reinforcement Learning with Linear Function Approximation

Anonymous authors

Paper under double-blind review

## Abstract

1 In the context of Markov Decision Processes (MDPs) with linear Bellman complete-  
 2 ness, a generalization of linear MDPs, we reconsider the learning capabilities of a  
 3 *greedy* algorithm. The motivation is that, when exploration is costly or dangerous,  
 4 an exploration-free approach may be preferable to optimistic or randomized solutions.  
 5 We show that, under a condition of sufficient diversity in the feature distribution, Least-  
 6 Squares Value Iteration (LSVI) can achieve sublinear regret. Specifically, we show that  
 7 the expected cumulative regret is at most  $\tilde{O}(H^3 \sqrt{dK/\lambda_0})$ , where  $K$  is the number  
 8 of episodes,  $H$  is the task horizon,  $d$  is the dimension of the feature map and  $\lambda_0$  is a  
 9 measure of feature diversity. We empirically validate our theoretical findings on syn-  
 10 thetic linear MDPs. Our analysis is a first step towards exploration-free reinforcement  
 11 learning in MDPs with large state spaces.

## 12 1 INTRODUCTION

13 Reinforcement Learning (RL) is one of the most popular approaches to sequential decision making  
 14 under uncertainty. In the last few years, RL in large state spaces has received a lot of attention both  
 15 in theory (Long & Han, 2023) and practice, with applications ranging from robotics (Singh et al.,  
 16 2022) to LLM finetuning (Ahmadian et al., 2024). One great potential of RL solutions, still largely  
 17 untapped, is their intrinsically *adaptive* nature: RL agents, once deployed, can improve over time  
 18 from interaction data. This requires a careful balancing of *exploitation* (taking decisions that are  
 19 known to be good) and *exploration* (taking decisions that may be even better, but of which little is  
 20 known).

21 This exploration-exploitation dilemma is well known in the RL literature since its beginnings (Sut-  
 22 ton & Barto, 2018) and is the main subject of study of the bandit literature (Lattimore & Szepesvári,  
 23 2020) and of a good part of RL theory (Agarwal et al., 2019). All agree on this basic principle:  
 24 that some form of exploration is necessary. A purely *greedy* agent can easily get stuck on a promis-  
 25 ing course of action, without ever discovering better but neglected alternatives. Some of the most  
 26 popular exploration strategies are based on the *optimism in the face of uncertainty* principle (Lai  
 27 & Robbins, 1985), of which (Azar et al., 2017) and (Jaksch et al., 2010) are notable applications  
 28 to RL, *posterior sampling* (Thompson, 1933), like (Osband et al., 2013), or simple noise injec-  
 29 tion (Haarnoja et al., 2018).

30 In practice, however, there are several reasons to avoid exploration in favor of a greedy approach. In  
 31 safety-critical applications, such as robotic (Brunke et al., 2022), explorative actions may be danger-  
 32 ous. In many cases, exploration for the sake of learning can also be considered unethical (Bird et al.,  
 33 2016), some prominent examples being drug trials, predictive policing, lending, resume screening,  
 34 and social media personalization. It is not hard to imagine that chatbots will incur in similar ethical  
 35 issues (Følstad et al., 2021). Furthermore, explorative solutions are more expensive to implement,  
 36 their behavior is less predictable, and their decisions less interpretable. Greedy approaches are not  
 37 only favored for the aforementioned reasons, but often are also surprisingly effective in practice (e.g.,

38 [Li et al., 2024](#)). Hence, even if theory clearly shows the necessity of exploration, common sense may  
39 suggest otherwise in many real-world scenarios.

40 To reconcile theory and practice, [Bastani et al. \(2021\)](#), closely followed by [Kannan et al. \(2018\)](#),  
41 proposed to study special conditions under which exploration-free learning *is* possible. They did  
42 so within the framework of *linear contextual bandits* ([Lattimore & Szepesvári, 2020](#), Chapter 19).  
43 In this model, at each timestep  $t$ , the agent observes a context  $X_t$  (e.g., data about the current  
44 user) and selects an action  $A_t$  (e.g., an item to recommend). The agent receives a reward that is  
45 *linear* in some context-action features. Clearly, some *structure* in the rewards (such as linearity) is  
46 necessary for exploration-free learning. If rewards of different actions are completely uncorrelated,  
47 active exploration is the only way to compare the value of different actions. On the other hand,  
48 if some structure is present, an action may reveal something about other actions, reducing or even  
49 removing the need for exploration. Indeed, [Bastani et al. \(2021\)](#) show that under sufficient *diversity*  
50 of contexts, exploration-free learning is possible in linear contextual bandits. In particular, they  
51 introduce a covariate-diversity assumption and prove that the regret of a simple greedy algorithm  
52 is sublinear. This does not mean that exploration is in general unnecessary for linear contextual  
53 bandits, but provides a possible characterization of tasks for which pure exploitation suffices.

54 Our purpose is to provide a similar characterization for Markov Decision Processes with structure,  
55 showing *when* exploration-free RL is possible. To leverage results from the linear contextual bandit  
56 literature, we examine MDPs with some kind of linear structure. These are commonly studied in  
57 the context of no-regret RL with linear function approximation. This line of work was pioneered  
58 by [Jin et al. \(2023\)](#), who first designed a no-regret algorithm for finite-horizon MDPs with linear  
59 rewards and transition probabilities, also known as low-rank MDPs ([Yang & Wang, 2019](#)). The  
60 algorithm is called LSVI-UCB and is based on the optimism principle. A follow-up work by [Zanette](#)  
61 [et al. \(2020\)](#) considers a more general class of “linear” MDPs where the class of linear action-  
62 value functions is closed under the Bellman optimality operator. This is the framework that we will  
63 adopt for our analysis, although we will use low-rank MDPs as numerical examples.<sup>1</sup> Nonlinear  
64 function approximation is also an active area of research (e.g., [Jin et al., 2021](#)). This is beyond  
65 the scope of this paper, but we believe that our analysis of linear MDPs is a necessary step in  
66 the study of exploration-free reinforcement learning in complex environments requiring general  
67 function approximation.

68 Our main contributions are as follows: we define a novel diversity condition, inspired by [Bastani](#)  
69 [et al. \(2021\)](#) and [Kannan et al. \(2018\)](#), for Markov Decision Processes with linear function approxi-  
70 mation, and present new insights into how feature coverage affects the performance of exploration-  
71 free reinforcement learning algorithms. We prove that a greedy algorithm (LSVI) achieves sublinear  
72 cumulative regret with high probability when the diversity condition is satisfied. We also establish an  
73 *any-time* bound on the expected cumulative regret. Finally, we empirically validate our theoretical  
74 findings on synthetic linear MDPs.

75 The paper is structured as follows. In Section 2 we present all the necessary preliminaries for  
76 understanding and developing the concepts discussed in this work. We begin by introducing Markov  
77 Decision Processes (MDPs), followed by the specific case of MDPs that satisfy the linear Bellman  
78 completeness condition, which is the setting of this work. We also consider the special case of  
79 low-rank MDPs. Section 3 describes the analyzed algorithm, outlines the assumptions required for  
80 our analysis, and presents the theoretical results. Section 4 provides more details on the theoretical  
81 analysis, where we state the key lemmas used in the proof of the main theorem, followed by a  
82 detailed proof of the latter. Other proofs can be found in the Appendix. In Section 5, we discuss  
83 related works, while Section 6 focuses on the experiments conducted to empirically validate our  
84 theoretical results.

---

<sup>1</sup>A more intuitive generalization of low-rank MDPs is linear realizability of action-value functions. However, this has so far proven to be much more challenging to analyze ([Weisz et al., 2023](#)).

Table 1: Notation

SYMBOL	DESCRIPTION
$[n]$	$\{1, \dots, n\}$
$\mathbb{I}$	Indicator function
$\lambda_{\min}(A)$	Minimum eigenvalue of matrix $A$
$\langle x, y \rangle$	Inner product, $\langle x, y \rangle = \sum_i x_i y_i$
$\ x\ _p$	p-norm of vector $x$
$\ x\ _A^2$	$x^\top A x$

## 85 2 PRELIMINARIES

86 In this section, we provide the necessary background on Markov decision processes and the linearity  
 87 assumption under which our work is conducted. Our notation is summarized in Table 1.

### 88 2.1 Markov Decision Processes

89 A finite-horizon Markov Decision Process (MDP, [Puterman, 1994](#)) is denoted by the tuple  $M =$   
 90  $(\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r, \mu)$ , where  $\mathcal{S}$  is the space of states,  $\mathcal{A}$  is the space of actions,  $H \in \mathbb{N}$  is the length of  
 91 each episode,  $\mathbb{P} = \{\mathbb{P}_h\}_{h=1}^H$  and  $r = \{r_h\}_{h=1}^H$  are, respectively, the state transition probabilities and  
 92 the reward functions. We assume that  $\mathcal{S}$  is a measurable space and  $\mathcal{A}$  has finite cardinality. For each  
 93 step  $h \in [H]$ ,  $\mathbb{P}_h(\cdot|s, a)$  denotes the transition kernel over the next states if we choose action  $a$  in  
 94 state  $s$ , and  $r_h : \mathcal{S} \times \mathcal{A} \rightarrow [-1, 1]$  is the deterministic reward function. Finally  $\mu$  is the starting-state  
 95 probability distribution over  $\mathcal{S}$ .

96 An agent interacts with the MDP as follows: an initial state  $s_1$  is drawn from  $\mu$ , then at each step  
 97  $h \in [H]$  the agent observes the state  $s_h$ , picks an action  $a_h$  and receives a reward  $r_h(s_h, a_h)$ .  
 98 The MDP evolves into a new state  $s_{h+1}$  that is drawn from the transition kernel  $\mathbb{P}_h(\cdot|s_h, a_h)$ . The  
 99 episode ends when state  $s_{H+1}$  is reached. A (deterministic) policy  $\pi$  of an agent is a function  
 100  $\pi : \mathcal{S} \times [H] \rightarrow \mathcal{A}$ , where  $\pi(s, h)$  is the action that the agent takes in state  $s$  at the  $h$ -th step of the  
 101 episode. We will abbreviate  $\pi(s, h)$  as  $\pi_h(s)$  in the following. For a policy  $\pi$ , for each  $h \in [H]$ ,  
 102 we can define the value function  $V_h^\pi : \mathcal{S} \rightarrow \mathbb{R}$ , which, given the current state at step  $h$ , returns the  
 103 cumulative expected reward following policy  $\pi$ :

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

104 where  $\mathbb{E}_\pi$  is short for  $a_h \sim \pi(\cdot|s_h), s_{h+1} \sim \mathbb{P}(\cdot|s_h, a_h), \dots, a_H \sim \pi(\cdot|s_H)$  conditional on  $\pi$ . We  
 105 also define the action-value function  $Q_h^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , which gives the expected value of cumula-  
 106 tive rewards when the agent starts from a given state-action pair at the  $h$ -th step and follows policy  
 107  $\pi$  afterwards. We have:

$$Q_h^\pi(s, a) := r_h(s, a) + \mathbb{E}_\pi \left[ \sum_{h'=h+1}^H r_{h'}(s_{h'}, \pi_{h'}(s_{h'})) \middle| s_h = s, a_h = a \right],$$

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

109 Finally, we can define the occupancy measure of the policy  $\pi$ :

$$\rho_h^\pi(s) := \mathbb{E}_{\pi, s_0 \sim \mu} [\mathbb{I}\{s_h = s\}].$$

110 There always exists an optimal deterministic policy  $\pi^*$  which gives the optimal value  $V_h^*(s) =$   
 111  $\sup_\pi V_h^\pi(s)$  for all  $s \in \mathcal{S}$  and  $h \in [H]$  ([Puterman, 1994](#)). Similarly,  $Q_h^*(s, a) = Q_h^{\pi^*}(s, a)$ .

112 In an episodic MDP, the agent aims to learn the optimal policy by interacting with the environment  
 113 over a series of  $K$  episodes. For each  $k \geq 1$ , an initial state  $s_1^k$  is drawn from  $\mu$  and the agent  
 114 chooses policy  $\pi_k$ . The difference in values between  $V_1^{\pi_k}(s_k)$  and  $V_1^*(s_k)$  is the instantaneous  
 115 regret, or suboptimality, of the agent at the  $k$ -th episode. Thus, after playing for  $K$  episodes, the  
 116 total regret is

$$R(K) := \sum_{k=1}^K \mathbb{E}_{s_k \sim \mu} [V_1^*(s_k) - V_1^{\pi_k}(s_k)].$$

117 We can also rewrite the total regret, by using a performance difference lemma (e.g., Proposition 29  
 118 from Papini et al. (2021a)), as follows:

$$R(K) := \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{s_h \sim \rho_h^{\pi_k}} [\Delta_h(s_h, \pi_k(s_h))], \quad (1)$$

119 where  $\Delta_h(s, a) := V_h^*(s) - Q_h^*(s, a)$  is the suboptimality gap.

## 120 2.2 Linear Bellman Completeness

121 We will consider a setting in which we have a set of features that satisfy the linear Bellman com-  
 122 pleteness condition, which we will refer to as linear MDPs for brevity. In this scenario we work with  
 123 a feature map  $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ . Let us first define the set of admissible parameters as:

$$\mathcal{W} = \{\mathbf{w} \in \mathbb{R}^d \text{ s.t. } |\langle \phi(s, a), \mathbf{w} \rangle| \leq H \ \forall s \in \mathcal{S}, \forall a \in \mathcal{A}\}.$$

124 We restrict our analysis to MDPs equipped with a feature map that satisfies the following:

125 **Assumption 2.1** (Linear Bellman completeness, Agarwal et al. (2021a)). *We say that the feature*  
 126 *map  $\phi$  satisfies the linear Bellman completeness property if, for all  $\theta \in \mathcal{W}$  and  $(s, a, h) \in \mathcal{S} \times \mathcal{A} \times$*   
 127  *$[H]$ , there exists  $\mathbf{w} \in \mathcal{W}$  such that:*

$$\mathbf{w}^\top \phi(s, a) = r(s, a) + \mathbb{E}_{s' \sim \mathbb{P}_h(s, a)} \max_{a'} \theta^\top \phi(s', a').$$

128 This condition implies that  $Q_h^*(s, a)$  is linear in  $\phi$ , i.e., there exists  $\theta_h^*$  such that  $Q_h^*(s, a) =$   
 129  $(\theta_h^*)^\top \phi(s, a)$  (Zanette et al., 2020, Lemma 6). This justifies the use of linear function approxi-  
 130 mation.

## 131 2.3 Low-Rank Markov Decision Processes

132 Although our theoretical results apply to general linear-Bellman-complete MDPs, we mention a  
 133 particular case in which Assumption 2.1 holds, *low-rank* Markov Decision Processes (Jin et al.,  
 134 2023). In this scenario, the transition kernel and the reward function are assumed to be linear w.r.t.  
 135 known state-action features.

136 Formally, a Markov Decision Process defined as  $M = (\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r)$ , with a feature map  $\phi :$   
 137  $\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ , is considered a *low-rank MDP* (Yang & Wang, 2019; Jin et al., 2023) if, for each  
 138 time step  $h \in [H]$ , there exist  $d$  signed measures  $\rho_h = (\rho_h^{(1)}, \dots, \rho_h^{(d)})$  over the state space  $\mathcal{S}$ , and a  
 139 vector  $\theta_h \in \mathbb{R}^d$ , such that, for any state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , the following holds:

$$\mathbb{P}_h(\cdot | s, a) = \langle \phi(s, a), \rho_h(\cdot) \rangle, \quad r_h(s, a) = \langle \phi(s, a), \theta_h \rangle.$$

140 A key characteristic of a low-rank MDP is that the action-value functions of *all* policies are linear  
 141 with respect to the same feature map  $\phi$  (Jin et al., 2023, Proposition 2.3). It is easy to show that all  
 142 low-rank MDPs are linear-Bellman-complete. The opposite is not true (Zanette et al., 2020).

### 3 GREEDY LEARNING

In this section, after reviewing the LSVI algorithm, we present our feature-diversity assumption and show how this is sufficient to achieve sublinear regret in an exploration-free manner.

#### 3.1 Algorithm

The algorithm we consider in our work is Least-Square Value Iteration (LSVI, Bradtke & Barto, 1996), a simple *greedy* algorithm, based on value-iteration, which finds the optimal Q-function by iterative application of Bellman’s optimality equation:

$$Q_h^*(s, a) = r_h(s, a) + \mathbb{E}_{s' \sim \mathbb{P}_h(\cdot | s, a)} \max_{a' \in \mathcal{A}} Q_{h+1}^*(s', a').$$

LSVI parametrizes  $Q_h^*(s, a)$  by a linear form and approximates the optimality equation with a regularized least-squares problem in which we solve for  $\mathbf{w}_h$ . The algorithm solves the following program at each stage of each episode:

$$\mathbf{w}_h \leftarrow \underset{\mathbf{w} \in \mathcal{W}}{\operatorname{argmin}} \sum_{\tau=1}^{k-1} [r_h(s_h^\tau, a_h^\tau) + \max_{a \in \mathcal{A}} Q_{h+1}(s_{h+1}^\tau, a) - \mathbf{w}^\top \phi(s_h^\tau, a_h^\tau)]^2 + \lambda \|\mathbf{w}\|^2.$$

---

#### Algorithm 1 LSVI

---

```

1: for episode  $k = 1, \dots, K$  do
2:   Observe the initial state  $s_1^k \sim \mu$ 
3:   for step  $h = H, \dots, 1$  do
4:      $\hat{\Sigma}_{k,h} = \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top + \lambda \cdot \mathbf{I}$ 
5:      $\tilde{\mathbf{w}}_h^k = \hat{\Sigma}_{k,h}^{-1} \sum_{\tau=1}^{k-1} \phi(s_h^\tau, a_h^\tau) [r_h(s_h^\tau, a_h^\tau) + \max_a Q_{h+1}^k(s_{h+1}^\tau, a)]$ 
6:      $\hat{\mathbf{w}}_h^k = \arg \min_{\mathbf{w} \in \mathcal{W}} \|\mathbf{w} - \tilde{\mathbf{w}}_h^k\|_{\hat{\Sigma}_{k,h}}$ 
7:      $Q_h^k = \langle \hat{\mathbf{w}}_h^k, \phi(\cdot, \cdot) \rangle$ 
8:   end for
9:   for step  $h = 1, \dots, H$  do
10:    take action  $a_h^k = \arg \max_{a \in \mathcal{A}} Q_h^k(s_h^k, a)$  and observe  $s_{h+1}^k$ 
11:   end for
12: end for
```

---

At a high level, each episode involves two main passes through all time-steps. The first *backward* pass (lines 3-8) updates  $\hat{\mathbf{w}}_h^k$  and  $\hat{\Sigma}_{k,h}$ , that are, respectively, the parameters we are trying to estimate and the covariance matrix, which are used to construct the action-value function  $Q_h^k$ . In the second pass (lines 9-11), the greedy policy is executed:  $a_h^k = \arg \max_{a \in \mathcal{A}} Q_h^k(s_h^k, a)$ , using the  $Q_h^k$  computed in the first pass. It’s important to note that  $Q_{H+1} \equiv 0$  since no reward is given after the  $H$ -th step. In the first episode ( $k = 1$ ), the summations in lines 4 and 5 run from  $\tau = 1$  to 0, meaning  $\hat{\Sigma}_{1,h} = \lambda \cdot \mathbf{I}$  and  $\hat{\mathbf{w}}_h^1 = 0$ . The inverse covariance matrix can be updated directly using Sherman-Morrison’s formula for improved computational complexity. Line 6 is a projection step ensuring  $\hat{\mathbf{w}}_h^k \in \mathcal{W}$ .<sup>2</sup>

#### 3.2 Assumptions

We will now outline the assumptions necessary for our regret analysis. The first is a technical one on the parameter set:

**Assumption 3.1.**  $\mathcal{W}$  is a convex set. Moreover, there exists a constant  $\phi_{\max}$  such that  $\|\phi(s, a)\|_2 \leq \phi_{\max}$  for all  $s, a$ , and a constant  $w_{\max}$  such that  $\|\mathbf{w}\|_2 \leq w_{\max}$  for all  $\mathbf{w} \in \mathcal{W}$ .

<sup>2</sup>It is more common to directly clip the Q-function estimate in  $[-H, H]$ . However, for technical reasons, we need to preserve the linearity of the estimator.

167 The most important assumption is the following, inspired by conceptually similar conditions pro-  
 168 posed by Bastani et al. (2021) and Kannan et al. (2018) for linear contextual bandits:

169 **Assumption 3.2. (Covariate Diversity).** *There exists a positive constant  $\lambda_0$  such that, for each*  
 170 *policy  $\pi$ ,  $\mathbf{w} \in \mathcal{W}$ , and for each  $h \in [H]$ ,*

$$\lambda_{\min} \left( \mathbb{E}_{s \sim \rho_h^\pi(s)} [\phi(s, \pi(s)) \phi(s, \pi(s))^\top \mathbb{I}\{\langle \phi(s, \pi(s)), \mathbf{w} \rangle \geq \max_{a \in \mathcal{A}} \langle \phi(s, a), \mathbf{w} \rangle\}] \right) \geq \lambda_0.$$

171 Intuitively, the feature vectors witnessed by the agent in “sensible” rounds must cover the whole  
 172 feature space. Fix a linear Q-function estimator. A round is “sensible” if the agent plays an action  
 173 that would appear optimal according to the Q-function estimate. It must hold true for all determin-  
 174 istic policies the agent may play, all linear Q-function estimators, and separately for each episode’s  
 175 timestep. This is a joint property of the MDP and of the feature map. It is encouraged by feature  
 176 maps showing great diversity across states, but also by strongly connected MDPs and starting-states  
 177 distributions with a large support. The constant  $\lambda_0$  is a measure of diversity. We expect exploration-  
 178 free learning to be easier when  $\lambda_0$  is larger.

179 A simple example where Assumption 3.2 holds is the following. For simplicity we consider two  
 180 actions and  $d = 1$ , but similar constructions can be made for a generic number of actions and a  
 181 larger feature dimension.

182 **Proposition 3.3** (Noisy features). *Let  $|\mathcal{A}| = 2$  and  $\phi(s, a) = f(s, a) + \eta(a)$  for some function*  
 183  *$f : \mathcal{S} \times \mathcal{A} \rightarrow [0, \sqrt{2}\sigma]$  and independent Gaussian noises  $\eta(a) \sim \mathcal{N}(0, \sigma^2)$ . Then Assumption 3.2*  
 184 *holds with  $\lambda_0 \geq 0.2\sigma^2$ .*

### 185 3.3 Regret of LSVI with Covariate Diversity

186 We now establish an upper bound on the cumulative regret of LSVI in the case of an MDP whose  
 187 representation satisfies both the Assumption 2.1 and Assumption 3.2.

188 **Theorem 3.4.** *Under Assumptions 2.1, 3.1, and 3.2, with probability  $1 - \delta$ , the cumulative regret of*  
 189 *LSVI is at most:*

$$R(K) = \mathcal{O} \left( H^3 \sqrt{\frac{dK}{\lambda_0}} \log(K/\delta) \right).$$

190 Notice that Algorithm 1 is not parametric in the failure probability  $\delta$ . By setting this free parameter  
 191 to  $\delta = 1/\sqrt{K}$ , by a standard argument, we obtain an upper bound on the *expected* regret, where the  
 192 extra expectation is over the random sequence of (deterministic) policies played by LSVI.

193 **Corollary 3.5.** *Under the same assumptions as Theorem 3.4, the expected cumulative regret of LSVI*  
 194 *is at most:*

$$\mathbb{E} [R(K)] = \mathcal{O} \left( H^3 \sqrt{\frac{dK}{\lambda_0}} \log(K) \right).$$

195 The result is still *any-time*, that is, the algorithm does not need to know the number of episodes  $K$   
 196 in advance.

197 Our regret upper bounds, scaling with  $\sqrt{d}$ , seem to contradict existing  $\Omega(d\sqrt{K})$  lower bounds  
 198 (cf. Zanette et al. (2020), Theorem 2). This may actually be possible under the non-standard As-  
 199 sumption 3.2. Anyway, notice that  $\lambda_0 \leq 1/d$ , the minimum eigenvalue of the covariance matrix  
 200 of a D-optimal design (Lattimore et al., 2020). Hence, linear dependence on the dimension of the  
 201 feature map is not avoided. If  $\lambda_0 \simeq 1/d$ , LSVI with covariate diversity has a better dependence than  
 202 LSVI-UCB ( $d\sqrt{d}$ ) and matches that of the computationally inefficient ELEANOR (Zanette et al.,  
 203 2020). This is possible thanks to the linearity of the Q-function estimates, while LSVI-UCB incurs  
 204 an extra  $\sqrt{d}$  factor due to its nonlinear exploration bonuses.



### 3.4 Misspecification

Our results extend to the case where the MDP is only *approximately* linear. In particular, we consider the notion of *low inherent Bellman error* introduced by Zanette et al. (2020):

**Assumption 3.6.**

$$\sup_{\mathbf{w}' \in \mathcal{W}} \inf_{\mathbf{w} \in \mathcal{W}} \sup_{s \in \mathcal{S}, a \in \mathcal{A}} \left| \langle \phi(s, a), \mathbf{w} \rangle - r_h(s, a) - \mathbb{E}_{s' \sim \mathbb{P}_h(\cdot | s, a)} \left[ \max_{a'} \langle \phi(s', a'), \mathbf{w}' \rangle \right] \right| \leq \zeta.$$

The constant  $\zeta$  measures the level of misspecification, and the linear Bellman completeness case we considered so far corresponds to  $\zeta = 0$ , no misspecification. The optimal action-value function is no longer linear, but is well approximated by a linear function (Zanette et al., 2020, Lemma 6). Our results generalize well to this misspecified setting.

**Theorem 3.7.** *If Assumptions 3.6 and 3.2 are satisfied, with probability  $1 - \delta$ , the cumulative regret of LSVI is at most:*

$$R(K) = \tilde{O} \left( H^3 \sqrt{\frac{dK}{\lambda_0}} + H^2 \zeta \frac{K}{\sqrt{\lambda_0}} + H^2 \zeta K \right).$$

With misspecification, the linear term in  $K$  is inevitable (Zanette et al., 2020), but is controlled by  $\zeta$ , which is supposed to be very small. In fact, our result seems to violate a fundamental  $\Omega(\zeta \sqrt{dK})$  lower bound (Lattimore et al., 2020). Again, this is not the case since  $\lambda_0 \leq 1/d$ , making the second term in the regret  $\tilde{O}(H^2 \zeta \sqrt{dK})$ .

## 4 ANALYSIS

In this section, we prove our main result, Theorem 3.4. We first provide two fundamental lemmas, whose proofs are deferred to Appendix A and B.

The first lemma provides an upper bound on the difference between the estimated  $Q$ -function at episode  $k$  and step  $h$ , and the actual optimal  $Q$ -function.

**Lemma 4.1.** *Assume  $\lambda_{\min}(\hat{\Sigma}_{k,h}) \geq \lambda_k$  for all  $k \geq 1$  and  $h \in [H]$ . Under Assumptions 2.1 and 3.1, with probability  $1 - \delta/2$ , for all  $k \geq 1, h \in [H], s \in \mathcal{S}, a \in \mathcal{A}$ :*

$$|\hat{Q}_h^k(s, a) - Q_h^*(s, a)| \leq (H - h) \left( \left( \sqrt{\beta_k(\delta)} \right) \frac{\phi_{\max}}{\sqrt{\lambda_k}} \right),$$

where

$$\sqrt{\beta_k(\delta)} := H\sqrt{A + B + C} + 1 + w_{\max},$$

and  $A := d \ln \left( 1 + \frac{\phi_{\max}^2 k}{d} \right)$ ,  $B := d \ln(w_{\max}^2 \phi_{\max}^2 k)$ ,  $C := \ln(2H\delta^{-1})$ .

Next, we show that the minimum eigenvalue of the sample covariance matrix at time step  $h$  until episode  $k$ ,  $\lambda_{\min}(\hat{\Sigma}_{k,h})$ , grows linearly with  $k$ . This will guarantee the convergence of our regression estimate.

**Lemma 4.2.** *Given Assumptions 3.1 and 3.2, the following holds for the minimum eigenvalue of the empirical covariance matrix for each  $h \in [H]$  and for each  $k \geq 1$ :*

$$\mathbb{P} \left[ \lambda_{\min}(\hat{\Sigma}_{k,h}) \geq \lambda + \lambda_0 k - 8\phi_{\max}^2 \sqrt{k \log(4dk/\delta)} \right] \geq 1 - \frac{\delta}{2}$$



232 *Proof of Theorem 3.4.* Given the representation of regret from Equation (1):

$$\begin{aligned}
R(K) &= \sum_{k=1}^K \mathbb{E}_{s \sim \rho_h^{\pi_k}} \left[ \sum_{h=1}^H \Delta_h(s, \pi_k(s)) \right] \\
&= \sum_{k=1}^K \mathbb{E}_{s \sim \rho_h^{\pi_k}} \left[ \sum_{h=1}^H Q_h^*(s, \pi_h^*(s)) - Q_h^*(s, \pi_h^k(s)) \right] \\
&= \sum_{k=1}^K \mathbb{E}_{s \sim \rho_h^{\pi_k}} \left[ \sum_{h=1}^H Q_h^*(s, \pi_h^*(s)) - \widehat{Q}_h^k(s, \pi_h^k(s)) + \widehat{Q}_h^k(s, \pi_h^k(s)) - Q_h^*(s, \pi_h^k(s)) \right] \\
&\leq \sum_{k=1}^K \mathbb{E}_{s \sim \rho_h^{\pi_k}} \left[ \sum_{h=1}^H Q_h^*(s, \pi_h^*(s)) - \widehat{Q}_h^k(s, \pi_h^*(s)) + \widehat{Q}_h^k(s, \pi_h^k(s)) - Q_h^*(s, \pi_h^k(s)) \right] \\
&\leq \sum_{k=1}^K \mathbb{E}_{s \sim \rho_h^{\pi_k}} \left[ \sum_{h=1}^H 2 \sup_{a \in \mathcal{A}} |\widehat{Q}_h^k(s, a) - Q_h^*(s, a)| \right] \\
&\leq \sum_{k=1}^K \sum_{h=1}^H 2(H-h) \sqrt{\beta_k(\delta)} \frac{\phi_{\max}}{\sqrt{\lambda_k}} \quad \left( \text{w.p. } 1 - \frac{\delta}{2} \right) \\
&\leq \sum_{k=1}^K H(H+1) \sqrt{\beta_k(\delta)} \frac{\phi_{\max}}{\sqrt{\lambda_k}} \\
&= \mathcal{O} \left( H^3 \sqrt{d} \log(K/\delta) \sum_{k=1}^K \frac{1}{\sqrt{\lambda_0 k}} \right) \quad \left( \text{w.p. } 1 - \frac{\delta}{2} \right) \\
&= \mathcal{O} \left( H^3 \sqrt{\frac{dK}{\lambda_0}} \log(K/\delta) \right),
\end{aligned}$$

233 The first inequality is by definition of the greedy algorithm (Alg. 1, line 10). The third inequality  
234 follows from Lemma 4.1 with a choice of  $\lambda_k$  provided by Lemma 4.2, and the final inequality is an  
235 elementary upper bound on the sum  $\sum_{h=1}^H 2(H-h)$ . The last two equalities result from bounding  
236  $\beta_k(\delta)$  by  $\beta_K(\delta)$ , and then substituting  $\sqrt{\beta_k(\delta)} = \mathcal{O}(H\sqrt{d\log(k/\delta)})$  as defined in Lemma 4.1.  
237 Additionally, we use the fact that  $1/(\sqrt{k} - \mathcal{O}(\sqrt{k\log(k)})) = \mathcal{O}(\sqrt{\log(k)/k})$  and this derives from  
238 the definition of  $\lambda_k$  in Lemma 4.2.  $\square$

## 239 5 RELATED WORKS

240 Our work is mainly inspired by the one of Bastani et al. (2021) on linear contextual bandits. They  
241 first introduced a covariate-diversity assumption that allows a greedy algorithm to achieve sublinear  
242 regret. In the same paper, they also proposed a *greedy-first* algorithm that operates greedily until it  
243 detects that convergence to the optimal policy is unlikely, at which point it begins exploring. This  
244 was shown to outperform existing exploration-based bandit algorithms. A similar analysis of greedy  
245 algorithms, using a slightly different diversity condition, was carried out by Kannan et al. (2018).  
246 To the best of our knowledge, we are the first to extend this kind of analysis to MDPs.

247 With similar motivations, but a radically different approach, Saleh et al. (2022) studied noise-free  
248 reinforcement learning in MDPs with Lipschitz-continuous transition models. They proposed a  
249 regularized policy gradient approach called truly deterministic policy optimization.<sup>3</sup> They proved

<sup>3</sup>The name is a reference to the more popular *deterministic policy gradient* algorithms (Silver et al., 2014; Lillicrap et al., 2016; Fujimoto et al., 2018). These optimize a deterministic parametric policy using data collected by a stochastic counterpart obtained by noise injection. For this reason, despite the name, they cannot be considered exploration-free. The reasons for deploying deterministic policies are similar to ours, but are only applied to the final product of learning and not to the learning process itself. See Montenegro et al. (2024) for a recent investigation of this approach.

monotonic improvement guarantees, but did not study sample complexity nor regret. Their experiments show promising results in robotics applications, but are essentially incomparable to ours. Also, because of the (deterministic) policy regularization, their algorithm may not be considered greedy.

## 5.1 Diversity Conditions

Besides (Bastani et al., 2021) and (Kannan et al., 2018), which serve as the foundation for our work, several papers have adopted covariate-diversity assumptions in linear contextual bandits for diverse reasons, often as mere technical assumptions (Foster et al., 2019; Chatterji et al., 2020; Ghosh et al., 2021; Hao et al., 2020; Wu et al., 2020; Tirinzoni et al., 2022), sometimes, with a representation learning perspective, as a characterization of “good” feature maps (Papini et al., 2021b; Tirinzoni et al., 2022; 2023). See (Papini et al., 2021b) for a discussion and comparison of the different conditions.

For linear MDPs, Papini et al. (2021a) proposed a diversity condition, called UNISOFT, under which LSVI-UCB and other optimistic algorithms achieve *constant* regret (under a minimum-gap assumption). In our notation, they require:

$$\lambda_{\min} \left( \mathbb{E}_{s \sim \rho_h^{\pi^*}(s)} [\phi(s, \pi^*(s)) \phi(s, \pi^*(s))^{\top}] \right) \geq \lambda^*. \quad (2)$$

It is sufficient to observe that  $\langle \phi(s, \pi^*(s)), \mathbf{w}^* \rangle \geq \max_a \langle \phi(s, a), \mathbf{w}^* \rangle$ , where  $\mathbf{w}^*$  is the linear parameter of  $Q^*$ , to see that Assumption 3.2 implies UNISOFT and  $\lambda_0 \leq \lambda^*$ . Intuitively, UNISOFT requires optimal trajectories to be informative, while we ask the same of all trajectories that are optimal *according to some linear estimate*. In fact, our design of Assumption 3.2 was partly inspired by UNISOFT.

In the theory of policy gradient algorithms, *concentrability coefficients* (Agarwal et al., 2021b) play a similar role than our covariate-diversity assumption: by assuming that the starting-state distribution covers well the subset of the state space visited by the optimal policy, they remove part of the challenges of exploration. This allows to study policy gradient algorithms with simple exploration strategies (e.g., noise injection) from the perspective of stochastic optimization. The algorithms considered in this line of work always employ stochastic policies, at least for data collection. Our covariate diversity assumption is also reminiscent of some *coverage ratios* used in offline RL (Uehara & Sun, 2022) with linear function approximation and of some notions of coverability (Xie et al., 2023).

## 5.2 Safe Exploration

A related body of literature focuses on developing reinforcement learning algorithms that prioritize controlled exploration for ethical and safety reasons, moving away from conventional exploratory methods. This challenge is well outlined in Amodei et al. (2016), which identifies several AI safety problems. These include issues like “avoiding side effects” and “reward hacking,” where agents can inadvertently perform harmful actions due to poorly designed objective functions, but also concerns regarding undesirable behavior during the learning process, a problem known as *safe exploration*.

The concept of safe exploration was first introduced by Moldovan & Abbeel (2012), who presents an algorithm for safe but potentially suboptimal exploration in Markov Decision Processes (MDPs). A key contribution of their work is a formal definition of safety that focuses on maintaining ergodicity with a controlled probability. Although the problem is NP-hard, the authors propose an approximation scheme that balances safety and performance.

A natural way to ensure safe exploration is by adding constraints to the MDP (Altman, 2021) and enforcing them during the learning process. A notable example is the Constrained Policy Optimization (CPO) algorithm by Achiam et al. (2017), that ensures near-satisfaction of safety constraints at each iteration.

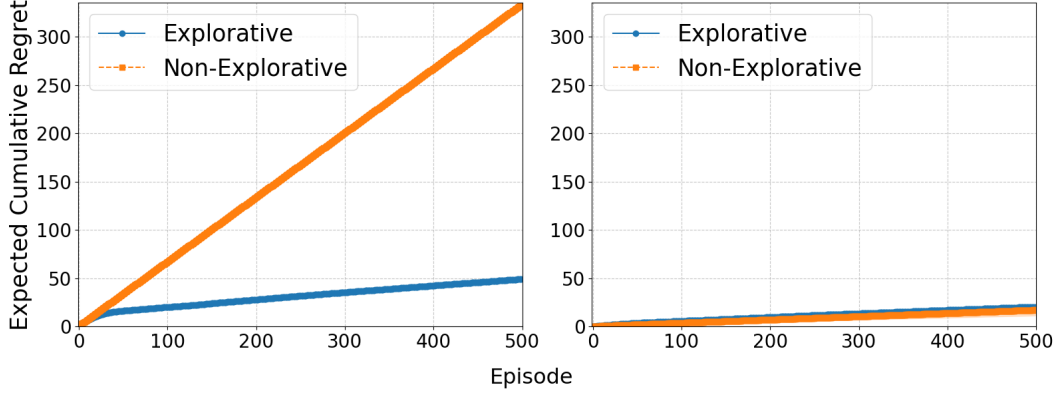


Figure 1: Expected cumulative regret of LSVI (without exploration) and LSVI-UCB (with exploration), with (right) and without (left) covariate diversity.

Safe exploration can also be defined in terms of stability, as in (Berkenkamp et al., 2017). A model-based predictive approach allows the agent to avoid exploratory actions that may lead to irrecoverable states.

Finally, with an appropriately designed reward function, safety during the learning process can be ensured by enforcing monotonic performance improvement (Papini et al., 2022).

## 6 EXPERIMENTS

In this section, we present the experimental results obtained from synthetic low-rank MDPs (as defined in Section 2.3). For each scenario, we evaluate the performance of LSVI and of an optimistic algorithm, LSVI-UCB (Jin et al., 2023).

**Synthetic Problems.** We define a randomly generated low-rank MDP with two distinct realizable linear representations, both obtained following Example 2.2 in Jin et al. (2023), each exhibiting different characteristics. The first representation satisfies covariate diversity by generating random features (cf. Proposition 3.3), while the second is specifically constructed to violate the diversity assumption, using orthogonal features to simulate an MDP without any correlation between actions (a tabular MDP). The purpose of these experiments is to demonstrate the varying behavior of LSVI-UCB and LSVI across different MDPs and to highlight the impact of covariate diversity, comparing the cases where it is satisfied and where it is not. We conduct our experiments in a setting where  $H = 3$ ,  $d = 10$  and  $K = 500$ . Each experiment is replicated 30 times under these same parameters. We plot cumulative regret normalized by  $V^*$ , averaged over the independent runs, with shaded areas corresponding to one standard deviation.

**Covariate vs Non-Covariate Diversity.** We construct two different environments with randomly generated parameters and compare, normalizing both with respect to their  $V^*$ , the expected cumulative regret obtained when using a representation that satisfies covariate diversity against one that does not. In the left image of Fig. 1, it is evident that the absence of covariate diversity causes the expected cumulative regret to increase linearly with LSVI (exploration-free). Conversely, the right image, depicting a setting satisfying covariate diversity, shows sub-linear curves for both LSVI (exploration-free) and LSVI-UCB (optimistic). Fig. 2 shows a close-up of the results under covariate diversity. Additionally, in Fig. 3, we show the average over 100 different MDPs with covariate diversity for both LSVI (non-explorative) and LSVI-UCB (explorative).

As shown by the experimental results, our findings align with the theory. Specifically, in both Fig. 1 and Fig. 3, we observe that the presence of covariate diversity makes exploration-free learning feasible. In Fig. 3, we average the (normalized) performance across 100 randomly generated MDPs that satisfy covariate diversity, with each parameter sampled from a uniform distribution between 0

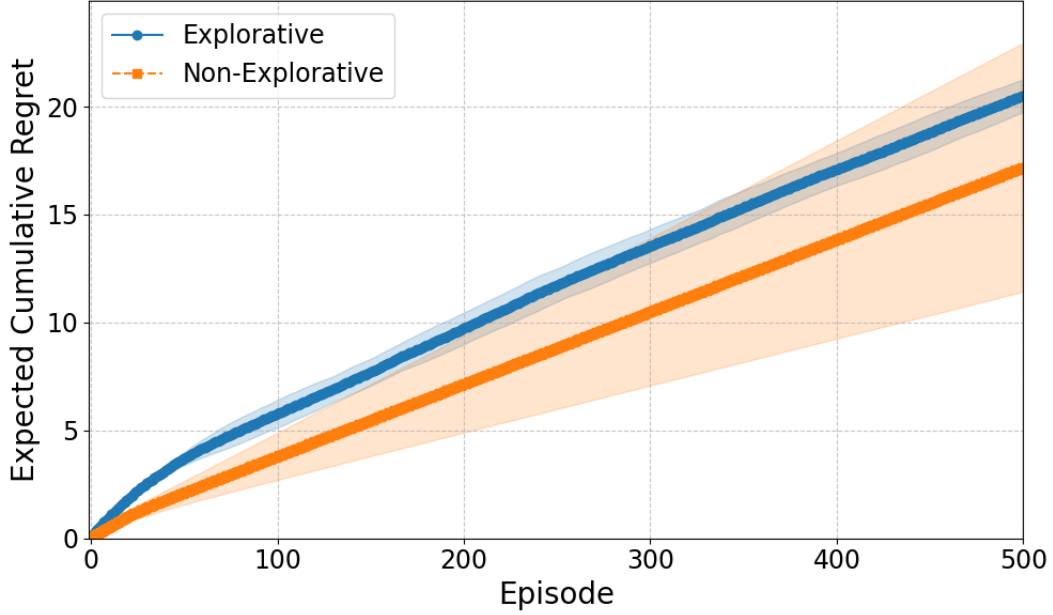


Figure 2: Expected cumulative regret of LSVI (without exploration) and LSVI-UCB (with exploration), with covariate diversity.

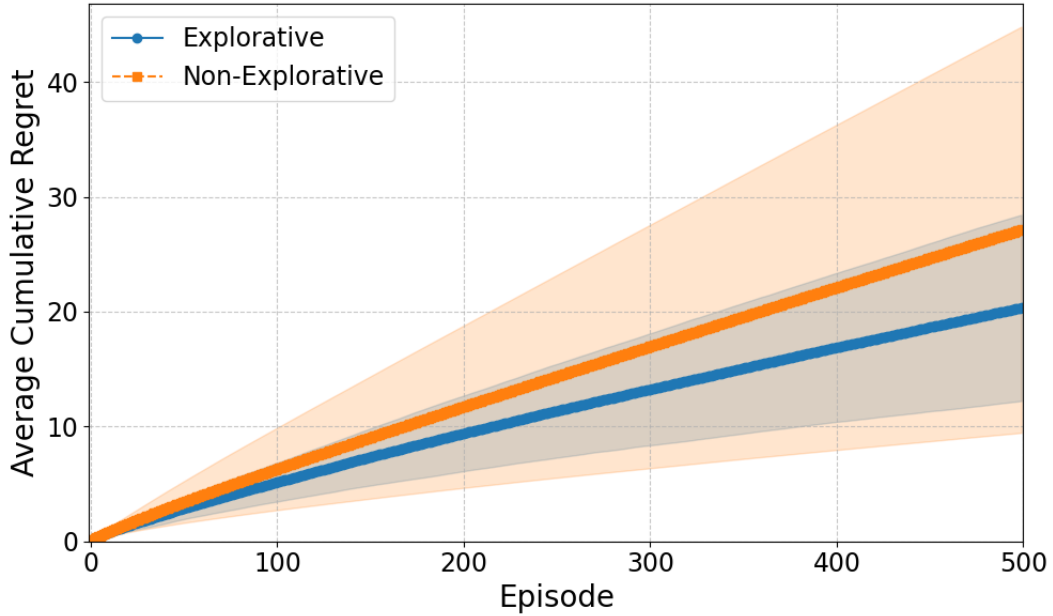


Figure 3: Average cumulative regret of LSVI (without exploration) and LSVI-UCB (with exploration), with covariate diversity over 100 MDP.

328 and 1, using parameters within this range can result in some MDPs having very small suboptimality  
 329 gaps, which makes the learning process harder. This variability is reflected in the large standard  
 330 deviation of Figure 3. Additionally, some of them may have a very small  $\lambda_0$ , resulting in weak  
 331 covariate diversity. Unfortunately, this condition is difficult to measure, as  $\lambda_0$  is defined for every  
 332 parameter  $w$  and policy  $\pi$ .

## 333 7 CONCLUSION

334 In this paper, we have proven that in the setting of Markov Decision Processes, under the assumption  
 335 of *linear Bellman completeness*, LSVI, a greedy algorithm, can achieve *sub-linear regret* if there is  
 336 sufficient *diversity* in the feature distribution, as defined by our proposed diversity condition. This  
 337 eliminates the need for explicit exploration for the agent to learn the optimal policy. Experimental  
 338 results are coherent with the theory.

339 Our results on linear function approximation pave the way for exploration-free RL in MDPs with  
 340 structure. Future work should focus on nonlinear function approximation in order to scale to com-  
 341 plex control problems. Some questions remain even in the linear realm. Is exploration-free learning  
 342 possible in  $Q^\pi$ -realizable MDPs (Weisz et al., 2023)? Since our covariate diversity condition is a  
 343 special case of UNISOFT (Papini et al., 2021a), it might be possible for LSVI to achieve *constant*  
 344 regret like LSVI-UCB under a minimum-gap assumption. From the perspective of representation  
 345 learning, our notion of feature diversity could be encouraged by some form of spectral regularization  
 346 as proposed by Tirinzoni et al. (2022) for UNISOFT. Our approach is inherently value-based, but  
 347 similar ideas could be applied to policy-based RL to reconcile theory with practical “deterministic  
 348 policy gradient” algorithms such as the one proposed by (Saleh et al., 2022). Finally, we may try  
 349 to develop a greedy-first exploration algorithm for MDPs following the example of Bastani et al.  
 350 (2021) on linear contextual bandits.

## 351 References

- 352 Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization.  
 353 In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference*  
 354 *on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of  
 355 *Proceedings of Machine Learning Research*, pp. 22–31. PMLR, 2017.
- 356 Alekh Agarwal, Nan Jiang, Sham M Kakade, and Wen Sun. Reinforcement learning: Theory and  
 357 algorithms. *CS Dept., UW Seattle, Seattle, WA, USA, Tech. Rep.*, 32:96, 2019.
- 358 Alekh Agarwal, Nan Jiang, Sham M. Kakade, and Wen Sun. *Reinforcement Learning: Theory and*  
 359 *Algorithms*. 2021a.
- 360 Alekh Agarwal, Sham M. Kakade, Jason D. Lee, and Gaurav Mahajan. On the theory of policy  
 361 gradient methods: Optimality, approximation, and distribution shift. *J. Mach. Learn. Res.*, 22:  
 362 98:1–98:76, 2021b.
- 363 Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin,  
 364 Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce-style optimization for learn-  
 365 ing from human feedback in llms. In *ACL (1)*, pp. 12248–12267. Association for Computational  
 366 Linguistics, 2024.
- 367 Eitan Altman. *Constrained Markov decision processes*. Routledge, 2021.
- 368 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul F. Christiano, John Schulman, and Dan Mané.  
 369 Concrete problems in AI safety. *CoRR*, abs/1606.06565, 2016.
- 370 Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforce-  
 371 ment learning. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pp. 263–272.  
 372 PMLR, 2017.
- 373 Hamsa Bastani, Mohsen Bayati, and Khashayar Khosravi. Mostly exploration-free algorithms for  
 374 contextual bandits. *Manag. Sci.*, 67(3):1329–1349, 2021.
- 375 Felix Berkenkamp, Matteo Turchetta, Angela P. Schoellig, and Andreas Krause. Safe model-based  
 376 reinforcement learning with stability guarantees. In *NIPS*, pp. 908–918, 2017.

- 377 Sarah Bird, Solon Barocas, Kate Crawford, Fernando Diaz, and Hanna Wallach. Exploring or  
378 exploiting? social and ethical implications of autonomous experimentation in ai. In *Workshop on*  
379 *Fairness, Accountability, and Transparency in Machine Learning*, 2016.
- 380 Steven J. Bradtke and Andrew G. Barto. Linear least-squares algorithms for temporal difference  
381 learning. *Mach. Learn.*, 22(1-3):33–57, 1996.
- 382 Lukas Brunke, Melissa Greeff, Adam W. Hall, Zhaocong Yuan, Siqi Zhou, Jacopo Panerati, and  
383 Angela P. Schoellig. Safe learning in robotics: From learning-based control to safe reinforcement  
384 learning. *Annu. Rev. Control. Robotics Auton. Syst.*, 5:411–444, 2022.
- 385 Niladri S. Chatterji, Vidya Muthukumar, and Peter L. Bartlett. OSOM: A simultaneously optimal  
386 algorithm for multi-armed and linear contextual bandits. In *AISTATS*, volume 108 of *Proceedings*  
387 *of Machine Learning Research*, pp. 1844–1854. PMLR, 2020.
- 388 Asbjørn Følstad, Theo Araujo, Effie Lai-Chong Law, Petter Bae Brandtzaeg, Symeon Papadopoulos,  
389 Lea Reis, Marcos Baez, Guy Laban, Patrick McAllister, Carolin Ischen, et al. Future directions  
390 for chatbot research: an interdisciplinary research agenda. *Computing*, 103(12):2915–2942, 2021.
- 391 Dylan J. Foster, Akshay Krishnamurthy, and Haipeng Luo. Model selection for contextual bandits.  
392 In *NeurIPS*, pp. 14714–14725, 2019.
- 393 Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in  
394 actor-critic methods. In *ICML*, volume 80 of *Proceedings of Machine Learning Research*, pp.  
395 1582–1591. PMLR, 2018.
- 396 Avishek Ghosh, Abishek Sankararaman, and Kannan Ramchandran. Problem-complexity adaptive  
397 model selection for stochastic linear bandits. In *AISTATS*, volume 130 of *Proceedings of Machine*  
398 *Learning Research*, pp. 1396–1404. PMLR, 2021.
- 399 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy  
400 maximum entropy deep reinforcement learning with a stochastic actor. In *ICML*, volume 80 of  
401 *Proceedings of Machine Learning Research*, pp. 1856–1865. PMLR, 2018.
- 402 Botao Hao, Tor Lattimore, and Csaba Szepesvári. Adaptive exploration in linear contextual bandit.  
403 In *AISTATS*, volume 108 of *Proceedings of Machine Learning Research*, pp. 3536–3545. PMLR,  
404 2020.
- 405 Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement  
406 learning. *J. Mach. Learn. Res.*, 11:1563–1600, 2010.
- 407 Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman eluder dimension: New rich classes of RL  
408 problems, and sample-efficient algorithms. In *NeurIPS*, pp. 13406–13418, 2021.
- 409 Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I. Jordan. Provably efficient reinforcement  
410 learning with linear function approximation. *Math. Oper. Res.*, 48(3):1496–1521, 2023.
- 411 Sampath Kannan, Jamie Morgenstern, Aaron Roth, Bo Waggoner, and Zhiwei Steven Wu. A  
412 smoothed analysis of the greedy algorithm for the linear contextual bandit problem. In *NeurIPS*,  
413 pp. 2231–2241, 2018.
- 414 Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances*  
415 *in applied mathematics*, 6(1):4–22, 1985.
- 416 Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- 417 Tor Lattimore, Csaba Szepesvári, and Gellért Weisz. Learning with good feature representations in  
418 bandits and in RL with a generative model. In *ICML*, volume 119 of *Proceedings of Machine*  
419 *Learning Research*, pp. 5662–5670. PMLR, 2020.

- 420 Zaile Li, Weiwei Fan, and L Jeff Hong. The (surprising) sample optimality of greedy procedures for  
421 large-scale ranking and selection. *Management Science*, 2024.
- 422 Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,  
423 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In *ICLR*  
424 (*Poster*), 2016.
- 425 Zhishuai Liu and Pan Xu. Distributionally robust off-dynamics reinforcement learning: Provable ef-  
426 ficiency with linear function approximation. In *International Conference on Artificial Intelligence*  
427 *and Statistics*, 2024.
- 428 Jihao Long and Jiequn Han. Reinforcement learning with function approximation: From linear to  
429 nonlinear. *CoRR*, abs/2302.09703, 2023.
- 430 Teodor Mihai Moldovan and Pieter Abbeel. Safe exploration in markov decision processes. In  
431 *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh,*  
432 *Scotland, UK, June 26 - July 1, 2012*. icml.cc / Omnipress, 2012.
- 433 Alessandro Montenegro, Marco Mussi, Alberto Maria Metelli, and Matteo Papini. Learning optimal  
434 deterministic policies with stochastic policy gradients. In *ICML*. OpenReview.net, 2024.
- 435 Ian Osband, Daniel Russo, and Benjamin Van Roy. (more) efficient reinforcement learning via  
436 posterior sampling. In *NIPS*, pp. 3003–3011, 2013.
- 437 Matteo Papini, Andrea Tirinzoni, Aldo Pacchiano, Marcello Restelli, Alessandro Lazaric, and Mat-  
438 teo Pirodda. Reinforcement learning in linear mdps: Constant regret and representation selection.  
439 In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wort-  
440 man Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference*  
441 *on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*,  
442 pp. 16371–16383, 2021a.
- 443 Matteo Papini, Andrea Tirinzoni, Marcello Restelli, Alessandro Lazaric, and Matteo Pirodda. Lever-  
444 aging good representations in linear contextual bandits. In Marina Meila and Tong Zhang (eds.),  
445 *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July*  
446 *2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 8371–8380.  
447 PMLR, 2021b.
- 448 Matteo Papini, Matteo Pirodda, and Marcello Restelli. Smoothing policies and safe policy gradients.  
449 *Mach. Learn.*, 111(11):4081–4137, 2022.
- 450 Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wi-  
451 ley Series in Probability and Statistics. Wiley, 1994.
- 452 Ehsan Saleh, Saba Ghaffari, Timothy Bretl, and Matthew West. Truly deterministic policy optimiza-  
453 tion. In *NeurIPS*, 2022.
- 454 David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin A. Riedmiller.  
455 Deterministic policy gradient algorithms. In *ICML*, volume 32 of *JMLR Workshop and Confer-*  
456 *ence Proceedings*, pp. 387–395. JMLR.org, 2014.
- 457 Bharat Singh, Rajesh Kumar, and Vinay Pratap Singh. Reinforcement learning in robotic applica-  
458 tions: a comprehensive survey. *Artif. Intell. Rev.*, 55(2):945–990, 2022.
- 459 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 460 William R Thompson. On the likelihood that one unknown probability exceeds another in view of  
461 the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.



- 462 Andrea Tirinzoni, Matteo Papini, Ahmed Touati, Alessandro Lazaric, and Matteo Pirodda. Scalable  
 463 representation learning in linear contextual bandits with constant regret guarantees. In *NeurIPS*,  
 464 2022.
- 465 Andrea Tirinzoni, Matteo Pirodda, and Alessandro Lazaric. On the complexity of representation  
 466 learning in contextual linear bandits. In *AISTATS*, volume 206 of *Proceedings of Machine Learning Research*, pp. 7871–7896. PMLR, 2023.
- 468 Joel A. Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of Com-*  
 469 *putational Mathematics*, 12(4):389–434, August 2011. ISSN 1615-3383. DOI: 10.1007/  
 470 s10208-011-9099-z.
- 471 Masatoshi Uehara and Wen Sun. Pessimistic model-based offline reinforcement learning under  
 472 partial coverage. In *ICLR*. OpenReview.net, 2022.
- 473 Gellért Weisz, András György, and Csaba Szepesvári. Online RL in linearly  $q^\pi$ -realizable mdps is  
 474 as easy as in linear mdps if you learn what to ignore. In *NeurIPS*, 2023.
- 475 Weiqiang Wu, Jing Yang, and Cong Shen. Stochastic linear contextual bandits with diverse contexts.  
 476 In *AISTATS*, volume 108 of *Proceedings of Machine Learning Research*, pp. 2392–2401. PMLR,  
 477 2020.
- 478 Tengyang Xie, Dylan J. Foster, Yu Bai, Nan Jiang, and Sham M. Kakade. The role of coverage in  
 479 online reinforcement learning. In *ICLR*. OpenReview.net, 2023.
- 480 Lin Yang and Mengdi Wang. Sample-optimal parametric q-learning using linearly additive features.  
 481 In *ICML*, volume 97 of *Proceedings of Machine Learning Research*, pp. 6995–7004. PMLR,  
 482 2019.
- 483 Andrea Zanette, Alessandro Lazaric, Mykel J. Kochenderfer, and Emma Brunskill. Learning near  
 484 optimal policies with low inherent bellman error. In *ICML*, volume 119 of *Proceedings of Ma-*  
 485 *chine Learning Research*, pp. 10978–10989. PMLR, 2020.

## Supplementary Materials

The following content was not necessarily subject to peer review.

### A Proof of Lemma 4.1

**Lemma 4.1.** Assume  $\lambda_{\min}(\hat{\Sigma}_{k,h}) \geq \lambda_k$  for all  $k \geq 1$  and  $h \in [H]$ . Under Assumptions 2.1 and 3.1, with probability  $1 - \delta/2$ , for all  $k \geq 1, h \in [H], s \in \mathcal{S}, a \in \mathcal{A}$ :

$$|\hat{Q}_h^k(s, a) - Q_h^*(s, a)| \leq (H - h) \left( \left( \sqrt{\beta_k(\delta)} \right) \frac{\phi_{\max}}{\sqrt{\lambda_k}} \right),$$

where

$$\sqrt{\beta_k(\delta)} := H\sqrt{A + B + C} + 1 + w_{\max},$$

and  $A := d \ln \left( 1 + \frac{\phi_{\max}^2 k}{d} \right)$ ,  $B := d \ln(w_{\max}^2 \phi_{\max}^2 k)$ ,  $C := \ln(2H\delta^{-1})$ .

We shall prove a more general version of Lemma 4.1 that takes misspecification into account.

Let the Bellman Error be defined as:

$$\mathcal{L}_h(\mathbf{w}; s, a, \mathbf{w}') := \left| \langle \phi(s, a), \mathbf{w} \rangle - r_h(s, a) - \mathbb{E}_{s' \sim \mathbb{P}_h(\cdot | s, a)} \left[ \max_{a'} \langle \phi(s', a'), \mathbf{w}' \rangle \right] \right|. \quad (3)$$

Assumption 3.6 can be rephrased as follows.

**Assumption A.1** ( $\zeta$ -Approximate Linear Bellman Completeness, Zanette et al. (2020)). For all  $h \in [H]$  and  $\mathbf{w}' \in \mathcal{W}$ , there exists a  $\mathbf{w} \in \mathcal{W}$  such that  $\sup_{s \in \mathcal{S}, a \in \mathcal{A}} \mathcal{L}_h(\mathbf{w}; s, a, \mathbf{w}') \leq \zeta$ .

We will denote:

$$\mathcal{T}_h(\mathbf{w}') := \arg \min_{\mathbf{w} \in \mathcal{W}} \sup_{s \in \mathcal{S}, a \in \mathcal{A}} \mathcal{L}_h(\mathbf{w}; s, a, \mathbf{w}'). \quad (4)$$

Clearly, by definition, for all  $s \in \mathcal{S}, a \in \mathcal{A}, h \in [H]$ :

$$\mathcal{L}_h(\mathcal{T}_h(\mathbf{w}'); s, a, \mathbf{w}') \leq \zeta. \quad (5)$$

Let  $\mathcal{Q} = \{ \langle \phi(\cdot, \cdot), \mathbf{w} \rangle \text{ s.t. } \mathbf{w} \in \mathcal{W} \}$  and  $\mathcal{V} = \{ \max_{a \in \mathcal{A}} Q(\cdot, a) \text{ s.t. } Q \in \mathcal{Q} \}$ . Let  $\hat{Q}_h^k(s, a) = \langle \phi(s, a), \hat{\mathbf{w}}_h^k \rangle$  and  $\hat{V}_h^k(s) = \max_{a \in \mathcal{A}} \langle \phi(s, a), \hat{\mathbf{w}}_h^k \rangle$ . Clearly  $\hat{Q}_h^k \in \mathcal{Q}$  and  $\hat{V}_h^k \in \mathcal{V}$  for all  $k, h$ .

**Proposition A.2** (Lemma 3 by Zanette et al. (2020)<sup>4</sup>). Fix  $h \in [H]$ . If  $\lambda = 1$ , with probability  $1 - \delta$ , for all  $V \in \mathcal{V}$ :

$$\left\| \sum_{t=1}^{k-1} \phi(s_h^t, a_h^t) \left( V(s_{h+1}^t) - \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s_h^t, a_h^t)} [V(s')] \right) \right\|_{\Sigma_{t,h}^{-1}} \leq H \sqrt{d \ln \left( 1 + \frac{\phi_{\max}^2 k}{d} \right) + d \ln(w_{\max}^2 \phi_{\max}^2 k) + \ln(\delta^{-1}) + 1}.$$

**Proposition A.3** (Lemma 8 by Zanette et al. (2020)). Let  $a_1, \dots, a_k \in \mathbb{R}^d$  and  $b_1, \dots, b_k \in \mathbb{R}$  such that  $|b_t| \leq \epsilon$  for all  $t \in [k]$ . Let  $\Sigma = \sum_{t=1}^k a_t a_t^\top + \lambda I$ . Then, for any  $k > 1$  and  $\lambda \geq 0$ :

$$\left\| \sum_{t=1}^k a_t b_t \right\|_{\Sigma^{-1}}^2 \leq k \epsilon^2. \quad (6)$$

<sup>4</sup>The extra  $H$  factor is due to the fact that we assume value functions to be in  $[0, H]$  rather than in  $[0, 1]$ .

507 **Lemma A.4.** Under Assumption A.1, if  $\lambda = 1$ , with probability  $1 - \delta$ , for all  $k \geq 1$  and  $h \in [H]$ :

$$\|\widehat{\mathbf{w}}_h^k - \mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k)\|_{\hat{\Sigma}_{k,h}} \leq \sqrt{\beta_k(\delta)} + \sqrt{k}\zeta,$$

508 where  $\sqrt{\beta_k(\delta)} := H\sqrt{d \ln\left(1 + \frac{\phi_{\max}^2 k}{d}\right) + d \ln(w_{\max}^2 \phi_{\max}^2 k) + \ln(H\delta^{-1}) + 1 + w_{\max}}$ .

509 *Proof.* First, we show that  $\|\widehat{\mathbf{w}}_h^k - \mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k)\|_{\hat{\Sigma}_{k,h}} \leq \|\widetilde{\mathbf{w}}_h^k - \mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k)\|_{\hat{\Sigma}_{k,h}}$ , that is, the projec-  
 510 tion step can only bring the estimated parameter closer to its target. Fix  $k, h$  and let  $\text{Proj}(\mathbf{w}) :=$   
 511  $\arg \max_{\mathbf{w}' \in \mathcal{W}} \|\mathbf{w} - \mathbf{w}'\|_{\hat{\Sigma}_{k,h}}$ . Since  $\mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k) \in \mathcal{W}$  by definition,  $\text{Proj}(\mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k)) = \mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k)$ .  
 512 Then:

$$\|\widehat{\mathbf{w}}_h^k - \mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k)\|_{\hat{\Sigma}_{k,h}} = \|\text{Proj}(\widetilde{\mathbf{w}}_h^k) - \text{Proj}(\mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k))\|_{\hat{\Sigma}_{k,h}} \leq \|\widetilde{\mathbf{w}}_h^k - \mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k)\|_{\hat{\Sigma}_{k,h}}, \quad (7)$$

513 where the last inequality is by contractivity of metric projections onto convex sets ( $\mathcal{W}$  is convex by  
 514 Asm. 3.1). We then proceed to upper bound  $\|\widetilde{\mathbf{w}}_h^k - \mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k)\|_{\hat{\Sigma}_{k,h}}$ . First notice that:

$$\widetilde{\mathbf{w}}_h^k = \hat{\Sigma}_{k,h}^{-1} \sum_{t=1}^{k-1} \phi(s_h^t, a_h^t) \left( r_h(s_h^t, a_h^t) + \max_{a' \in \mathcal{A}} \langle \phi(s_{h+1}^t, a'), \widehat{\mathbf{w}}_{h+1}^k \rangle \right) \quad (8)$$

$$= \hat{\Sigma}_{k,h}^{-1} \sum_{t=1}^{k-1} \phi(s_h^t, a_h^t) \left( r_h(s_h^t, a_h^t) + \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s_h^t, a_h^t)} \left[ \max_{a' \in \mathcal{A}} \langle \phi(s', a'), \widehat{\mathbf{w}}_{h+1}^k \rangle \right] \right) + (A) \quad (9)$$

$$= \hat{\Sigma}_{k,h}^{-1} \sum_{t=1}^{k-1} \phi(s_h^t, a_h^t) \langle \phi(s_h^t, a_h^t), \mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k) \rangle + (A) + (B) \quad (10)$$

$$= \mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k) + (A) + (B) + (C), \quad (11)$$

515 where

$$\begin{aligned} (A) &:= \hat{\Sigma}_{k,h}^{-1} \sum_{t=1}^{k-1} \phi(s_h^t, a_h^t) \left( \max_{a' \in \mathcal{A}} \langle \phi(s_{h+1}^t, a'), \widehat{\mathbf{w}}_{h+1}^k \rangle - \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s_h^t, a_h^t)} \left[ \max_{a' \in \mathcal{A}} \langle \phi(s', a'), \widehat{\mathbf{w}}_{h+1}^k \rangle \right] \right) \\ &= \hat{\Sigma}_{k,h}^{-1} \sum_{t=1}^{k-1} \phi(s_h^t, a_h^t) \left( \widehat{V}_{h+1}^k(s_{h+1}^t) - \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s_h^t, a_h^t)} \left[ \widehat{V}_{h+1}^k(s') \right] \right), \end{aligned} \quad (12)$$

516 and

$$(B) := \hat{\Sigma}_{k,h}^{-1} \sum_{t=1}^{k-1} \phi(s_h^t, a_h^t) \left( r_h(s_h^t, a_h^t) + \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s_h^t, a_h^t)} \left[ \max_{a' \in \mathcal{A}} \langle \phi(s', a'), \widehat{\mathbf{w}}_{h+1}^k \rangle \right] - \langle \phi(s_h^t, a_h^t), \mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k) \rangle \right), \quad (13)$$

517 and

$$(C) := \hat{\Sigma}_{k,h}^{-1} (-\lambda \mathbf{I} \mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k)). \quad (14)$$

518 By the triangular inequality,  $\|\widetilde{\mathbf{w}}_h^k - \mathcal{T}_h(\widehat{\mathbf{w}}_{h+1}^k)\|_{\hat{\Sigma}_{k,h}} \leq \|(A)\|_{\hat{\Sigma}_{k,h}} + \|(B)\|_{\hat{\Sigma}_{k,h}} + \|(C)\|_{\hat{\Sigma}_{k,h}}$ .

519 By Proposition A.2, if  $\lambda = 1$ , since  $\widehat{V}_{h+1}^k \in \mathcal{V}$ , the following holds with probability  $1 - \delta$  for all  
 520  $h \in [H]$ :

$$\|(A)\|_{\hat{\Sigma}_{k,h}} = \left\| \sum_{t=1}^{k-1} \phi(s_h^t, a_h^t) \left( \widehat{V}_{h+1}^k(s_{h+1}^t) - \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s_h^t, a_h^t)} \left[ \widehat{V}_{h+1}^k(s') \right] \right) \right\|_{\hat{\Sigma}_{k,h}^{-1}} \quad (15)$$

$$\leq H \sqrt{d \ln\left(1 + \frac{\phi_{\max}^2 k}{d}\right) + d \ln(w_{\max}^2 \phi_{\max}^2 k) + \ln(H\delta^{-1}) + 1}. \quad (16)$$

521 By Proposition A.3 and Equation (5):

$$\begin{aligned} \|(B)\|_{\hat{\Sigma}_{k,h}} &= \left\| \sum_{t=1}^{k-1} \phi(s_h^t, a_h^t) \left( r_h(s_h^t, a_h^t) + \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s_h^t, a_h^t)} \left[ \max_{a' \in \mathcal{A}} \langle \phi(s', a'), \hat{\mathbf{w}}_{h+1}^k \rangle \right] - \langle \phi(s_h^t, a_h^t), \mathcal{T}_h(\hat{\mathbf{w}}_{h+1}^k) \rangle \right) \right\|_{\hat{\Sigma}_{k,h}^{-1}} \\ &\leq \left\| \sum_{t=1}^{k-1} \phi(s_h^t, a_h^t) \mathcal{L}(\mathcal{T}_h(\hat{\mathbf{w}}_{h+1}^k); s_h^t, a_h^t, \hat{\mathbf{w}}_{h+1}^k) \right\|_{\hat{\Sigma}_{k,h}^{-1}} \end{aligned} \quad (17)$$

$$\leq \sqrt{k} \zeta. \quad (18)$$

522 Finally:

$$\|(C)\|_{\hat{\Sigma}_{k,h}} = \lambda \left\| \Sigma_{k,h}^{-1} \mathcal{T}_h(\hat{\mathbf{w}}_{h+1}^k) \right\|_{\hat{\Sigma}_{k,h}} \quad (19)$$

$$= \lambda \left\| \mathcal{T}_h(\hat{\mathbf{w}}_{h+1}^k) \right\|_{\hat{\Sigma}_{k,h}^{-1}} \quad (20)$$

$$\leq \sqrt{\lambda} \left\| \mathcal{T}_h(\hat{\mathbf{w}}_{h+1}^k) \right\| \quad (21)$$

$$\leq \sqrt{\lambda} w_{\max}. \quad (22)$$

523  $\square$

524 **Lemma A.5.** Assume  $\lambda_{\min}(\hat{\Sigma}_{k,h}) \geq \lambda_k$  for all  $k \geq 1$  and  $h \in [H]$ . Under Assumption A.1, with  
525 probability  $1 - \delta$ , for all  $k \geq 1, h \in [H], s \in \mathcal{S}, a \in \mathcal{A}$ :

$$|\hat{Q}_h^k(s, a) - Q_h^*(s, a)| \leq (H - h) \left( \left( \sqrt{\beta_k(\delta)} + \sqrt{k} \zeta \right) \frac{\phi_{\max}}{\sqrt{\lambda_k}} + \zeta \right),$$

526 where  $\sqrt{\beta_k}$  is defined in Lemma A.4.

527 *Proof.* First:

$$|\hat{Q}_h^k(s, a) - \mathcal{T}_h \hat{Q}_{h+1}^k(s, a)| = \left| \langle \phi(s, a), \hat{\mathbf{w}}_h^k \rangle - r_h(s, a) - \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s, a)} \left[ \max_{a' \in \mathcal{A}} \langle \phi(s', a'), \hat{\mathbf{w}}_{h+1}^k \rangle \right] \right| \quad (23)$$

$$= \left| \langle \phi(s, a), \hat{\mathbf{w}}_h^k - \mathcal{T}_h(\hat{\mathbf{w}}_{h+1}^k) \rangle + \langle \phi(s, a), \mathcal{T}_h(\hat{\mathbf{w}}_{h+1}^k) \rangle - r_h(s, a) - \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s, a)} \left[ \max_{a' \in \mathcal{A}} \langle \phi(s', a'), \hat{\mathbf{w}}_{h+1}^k \rangle \right] \right| \quad (24)$$

$$\leq |\langle \phi(s, a), \hat{\mathbf{w}}_h^k - \mathcal{T}_h(\hat{\mathbf{w}}_{h+1}^k) \rangle| + \zeta \quad (\text{by Equation 5}) \quad (25)$$

$$\leq \|\phi(s, a)\|_{\hat{\Sigma}_{k,h}^{-1}} \|\hat{\mathbf{w}}_h^k - \mathcal{T}_h(\hat{\mathbf{w}}_{h+1}^k)\|_{\hat{\Sigma}_{k,h}} + \zeta \quad (26)$$

$$\leq \frac{\phi_{\max}}{\sqrt{\lambda_{\min}(\hat{\Sigma}_{k,h})}} \|\hat{\mathbf{w}}_h^k - \mathcal{T}_h(\hat{\mathbf{w}}_{h+1}^k)\|_{\hat{\Sigma}_{k,h}} + \zeta \quad (27)$$

$$\leq \frac{\phi_{\max}}{\sqrt{\lambda_{\min}(\hat{\Sigma}_{k,h})}} \left( \sqrt{\beta_k(\delta)} + \sqrt{k} \zeta \right) + \zeta \quad (28)$$

$$\leq \frac{\phi_{\max}}{\sqrt{\lambda_k}} \left( \sqrt{\beta_k(\delta)} + \sqrt{k} \zeta \right) + \zeta := \varepsilon_k. \quad (29)$$

528 The rest of the proof is by (backward) induction. Note that  $Q_{H+1}^* = \hat{Q}_{H+1}^k = 0$ . Then,  $\mathcal{T}_H \hat{Q}_{H+1}^k =$   
529  $\mathcal{T}_H Q_{H+1}^*(s, a) = Q_H^*(s, a)$ . By Equation (29):

$$|\hat{Q}_H^k(s, a) - Q_H^*(s, a)| = |\hat{Q}_H^k(s, a) - \mathcal{T}_H \hat{Q}_{H+1}^k| \leq \varepsilon_k. \quad (30)$$

530 This is our base case ( $h = H$ ). The inductive hypothesis is:

$$|\widehat{Q}_{h+1}^k(s, a) - Q_{h+1}^*(s, a)| \leq (H - h - 1)\varepsilon_k. \quad (31)$$

531 Then:

$$|\widehat{Q}_h^k(s, a) - Q_h^*(s, a)| = |\widehat{Q}_h^k(s, a) - \mathcal{T}_h \widehat{Q}_{h+1}^k(s, a) + \mathcal{T}_h \widehat{Q}_{h+1}^k(s, a) - Q_h^*(s, a)| \quad (32)$$

$$\leq |\widehat{Q}_h^k(s, a) - \mathcal{T}_h \widehat{Q}_{h+1}^k(s, a)| + |\mathcal{T}_h \widehat{Q}_{h+1}^k(s, a) - Q_h^*(s, a)| \quad (33)$$

$$\leq \varepsilon_k + |\mathcal{T}_h \widehat{Q}_{h+1}^k(s, a) - Q_h^*(s, a)| \quad (34)$$

$$= \varepsilon_k + |\mathcal{T}_h \widehat{Q}_{h+1}^k(s, a) - \mathcal{T}_h Q_{h+1}^*(s, a)| \quad (35)$$

$$\leq \varepsilon_k + |\widehat{Q}_{h+1}^k(s, a) - Q_{h+1}^*(s, a)| \quad (36)$$

$$\leq \varepsilon_k + (H - h - 1)\varepsilon_k \quad (37)$$

$$= (H - h)\varepsilon_k, \quad (38)$$

532 where the inequalities are, in order: by triangular inequality, by Equation (29), by the contraction  
533 property of Bellman's operator, by the induction hypothesis.  $\square$

534 Lemma 4.1 is just a special case of Lemma A.5 when  $\zeta = 0$ .

## 535 B Proof of Lemma 4.2

536 **Lemma 4.2.** *Given Assumptions 3.1 and 3.2, the following holds for the minimum eigenvalue of the*  
537 *empirical covariance matrix for each  $h \in [H]$  and for each  $k \geq 1$ :*

$$\mathbb{P} \left[ \lambda_{\min}(\hat{\Sigma}_{k,h}) \geq \lambda + \lambda_0 k - 8\phi_{\max}^2 \sqrt{k \log(4dk/\delta)} \right] \geq 1 - \frac{\delta}{2}$$

538 *Proof.* Let  $\pi_\tau$  be the policy played by Algorithm 1 in the  $\tau$ -th episode. We can rewrite the design  
539 matrix as:

$$\hat{\Sigma}_{k,h} = \sum_{\tau=1}^k \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top + \lambda \mathbf{I} \quad (39)$$

$$= \lambda \mathbf{I} + \sum_{\tau=1}^k \mathbb{E}_{s \sim \rho_h^{\pi_\tau}} [\phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top] - \sum_{\tau=1}^k \left( \mathbb{E}_{s \sim \rho_h^{\pi_\tau}} [\phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top] - \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top \right) \quad (40)$$

540 If  $X_\tau = \mathbb{E}_{s \sim \rho_h^{\pi_\tau}} [\phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top] - \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top$ , then we have that  $\mathbb{E}_\tau[X_\tau] = 0$ . Also  
541 since  $X_\tau$  is symmetric:

$$X_\tau^2 \leq \lambda_{\max}(X_\tau^2) \mathbf{I} \leq \|X_\tau\|^2 \mathbf{I} \leq 4\phi_{\max}^4 \mathbf{I} \quad (41)$$

542 Hence from the matrix Azuma inequality by Tropp (2011), with probability  $1 - \delta_k$ , for all  $k \geq 1$ :

$$\lambda_{\max} \left( \sum_{\tau=1}^k X_\tau \right) \leq 4\phi_{\max}^2 \sqrt{2k \log d/\delta_k} \quad (42)$$

543 We set  $\delta_k = \delta/(2k^2)$  and perform a union bound over time. Finally with probability at least  $1 - \delta$   
544 for all  $k \geq 1$ :

$$\lambda_{\max} \left( \sum_{\tau=1}^k X_\tau \right) \leq 4\phi_{\max}^2 \sqrt{2k \log(4dk^2/\delta)} \leq 8\phi_{\max}^2 \sqrt{k \log(4dk/\delta)}. \quad (43)$$

545 Now let  $\pi = \pi_k$ . By definition of Algorithm 1:

$$\sum_{\tau=1}^k \mathbb{E}_{s \sim \rho_h^{\pi^\tau}} [\phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top] = \sum_{\tau=1}^k \mathbb{E}_{s \sim \rho_h^{\pi^\tau}} [\phi(s_h^\tau, \pi_h^\tau(s)) \phi(s_h^\tau, \pi_h^\tau(s))^\top] \quad (44)$$

$$= \sum_{\tau=1}^k \mathbb{E}_{s \sim \rho_h^{\pi^\tau}} [\phi(s_h^\tau, \pi_h^\tau(s)) \phi(s_h^\tau, \pi_h^\tau(s))^\top \mathbb{I}\{\langle \phi(s, \pi_h^\tau(s)), \hat{\mathbf{w}}_h^k \rangle \geq \max_{a \in \mathcal{A}} \langle \phi(s, a), \hat{\mathbf{w}}_h^k \rangle\}]. \quad (45)$$

546 So, by Assumption 3.2:

$$\lambda_{\min}(\hat{\Sigma}_{k,h}) \geq \lambda \mathbf{I} + \lambda_0 k - 8\phi_{\max}^2 \sqrt{k \log(4dk/\delta)}. \quad (46)$$

547  $\square$

## 548 C Other Proofs

549 **Proposition 3.3** (Noisy features). *Let  $|\mathcal{A}| = 2$  and  $\phi(s, a) = f(s, a) + \eta(a)$  for some function*  
 550  *$f : \mathcal{S} \times \mathcal{A} \rightarrow [0, \sqrt{2}\sigma]$  and independent Gaussian noises  $\eta(a) \sim \mathcal{N}(0, \sigma^2)$ . Then Assumption 3.2*  
 551 *holds with  $\lambda_0 \geq 0.2\sigma^2$ .*

552 *Proof.* Let  $\mathcal{A} = \{a_1, a_2\}$ ,  $w \in \mathcal{W}$ , and  $\pi$  be any deterministic policy:

$$\lambda_0 = \mathbb{E}_{s \sim \rho_h^\pi(s)} [\phi(s, \pi(s))^2 \mathbb{I}\{\phi(s, \pi(s))w \geq \max_{a \in \mathcal{A}} \{\phi(s, a)w\}\}] \quad (47)$$

$$= \mathbb{E}_{s \sim \rho_h^\pi(s)} [(f(s, a) + \eta(a))^2 \mathbb{I}\{\phi(s, \pi(s))w \geq \max_{a \in \mathcal{A}} \{\phi(s, a)w\}\}] \quad (48)$$

$$\geq \mathbb{E}_{s \sim \rho_h^\pi(s)} [\eta(a)^2 \mathbb{I}\{\phi(s, \pi(s))w \geq \max_{a \in \mathcal{A}} \{\phi(s, a)w\}\}] \quad (49)$$

$$= \sigma^2 \mathbb{P}(\phi(s, \pi(s))w \geq \max_{a \in \mathcal{A}} \{\phi(s, a)w\}) \quad (50)$$

$$= \sigma^2 \mathbb{P}((f(s, \pi(s)) + \eta(\pi(s)))w \geq \max_{a \in \mathcal{A}} \{(f(s, a) + \eta(a))w\}). \quad (51)$$

553 Fix a state  $s$  and let  $\pi(s) = a_1$  and  $w > 0$  without loss of generality. Then

$$\mathbb{P}((f(s, a_1) + \eta(a_1))w \geq \max_{a \in \mathcal{A}} \{(f(s, a) + \eta(a))w\}) \geq \mathbb{P}((f(s, a_1) + \eta(a_1))w \geq (f(s, a_2) + \eta(a_2))w) \quad (52)$$

$$= \mathbb{P}(f(s, a_1) + \eta(a_1) \geq f(s, a_2) + \eta(a_2)) \quad (53)$$

$$\geq \mathbb{P}(\eta(a_1) - \eta(a_2) > \sqrt{2}\sigma) \quad (54)$$

$$= \mathbb{P}(X > \sqrt{2}\sigma) > 0.2, \quad (55)$$

554 where  $X \sim \mathcal{N}(0, 2\sigma^2)$ .  $\square$

555 **Corollary 3.5.** *Under the same assumptions as Theorem 3.4, the expected cumulative regret of LSVI*  
 556 *is at most:*

$$\mathbb{E}[R(K)] = \mathcal{O}\left(H^3 \sqrt{\frac{dK}{\lambda_0}} \log(K)\right).$$

557 *Proof.* We define an event  $P$ , under which the cumulative regret is bounded as described in Theorem  
 558 3.4, that occurs with probability  $1 - \delta$ :

$$\begin{aligned}\mathbb{E}[R(K)] &= \mathbb{E}[R(K)\mathbb{I}(P)] + \mathbb{E}[R(K)\mathbb{I}(\bar{P})] \\ &\leq \mathcal{O}\left(H^3\sqrt{\frac{dK}{\lambda_0}}\log(K/\delta)\right)(1 - \delta) + 2HK\delta \\ &\leq \mathcal{O}\left(H^3\sqrt{\frac{dK}{\lambda_0}}\log(K^{3/2})\right) + 2H\sqrt{K} \\ &= \mathcal{O}\left(H^3\sqrt{\frac{dK}{\lambda_0}}\log(K)\right)\end{aligned}$$

559 The first inequality is by Theorem 3.4, and by upper bounding the regret with the trivial  $2HK$  when  
 560  $P$  does not hold. The second inequality follows by  $1 - \delta \leq 1$  and by setting  $\delta = 1/\sqrt{K}$ .  $\square$

561 **Theorem 3.7.** *If Assumptions 3.6 and 3.2 are satisfied, with probability  $1 - \delta$ , the cumulative regret*  
 562 *of LSVI is at most:*

$$R(K) = \tilde{\mathcal{O}}\left(H^3\sqrt{\frac{dK}{\lambda_0}} + H^2\zeta\frac{K}{\sqrt{\lambda_0}} + H^2\zeta K\right).$$

563 *Proof.* This is a simple variant of Theorem 3.4 that uses the more general Lemma A.5 in place of  
 564 Lemma 4.1. Logarithmic terms are omitted for brevity.  $\square$

## 565 D Code and Computing Infrastructure

566 The experimental results presented in Section 6 were obtained using Python. The code, a modified  
 567 version of the official implementation from Liu & Xu (2024), is attached to the paper and was  
 568 executed on [Kaggle.com](https://www.kaggle.com), a cloud computing platform.