# Instance-Dependent Near-Optimal Policy Identification in Linear MDPs via Online Experiment Design

Andrew Wagenmaker\*

Kevin Jamieson<sup>†</sup>

July 7, 2022

#### Abstract

While much progress has been made in understanding the minimax sample complexity of reinforcement learning (RL)—the complexity of learning on the "worst-case" instance—such measures of complexity often do not capture the true difficulty of learning. In practice, on an "easy" instance, we might hope to achieve a complexity far better than that achievable on the worst-case instance. In this work we seek to understand the "instance-dependent" complexity of learning near-optimal policies (PAC RL) in the setting of RL with linear function approximation. We propose an algorithm, PEDEL, which achieves a fine-grained instance-dependent measure of complexity, the first of its kind in the RL with function approximation setting, thereby capturing the difficulty of learning on each particular problem instance. Through an explicit example, we show that PEDEL yields provable gains over low-regret, minimax-optimal algorithms and that such algorithms are unable to hit the instance-optimal rate. Our approach relies on a novel online experiment design-based procedure which focuses the exploration budget on the "directions" most relevant to learning a near-optimal policy, and may be of independent interest.

# 1 Introduction

In the PAC (Probably Approximately Correct) reinforcement learning (RL) setting, an agent is tasked with exploring an unknown environment in order to learn a policy which maximizes the amount of reward collected. In general, we are interested in learning such a policy using as few interactions with the environment (as small sample complexity) as possible. We might hope that the number of samples needed would scale with the difficulty of identifying a near-optimal policy in our particular environment. For example, in a "hard" environment, we would expect that more samples might be required, while in an "easy" environment, fewer samples may be needed.

The RL community has tended to focus on developing algorithms which have near-optimal worst-case sample complexity—sample complexities that are only guaranteed to be optimal on "hard" instances. Such algorithms typically have complexities which scale, for example, as  $\mathcal{O}(\text{poly}(d,H)/\epsilon^2)$ , for d the dimensionality of the environment, H the horizon, and  $\epsilon$  the desired level of optimality. While we may be able to show this complexity is optimal on a hard instance, it is unable to distinguish between "hard" and "easy" problems. The scaling is identical for two environments as long as the dimensionality and horizon of each are the same—no consideration is given to the actual difficulty of the problem—and we therefore have no guarantee that our algorithm is solving the problem with complexity scaling as the actual difficulty. Indeed, as recent work has shown

<sup>\*</sup>University of Washington, Seattle. Email: ajwagen@cs.washington.edu

<sup>&</sup>lt;sup>†</sup>University of Washington, Seattle. Email: jamieson@cs.washington.edu

(Wagenmaker et al., 2021b), this is not simply an analysis issue: worst-case optimal algorithms can be very suboptimal on "easy" instances.

Towards developing algorithms which overcome this, we might instead consider the *instance-dependent* difficulty—the hardness of solving a particular problem instance—and hope to obtain a sample complexity scaling with this instance-dependent difficulty, thereby guaranteeing that we solve "easy" problems using only a small number of samples, but still obtain the worst-case optimal rate on "hard" problems. While progress has been made in understanding the instance-dependent complexity of learning in RL, the results are largely limited to environments with a finite number of states and actions. In practice, real-world RL problems often involve large (even infinite) state-spaces and, in order to solve such problems, we must generalize across states. To handle such settings, the RL community has turned to function approximation-based methods, which allow for provable learning in large state-space environments. However, while worst-case optimal results have been shown, little is understood on the instance-dependent complexity of learning in these settings.

In this work we aim to bridge this gap. We consider, in particular, the linear MDP setting, and develop an algorithm which provably learns a near-optimal policy with sample complexity scaling as the difficulty of each individual instance. Furthermore, by comparing to our instance-dependent measure of complexity, we show that low-regret algorithms are provably suboptimal for PAC RL in function approximation settings. Our algorithm relies on a novel online experiment design-based procedure—adapting classical techniques from linear experiment design to settings where *navigation* is required to measure a particular covariate—which may be of independent interest.

# 1.1 Contributions

Our contributions are as follows:

• We propose an algorithm, PEDEL, which learns an  $\epsilon$ -optimal policy with instance-dependent sample complexity scaling as (up to H factors):

$$\sum_{h=1}^{H} \inf_{\pi_{\text{exp}}} \max_{\pi \in \Pi} \frac{\|\phi_{\pi,h}\|_{\mathbf{\Lambda}_{\pi_{\text{exp}},h}^{-1}}^2}{(V_0^{\star} - V_0^{\pi})^2 \vee \epsilon^2} \cdot \left(d + \log \frac{1}{\delta}\right)$$

for  $\phi_{\pi,h}$  the "average feature vector" of policy  $\pi$ ,  $\Lambda_{\pi_{\exp},h}$  the expected covariance of the policy  $\pi_{\exp}$ , and  $V_0^{\star} - V_0^{\pi}$  the "policy gap". We show that PEDEL also has worst-case optimal dimension-dependence—its sample complexity never exceeds  $\widetilde{\mathcal{O}}(d^2H^7/\epsilon^2)$ —but that on "easy" instances it achieves complexity much smaller than the worst-case optimal rate.

- It is well-known that low-regret algorithms achieve the worst-case optimal rate for PAC RL. We construct an explicit example, however, where the instance-dependent complexity of PEDEL improves on the complexity of any low-regret algorithm by a factor of the dimensionality, providing the first evidence that low-regret algorithms are provably suboptimal on "easy" instances for PAC RL in function approximation settings.
- We develop a general experiment design-based approach to exploration in MDPs, which allows us to focus our exploration in the directions most relevant to learning near-optimal policies. Our approach is based on the key observation that, while solving an experiment design in an MDP would require knowledge of the MDP dynamics, we can approximately solve one without knowledge of the dynamics by running a regret minimization algorithm on a carefully chosen reward function, inducing the correct exploration. We apply our experiment design approach to

efficiently explore our MDP so as to identify near-optimal policies, but show that it can also be used to collect observations minimizing much more general experiment design objective functions.

# 2 Related Work

The sample complexity of RL has been studied for decades (Kearns & Singh, 1998; Brafman & Tennenholtz, 2002; Kakade, 2003). The two primary problems considered are the regret minimization problem (where the goal is to obtain large online reward) and the PAC policy identification problem (where the goal is to find a near-optimal policy using as few samples as possible), which is the focus of this work. In the tabular RL setting, the question of obtaining worst-case optimal algorithms is nearly closed (Dann & Brunskill, 2015; Dann et al., 2019; Ménard et al., 2020; Zhang et al., 2020). As such, in this section we focus primarily on results in the RL with function approximation literature, as well as results on instance-dependent RL.

Sample-Efficient RL with Linear Function Approximation. To generalize beyond MDPs with a finite number of states and acions, the RL community has considered function approximation, replacing the tabular model with more powerful settings that allow for generalization across states. Such settings have been considered in classical works (Baird, 1995; Bradtke & Barto, 1996; Sutton et al., 1999; Melo & Ribeiro, 2007), yet these works do not provide polynomial sample complexities. More recently, there has been intense interest in obtaining polynomial complexities for general function classes (Jiang et al., 2017; Du et al., 2021; Jin et al., 2021; Foster et al., 2021), and, in particular, linear function classes (Yang & Wang, 2019; Jin et al., 2020; Wang et al., 2019; Du et al., 2019; Zanette et al., 2020a,b; Ayoub et al., 2020; Jia et al., 2020; Weisz et al., 2021; Zhou et al., 2020, 2021; Zhang et al., 2021; Wang et al., 2021).

In the linear MDP setting, the state-of-the-art in PAC RL is the work of Wagenmaker et al. (2022), which proposes a computationally efficient algorithm achieving a complexity of  $\mathcal{O}(d^2H^5/\epsilon^2)$  for the more general reward-free RL problem, and shows a matching lower bound of  $\Omega(d^2H^2/\epsilon^2)$  for the PAC RL problem. While this result obtains tight dimension-dependence, it is still worst-case, and offers no insight on the instance-dependent complexity. Other works of note in this category are (Jin et al., 2020; Zanette et al., 2020b; Zhou et al., 2020), which establish regret guarantees in the setting of linear MDPs and the related setting of linear mixture MDPs. Jin et al. (2020) and Zanette et al. (2020b) obtain regret guarantees of  $\mathcal{O}(\sqrt{d^3H^4K})$  and  $\mathcal{O}(\sqrt{d^2H^4K})$ , respectively, though the approach of Zanette et al. (2020b) is computationally inefficient. Via an online-to-batch conversion (Jin et al., 2018), these algorithms achieve PAC complexities of  $\mathcal{O}(d^3H^4/\epsilon^2)$  and  $\mathcal{O}(d^2H^4/\epsilon^2)$ . In the setting of linear mixture MDPs, Zhou et al. (2020) show a regret bound of  $\mathcal{O}(\sqrt{d^2H^3K})$  and a matching lower bound, yielding the first provably tight and computationally efficient algorithms for RL with function approximation.

Instance-Dependent RL. Much of the recent work on instance-dependent RL has focused on the tabular setting. Ok et al. (2018) provide an algorithm which achieves asymptotically optimal instance-dependent regret, yet it is computationally inefficient. Simchowitz & Jamieson (2019) show that standard optimistic algorithms achieve regret bounded as  $\mathcal{O}(\sum_{s,a,h} \frac{\log K}{\Delta_h(s,a)})$ , for  $\Delta_h(s,a)$  the value-function gap, a result later refined by (Xu et al., 2021; Dann et al., 2021). Obtaining instance-dependent guarantees for policy identification has proved more difficult, yet a variety of results do exist (Zanette et al., 2019; Jonsson et al., 2020; Marjani & Proutiere, 2020; Marjani et al., 2021). In the tabular setting, the most comparable work to ours is that of Wagenmaker

et al. (2021b), which propose a refined instance-dependent measure of complexity, the gap-visitation complexity, and show that it is possible to learn an  $\epsilon$ -optimal policy with complexity scaling as the gap-visitation complexity. While the gap-visitation is shown to be tight in certain settings, no general lower-bounds exist. Towards obtaining sharp guarantees, Tirinzoni et al. (2022) show that in the simpler setting of deterministic MDPs, a quantity similar in spirit to the gap-visitation complexity is tight, providing matching upper and lower bounds.

In the setting of RL with function approximation, to our knowledge, only two existing works obtain guarantees that would be considered "instance-dependent". Wagenmaker et al. (2021a) show a "first-order" regret bound of  $\mathcal{O}(\sqrt{d^3H^3V_0^*K})$ , where  $V_0^*$  is the value of the *optimal* policy on the particular MDP under consideration. He et al. (2020) show that standard optimistic algorithms achieve regret guarantees of  $\mathcal{O}(\frac{d^3H^5\log K}{\Delta_{\min}})$  and  $\mathcal{O}(\frac{d^2H^5\log^3 K}{\Delta_{\min}})$  in the settings of linear MDPs and linear mixture MDPs, respectively, for  $\Delta_{\min}$  the minimum value-function gap. While both these works do obtain instance-dependent results, the instance-dependence is rather coarse, depending on only a single parameter  $(V_0^*$  or  $\Delta_{\min})$ —our goal will instead be to obtain more refined instance-dependent guarantees.

**Experiment Design in Sequential Environments.** Experiment design is a well-developed subfield of statistics, and a full survey is beyond the scope of this work (see Pukelsheim (2006) for an overview). We highlight several works on experiment design in sequential environments that are particularly relevant. First, the work of Fiez et al. (2019) achieves the instance-optimal rate for best-arm identification in linear bandits and relies on an adaptive experiment design-based. Their approach, as well as the related work of Soare et al. (2014), provides inspiration for our algorithm—in some sense PEDEL can be seen as a generalization of the RAGE algorithm to problems with horizon greater than 1. Second, the work of Wagenmaker et al. (2021c) provides an experiment design-based algorithm in the setting of linear dynamical systems, and show that it hits the optimal instance-dependent rate for learning in such systems. While their results are somewhat more general, they specialize to the problem of identifying a near-optimal controller for the LQR problem—thereby solving the PAC RL problem optimally in the special case of quadratic losses and linear dynamical systems. It is not clear, however, if their approach generalizes beyond linear dynamical systems. Finally, while the current work was in preparation, Mutny et al. (2022) proposed an approach to solving experiment design problems in MDPs. To our knowledge, this is the only existing work that directly considers the problem of experiment design in MDPs. However, they make the simplifying assumption that the transition dynamics are known, which essentially reduces their problem to a computational one—in contrast, our approach handles the much more difficult setting of unknown dynamics, and shows that efficient experiment design is possible even in this more difficult setting.

# 3 Preliminaries

We let  $\|\phi\|_{\mathbf{A}}^2 = \phi^{\top} \mathbf{A} \phi$ ,  $\|\cdot\|_{\text{op}}$  denote the matrix operator norm (matrix 2-norm), and  $\|\cdot\|_{\text{F}}$  denote the Frobenius norm. Given some norm  $\|\cdot\|$ ,  $\|\cdot\|_{*}$  denotes the dual norm.  $\mathbb{S}^{d}_{+}$  denotes the set of PSD matrices in  $\mathbb{R}^{d \times d}$ .  $\widetilde{\mathcal{O}}(\cdot)$  hides absolute constants and log factors of the arguments.  $\lesssim$  denotes inequality up to constants.  $\mathbb{E}_{\pi}$  and  $\mathbb{P}_{\pi}$  denote the expectation and probability measure induced by playing some policy  $\pi$  in our MDP. We let  $\phi_{h,\tau} := \phi(s_{h,\tau}, a_{h,\tau})$  denote the feature vector encountered at step h of episode  $\tau$  (and similarly define  $r_{h,\tau}$ ).

Markov Decision Processes. In this work, we study episodic, finite-horizon, time inhomogeneous Markov Decision Processes (MDPs), denoted by a tuple,  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, \{P_h\}_{h=1}^H, \{\nu_h\}_{h=1}^H)$ . We let  $\mathcal{S}$  denote the state space,  $\mathcal{A}$  the action space, H the horizon,  $\{P_h\}_{h=1}^H$  the transition kernel, and  $\{\nu_h\}_{h=1}^H$  the reward distribution, where  $P_h(\cdot|s,a) \in \Delta_{\mathcal{S}}$  denotes the distribution over the next state when playing action a in state s at step h, and  $\nu_h(s,a) \in \Delta_{[0,1]}$  denotes the corresponding distribution over reward. We overload notation and let  $\nu_h(s,a)$  also refer to the expected reward. We assume that every episode starts in state  $s_1$ , and that  $\{P_h\}_{h=1}^H$  and  $\{\nu_h\}_{h=1}^H$  are initially unknown.

Let  $\pi = \{\pi_h\}_{h=1}^H$  denote a policy mapping states to distributions over actions,  $\pi_h : \mathcal{S} \to \triangle_{\mathcal{A}}$ . When  $\pi$  is deterministic, we let  $\pi_h(s)$  denote the action policy  $\pi$  takes at (s,h). An episode begins at state  $s_1$ . The agent takes action  $a_1 \sim \pi_1(s_1)$ , transitions to state  $s_2 \sim P_1(\cdot|s_1,a_1)$ , and receives reward  $r_1(s_1,a_1) \sim \nu_1(s_1,a_1)$ . In  $s_2$ , the agent chooses a new action  $a_2 \sim \pi_2(s_2)$ , and the process repeats. After H steps, the episode terminates, and the agent restarts at  $s_1$ .

In general, we are interested in learning policies that collect a large amount of reward. We can quantify the performance of a policy in terms of the value function. In particular, the Q-value function,  $Q_h^{\pi}(s,a)$ , denotes the expected reward that will be obtained if we are in state s at step h, play action a, and then play policy  $\pi$  for the remainder of the episode. Formally,  $Q_h^{\pi}(s,a) := \mathbb{E}_{\pi}[\sum_{h'=h}^{H} r_{h'}(s_{h'},a_{h'})|s_h = s, a_h = a]$ . The value function is similarly defined as  $V_h^{\pi}(s) := \mathbb{E}_{\pi}[\sum_{h'=h}^{H} r_{h'}(s_{h'},a_{h'})|s_h = s]$ . For deterministic policies,  $V_h^{\pi}(s) = Q_h^{\pi}(s,\pi_h(s))$ . We denote the optimal Q-value function by  $Q_h^{\star}(s,a) = \sup_{\pi} Q_h^{\pi}(s,a)$  and the optimal value function by  $V_h^{\star}(s) = \sup_{\pi} V_h^{\pi}(s)$ , where the suprema is taken over all policies, both deterministic and stochastic. We define the value of a policy as  $V_0^{\pi} = V_1^{\pi}(s_1)$ —the expected reward policy  $\pi$  achieves over an entire episode—and say a policy  $\pi$  is optimal if  $V_0^{\pi} = V_0^{\star}$ . For some set of policies  $\Pi$  (which may not contain an optimal policy), we let  $V_0^{\star}(\Pi) := \sup_{\pi \in \Pi} V_0^{\pi}$ .

**PAC Reinforcement Learning.** In PAC RL, the goal is to identify some policy  $\widehat{\pi}$  using as few episodes as possible, such that, with probability at least  $1 - \delta$ ,

$$V_0^{\star} - V_0^{\widehat{\pi}} \le \epsilon.$$

We say that such a policy is  $\epsilon$ -optimal, and an algorithm with such a guarantee on every environment and reward function is  $(\epsilon, \delta)$ -PAC. We will also refer to this problem as "policy identification".

# 3.1 Linear MDPs

In this work, we are interested in the setting where the state space could be infinite, and the learner must generalize across states. In particular, we consider the linear MDP model defined as follows.

**Definition 3.1** (Linear MDPs (Jin et al., 2020)). We say that an MDP is a *d-dimensional linear MDP*, if there exists some (known) feature map  $\phi(s,a): \mathcal{S} \times \mathcal{A} \to \mathbb{R}^d$ , H (unknown) signed vector-valued measures  $\boldsymbol{\mu}_h \in \mathbb{R}^d$  over  $\mathcal{S}$ , and H (unknown) reward vectors  $\boldsymbol{\theta}_h \in \mathbb{R}^d$ , such that:

$$P_h(\cdot|s,a) = \langle \phi(s,a), \mu_h(\cdot) \rangle, \quad \mathbb{E}[\nu_h(s,a)] = \langle \phi(s,a), \theta_h \rangle.$$

We will assume  $\|\phi(s, a)\|_2 \le 1$  for all s, a; and for all h,  $\||\mu_h|(\mathcal{S})\|_2 = \|\int_{s \in \mathcal{S}} |\mathrm{d}\mu_h(s)|\|_2 \le \sqrt{d}$  and  $\|\theta_h\|_2 \le \sqrt{d}$ .

Linear MDPs encompass, for example, tabular MDPs, but can also model more complex settings, such as feature spaces corresponding to the simplex (Jin et al., 2020), or the linear bandit problem.

Critically, linear MDPs allow for infinite state-spaces, as well as generalization across states—rather than learning the behavior in particular states, we can learn in the d-dimensional ambient space. Note that the standard definition of linear MDPs, for example as given in Jin et al. (2020), assumes the rewards are deterministic, while we assume the rewards are random but that their means are linear. We still assume, however, that the random rewards,  $r_h(s, a)$ , are contained in [0, 1] almost surely.

For a given policy  $\pi$ , we define the feature-visitation at step h, the expected feature vector policy  $\pi$  encounters at step h, as  $\phi_{\pi,h} := \mathbb{E}_{\pi}[\phi(s_h, a_h)]$ . Note that this is a direct generalization of state-visitations in tabular RL—if our MDP is in fact tabular,  $[\phi_{\pi,h}]_{s,a} = \mathbb{P}_{\pi}[s_h = s, a_h = a]$ , so the feature visitation vector corresponds directly to the state visitations. Note also that we can write the value of a policy as  $V_0^{\pi} = \sum_{h=1}^{H} \langle \phi_{\pi,h}, \boldsymbol{\theta}_h \rangle$ . Denote the average feature vector induced by  $\pi$  in a particular state s as  $\phi_{\pi,h}(s) = \mathbb{E}_{a \sim \pi_h(\cdot|s)}[\phi(s,a)]$ . We also define  $\Lambda_{\pi,h} := \mathbb{E}_{\pi}[\phi(s_h, a_h)\phi(s_h, a_h)^{\top}]$ , the expected covariance of policy  $\pi$  at step h,  $\lambda_{\min,h}^{\star} = \sup_{\pi} \lambda_{\min}(\Lambda_{\pi,h})$  the largest achievable minimum eigenvalue at step h, and  $\lambda_{\min,h}^{\star} = \min_{h} \lambda_{\min,h}^{\star}$ . We will make the following assumption.

**Assumption 1** (Full Rank Covariates). In our MDP,  $\lambda_{\min}^{\star} > 0$ .

We remark that Assumption 1 is analogous to other explorability assumptions found in the RL with function approximation literature (Zanette et al., 2020c; Hao et al., 2021; Agarwal et al., 2021).

To reduce uncertainty in directions of interest, we will be interested in optimizing over the set of all realizable covariance matrices on our particular MDP. To this end, define

$$\mathbf{\Omega}_h := \{ \mathbb{E}_{\pi \sim \omega} [\mathbf{\Lambda}_{\pi,h}] : \omega \in \mathbf{\Omega}_{\pi} \}$$
(3.1)

for  $\Omega_{\pi}$  the set of all valid distributions over Markovian policies (both deterministic and stochastic).  $\Omega_h$  is, then, the set of all covariance matrices realizable by distributions over policies at step h.

# 4 Near-Optimal Policy Identification in Linear MDPs

We are now ready to state our algorithm, PEDEL.

**Pedel Description.** PEDEL is a *policy-elimination*-style algorithm. It takes as input some set of policies,  $\Pi$ , and proceeds in epochs, maintaining a set of *active* policies,  $\Pi_{\ell}$ , such that all  $\pi \in \Pi_{\ell}$  are guaranteed to satisfy  $V_0^{\pi} \geq V_0^{\star}(\Pi) - 4\epsilon_{\ell}$ , for  $\epsilon_{\ell} = 2^{-\ell}$ . After running for  $\lceil \log \frac{4}{\epsilon} \rceil$  epochs, it returns any of the remaining active policies, which will be guaranteed to have value at least  $V_0^{\star}(\Pi) - \epsilon$ .

In order to ensure  $\Pi_{\ell}$  only contains  $4\epsilon_{\ell}$ -optimal policies, sufficient exploration must be performed at every epoch to refine the estimate of each policy's value. While works such as Wagenmaker et al. (2022) have demonstrated how to efficiently traverse a linear MDP and collect the necessary observations, existing exploration procedures are unable to obtain the instance-dependent complexity we desire. To overcome this, PEDEL relies on a novel online experiment design procedure to ensure exploration is focused only on the directions necessary to evaluate the current set of active policies.

In particular, one can show that, if we have collected some covariates  $\Lambda_{h,\ell}$ , the uncertainty in our estimate of the value of policy  $\pi$  at step h scales as  $\|\hat{\phi}_{\pi,h}^{\ell}\|_{\Lambda_{h,\ell}^{-1}}$ , for  $\hat{\phi}_{\pi,h}^{\ell}$  the estimated feature-visitation for policy  $\pi$  at epoch  $\ell$ . To reduce our uncertainty at each round, we would therefore like to collect covariates such that  $\|\hat{\phi}_{\pi,h}^{\ell}\|_{\Lambda_{h,\ell}^{-1}} \lesssim \epsilon_{\ell}$ . Collecting covariates which satisfy this using the minimum number of episodes of exploration possible involves solving the experiment

## Algorithm 1 Policy Learning via Experiment Design in Linear MDPs (Pedel)

```
1: input: tolerance \epsilon, confidence \delta, policy set \Pi
 2: \ell_0 \leftarrow \lceil \log_2 \frac{d^{3/2}}{H} \rceil, \Pi_{\ell_0} \leftarrow \Pi, \widehat{\phi}_{\pi,1}^1 \leftarrow \mathbb{E}_{a \sim \pi_1(\cdot | s_1)}[\phi(s_1, a)], \forall \pi \in \Pi
 3: for \ell = \ell_0, \ell_0 + 1, \dots, \lceil \log \frac{4}{\epsilon} \rceil do
                 \epsilon_{\ell} \leftarrow 2^{-\ell}, \, \beta_{\ell} \leftarrow 64H^4 \log \frac{4H^2|\Pi_{\ell}|\ell^2}{\delta}
 5:
                         Solve (4.1) by running Algorithm 2, collect data \{(\phi_{h,\tau}, r_{h,\tau}, s_{h+1,\tau})\}_{\tau=1}^{K_{h,\ell}} such that:
 6:
                                                \max_{\pi \in \Pi_{\delta}} \|\widehat{\phi}_{\pi,h}^{\ell}\|_{\mathbf{\Lambda}_{h,\ell}^{-1}}^{2} \leq \epsilon_{\ell}^{2}/\beta_{\ell} \quad \text{for} \quad \mathbf{\Lambda}_{h,\ell} \leftarrow \sum_{\tau=1}^{K_{h,\ell}} \phi_{h,\tau} \phi_{h,\tau}^{\top} + 1/d \cdot I
                                                                                          // Estimate feature-visitations for active policies
 7:
                                  \widehat{\boldsymbol{\phi}}_{\pi,h+1}^{\ell} \leftarrow \left( \sum_{\tau=1}^{K_{h,\ell}} \boldsymbol{\phi}_{\pi,h+1}(s_{h+1,\tau}) \boldsymbol{\phi}_{h,\tau}^{\mathsf{T}} \boldsymbol{\Lambda}_{h,\ell}^{-1} \right) \widehat{\boldsymbol{\phi}}_{\pi,h}^{\ell}
 8:
                         \widehat{m{	heta}}_h^\ell \leftarrow m{\Lambda}_{h.\ell}^{-1} \sum_{	au=1}^{K_{h,\ell}} m{\phi}_{h,	au} r_{h,	au}
                                                                                                                                                                       // Estimate reward vectors
 9:
                 // Remove provably suboptimal policies from active policy set
10:
                               \Pi_{\ell+1} \leftarrow \Pi_{\ell} \setminus \left\{ \pi \in \Pi_{\ell} : \widehat{V}_0^{\pi} < \sup_{\pi' \in \Pi_{\ell}} \widehat{V}_0^{\pi'} - 2\epsilon_{\ell} \right\} \quad \text{for} \quad \widehat{V}_0^{\pi} := \sum_{h=1}^{H} \langle \widehat{\boldsymbol{\phi}}_{\pi,h}^{\ell}, \widehat{\boldsymbol{\theta}}_{h}^{\ell} \rangle
                 if |\Pi_{\ell+1}| = 1 then return \pi \in \Pi_{\ell+1}
11:
12: return any \pi \in \Pi_{\ell+1}
```

design:

$$\inf_{\mathbf{\Lambda}_{\exp} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi_\ell} \| \widehat{\boldsymbol{\phi}}_{\pi,h}^{\ell} \|_{\mathbf{\Lambda}_{\exp}^{-1}}^2. \tag{4.1}$$

Note that this design has the form of an XY-experiment design (Soare et al., 2014). Solving (4.1) will produce covariance  $\Lambda_{\rm exp}$  which reduces uncertainty in relevant feature directions. However, to solve this design we require knowledge of which covariance matrices are realizable on our particular MDP. In general we do not know the MDP's dynamics, and therefore do not have access to this knowledge. To overcome this and solve (4.1), in Section 5 we provide an algorithm, Algorithm 2, that is able to solve (4.1) in an online manner without knowledge of the MDP dynamics by running a low-regret algorithm on a carefully chosen reward function.

Estimating Feature-Visitations. We remark briefly on the estimation of the feature-visitations on Line 8. If we assume that  $\{\phi_{h,\tau}\}_{\tau=1}^{K_{h,\ell}}$  is fixed and that all randomness is due to  $s_{h+1,\tau}$ , then it is easy to see that, using the structure present in linear MDPs as given in Definition 3.1,

$$\mathbb{E}\left[\sum_{\tau=1}^{K_{h,\ell}} \boldsymbol{\phi}_{\pi,h+1}(s_{h+1,\tau}) \boldsymbol{\phi}_{h,\tau}^{\top} \boldsymbol{\Lambda}_{h,\ell}^{-1}\right] = \sum_{\tau=1}^{K_{h,\ell}} \left(\int \boldsymbol{\phi}_{\pi,h+1}(s) d\boldsymbol{\mu}_{h}(s)^{\top} \boldsymbol{\phi}_{h,\tau}\right) \boldsymbol{\phi}_{h,\tau}^{\top} \boldsymbol{\Lambda}_{h,\ell}^{-1}$$
$$= \int \boldsymbol{\phi}_{\pi,h+1}(s) d\boldsymbol{\mu}_{h}(s)^{\top}.$$

By Definition 3.1, we have  $\phi_{\pi,h+1} = (\int \phi_{\pi,h+1}(s) d\mu_h(s)^{\top}) \phi_{\pi,h}$ . Comparing these, we see that our estimator of  $\phi_{\pi,h+1}$  on Line 8 is (conditioned on  $\{\phi_{h,\tau}\}_{\tau=1}^{K_{h,\ell}}$ ) unbiased, assuming  $\hat{\phi}_{\pi,h}^{\ell} \approx \phi_{\pi,h}$ .

#### 4.1 Main Results

We have the following result on the performance of Pedel.

**Theorem 1.** Consider running PEDEL with some set of Markovian policies  $\Pi$  on any linear MDP satisfying Definition 3.1 and Assumption 1. Then with probability at least  $1 - \delta$ , PEDEL outputs a policy  $\widehat{\pi} \in \Pi$  such that  $V_0^{\widehat{\pi}} \geq V_0^*(\Pi) - \epsilon$ , and runs for at most

$$C_0 H^4 \cdot \sum_{h=1}^{H} \inf_{\boldsymbol{\Lambda}_{\text{exp}} \in \boldsymbol{\Omega}_h} \max_{\pi \in \Pi} \frac{\|\boldsymbol{\phi}_{\pi,h}\|_{\boldsymbol{\Lambda}_{\text{exp}}^{-1}}^2}{\max\{V_0^{\star}(\Pi) - V_0^{\pi}, \Delta_{\min}^{\Pi}, \epsilon\}^2} \cdot \left(\log|\Pi| + \log\frac{1}{\delta}\right) + C_1$$

episodes, with  $C_0 = \log \frac{1}{\epsilon} \cdot \operatorname{poly} \log(H, \log \frac{1}{\epsilon})$ ,  $C_1 = \operatorname{poly} \left(d, H, \frac{1}{\lambda_{\min}^*}, \log \frac{1}{\delta}, \log \frac{1}{\epsilon}, \log |\Pi|\right)$ ,  $\Delta_{\min}^{\Pi} := V_0^{\star}(\Pi) - \max_{\pi \in \Pi: V_0^{\pi} < V_0^{\star}(\Pi)} V_0^{\pi}$ , and  $\Omega_h$  the set of covariance matrices realizable on our MDP, as defined in (3.1).

The proof of Theorem 1 is given in Appendix B. Theorem 1 quantifies, in a precise instance-dependent way, the complexity of identifying a policy  $\widehat{\pi}$  with value at most a factor of  $\epsilon$  from the value of the optimal policy in  $\Pi$ . In particular, it trades off between the difficulty of showing a policy  $\pi$  is suboptimal—the "policy gap",  $V_0^{\star}(\Pi) - V_0^{\pi}$ —and the difficulty of exploring in the direction necessary to reduce the uncertainty on policy  $\pi$ ,  $\|\phi_{\pi,h}\|_{\Lambda_{\exp}^{-1}}$ . Rather than scaling with factors such as d and  $\epsilon$ , our complexity measure scales with instance-dependent quantities—the covariance matrices we can obtain on our particular MDP, the feature vectors we expect to observe on our MDP, and the policy gaps on our MDP.

Theorem 1 holds for an arbitrary set of policies, yet, in general, we are interested in learning a policy which has value within a factor of  $\epsilon$  of the value of the *optimal* policy on the MDP,  $V_0^{\star}$ . Such a guarantee is immediately attainable by applying Theorem 1 with a policy set  $\Pi$  such that  $\sup_{\pi \in \Pi} V_0^{\pi} \geq V_0^{\star} - \epsilon$ . The following result shows that it is possible to construct such a set of policies, and therefore learn a *globally* near-optimal policy.

Corollary 1. There exists a set of policies  $\Pi_{\epsilon}$  such that  $\log |\Pi_{\epsilon}| \leq \widetilde{\mathcal{O}}(dH^2 \cdot \log 1/\epsilon)$  and, for any linear MDP satisfying Definition 3.1,  $\sup_{\pi \in \Pi_{\epsilon}} V_0^{\pi} \geq V_0^{\star} - \epsilon$ . If we run PEDEL with  $\Pi \leftarrow \Pi_{\epsilon}$ , then with probability at least  $1 - \delta$ , it returns a policy  $\widehat{\pi}$  such that  $V_0^{\widehat{\pi}} \geq V_0^{\star} - 2\epsilon$ , and runs for at most

$$C_0 H^4 \cdot \sum_{h=1}^{H} \inf_{\mathbf{\Lambda}_{\exp} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi_{\epsilon}} \frac{\|\phi_{\pi,h}\|_{\mathbf{\Lambda}_{\exp}^{-1}}^2}{\max\{V_0^{\star} - V_0^{\pi}, \epsilon\}^2} \cdot \left(dH^2 + \log \frac{1}{\delta}\right) + C_1$$

episodes, for  $C_0 = \text{poly} \log(d, H, \frac{1}{\epsilon})$ .

While Theorem 1 and Corollary 1 quantify the instance-dependent complexity of learning, it is natural to ask what the *worst-case* complexity of PEDEL is. The following result provides such a bound.

Corollary 2. For any linear MDP satisfying Definition 3.1,  $\inf_{\Lambda_{\exp} \in \Omega_h} \max_{\pi \in \Pi_{\epsilon}} \|\phi_{\pi,h}\|_{\Lambda_{\exp}^{-1}}^2 \leq d$ , so the sample complexity of Algorithm 1 when run with  $\Pi \leftarrow \Pi_{\epsilon}$  is no larger than

$$\widetilde{\mathcal{O}}\left(\frac{dH^5(dH^2 + \log 1/\delta)}{\epsilon^2} + C_1\right).$$

Corollary 2 shows that PEDEL has worst-case optimal dimension dependence, matching the lower bound of  $\Omega(d^2H^2/\epsilon^2)$  given in Wagenmaker et al. (2022), up to H and log factors<sup>1</sup>.

 $<sup>^{1}</sup>$ We remark that the focus of this work is on instance-dependence and dimension-dependence, not in optimizing H factors, and we leave improving our H dependence for future work.

Remark 4.1 (Performance on Linear Contextual Bandits). Corollary 1 applies directly to linear contextual bandits by setting  $H=1^2$ . To our knowledge, this is the first instance-dependent result on PAC policy identification in linear contextual bandits. Furthermore, Corollary 2 shows that we also obtain a worst-case complexity of  $\widetilde{\mathcal{O}}(d^2/\epsilon^2)$  on linear contextual bandits, which is the optimal rate (Wagenmaker et al., 2022).

### 4.2 Low-Regret Algorithms are Suboptimal for PAC RL in Large State-Spaces

We next show that there are problems on which the instance-dependent complexity of PEDEL improves on the worst-case lower bound shown in Wagenmaker et al. (2022), thereby demonstrating that we do indeed obtain favorable complexities on "easy" instances.

**Proposition 2.** For any d > 2, there exists a d-dimensional linear MDP with H = 2 such that with probability  $1 - \delta$ , PEDEL identifies an  $\epsilon$ -optimal policy on this MDP after running for only  $\widetilde{\mathcal{O}}\left(\frac{\log d/\delta}{\epsilon^2} + \operatorname{poly}(d, \log \frac{1}{\delta}, \log \frac{1}{\epsilon})\right)$  episodes.

The complexity given in Proposition 2 is a factor of  $d^2$  better than the worst-case lower bound of  $\Omega(d^2/\epsilon^2)$ . While this shows that PEDEL yields a significant improvement over existing worst-case lower bounds on favorable instances, it is natural to ask whether the same complexity is attainable with existing algorithms, perhaps by applying a tighter analysis. Towards answering this, we will consider a class of low-regret algorithms and an online-to-batch learning protocol.

**Definition 4.1** (Low-Regret Algorithm). We say that an algorithm is a *low-regret algorithm* if its expected regret is bounded as, for all K:

$$\mathbb{E}[\mathcal{R}_K] = \sum_{k=1}^K \mathbb{E}[V_0^* - V_0^{\pi_k}] \le \mathcal{C}_1 K^\alpha + \mathcal{C}_2$$

for some constants  $C_1, C_2$ , and  $\alpha \in (0, 1)$ .

**Protocol 4.1** (Online-to-Batch Learning). The online-to-batch protocol proceeds as follows:

- 1. The learner plays a low-regret algorithm satisfying Definition 4.1 for K episodes.
- 2. The learner stops at a (possibly random) time K, and, using the observations it has collected in any way it wishes, outputs a policy  $\widehat{\pi}$  it believes is  $\epsilon$ -optimal.

In general, by applying online-to-batch learning, one can convert a regret guarantee of  $C_1K^{\alpha} + C_2$  to a PAC complexity of  $\mathcal{O}((\frac{C_1}{\epsilon})^{\frac{1}{1-\alpha}} + \frac{C_2}{\epsilon})$  (Jin et al., 2018), allowing low-regret algorithms such as that of Zanette et al. (2020b) to obtain the minimax-optimal PAC complexity of  $\mathcal{O}(d^2H^4/\epsilon^2)$ . The following result shows, however, that this protocol is unable to obtain the instance-optimal rate.

**Proposition 3.** On the instance of Proposition 2, for small enough  $\epsilon$ , any learner that is  $(\epsilon, \delta)$ -PAC and follows Protocol 4.1 with stopping time K must have  $\mathbb{E}[K] \geq \Omega(\frac{d \cdot \log 1/\delta}{\epsilon^2})$ .

Together, Proposition 2 and Proposition 3 show that running a low-regret algorithm to learn a near-optimal policy in a linear MDP is provably suboptimal—at least a factor of d worse than the instance-dependent rate obtained by PEDEL. While a similar observation was recently made in the setting of tabular MDPs (Wagenmaker et al., 2021b), to our knowledge, this is the first such result in the RL with function approximation setting, implying that, in this setting, low-regret algorithms are insufficient for obtaining optimal PAC sample complexity. As standard optimistic algorithms are also low-regret, this result implies that all such optimistic algorithms are also suboptimal.

<sup>&</sup>lt;sup>2</sup>We describe the exact mapping to linear contextual bandits in Appendix B.3.

#### 4.3 Tabular and Deterministic MDPs

To relate our results to existing results on instance-dependent RL, we next turn to the setting of tabular MDPs, where it is assumed that  $S := |\mathcal{S}| < \infty, A := |\mathcal{A}| < \infty$ . Define:

$$\Delta_h(s, a) = V_h^{\star}(s) - Q_h^{\star}(s, a), \quad w_h^{\pi}(s, a) = \mathbb{P}_{\pi}[s_h = s, a_h = a].$$

 $\Delta_h(s,a)$  denotes the value-function gap, and quantifies the suboptimality of playing action a in state s at step h and then playing the optimal policy, as compared to taking the optimal action in (s,h).  $w_h^{\pi}(s,a)$  denotes the state-action visitation distribution for policy  $\pi$ , and quantifies how likely policy  $\pi$  is to reach (s,a) at step h. Note that  $[\phi_{\pi,h}]_{(s,a)} = w_h^{\pi}(s,a)$ . We obtain the following corollary.

Corollary 3. In the setting of tabular MDPs, PEDEL outputs an  $\epsilon$ -optimal policy with probability at least  $1 - \delta$ , and has sample complexity bounded as

$$\widetilde{\mathcal{O}}\left(\sum_{h=1}^{H}\inf_{\pi \in \Pi}\max_{s,a}\max_{s,a}\frac{H^{4}}{w_{h}^{\pi \exp}(s,a)}\min\left\{\frac{1}{w_{h}^{\pi}(s,a)\Delta_{h}(s,a)^{2}},\frac{w_{h}^{\pi}(s,a)}{\Delta_{\min}(\Pi)^{2}},\frac{w_{h}^{\pi}(s,a)}{\epsilon^{2}}\right\}\cdot\left(SH+\log\frac{1}{\delta}\right)+C_{1}\right),$$

for 
$$C_1 = \text{poly}(S, A, H, \frac{1}{\min_h \min_s \sup_{\pi} w_h^{\pi}(s)}, \log \frac{1}{\delta}, \log \frac{1}{\epsilon})$$
 and  $\Pi$  the set of all deterministic policies.

For tabular MDPs, the primary comparable result on instance-dependent policy identification is that obtained by Wagenmaker et al. (2021b), which introduces a different measure of complexity, the *gap-visitation complexity*, and an algorithm, Moca, with sample complexity scaling as the gap-visitation complexity. The following result shows that the complexity Pedel obtains on tabular MDPs and the gap-visitation complexity do not have a clear ordering.

**Proposition 4.** Fix any  $\epsilon \in (0, 1/2)$  and  $S \ge \log_2(1/\epsilon)$ . Then there exist tabular MDPs  $\mathcal{M}_1$  and  $\mathcal{M}_2$ , each with H = 2, S states, and  $\mathcal{O}(S)$  actions, such that:

- On  $\mathcal{M}_1$ , the complexity bound of PEDEL given in Corollary 3 scales as  $\operatorname{poly}(S, \log 1/\delta)$ , while the gap-visitation complexity scales as  $\Omega(1/\epsilon^2)$ .
- On  $\mathcal{M}_2$ , the complexity bound of PEDEL given in Corollary 3 scales as  $\Omega(1/\epsilon^2)$ , while the gap-visitation complexity scales as poly(S, log  $1/\delta$ ).

The lack of ordering between the two complexity measures arises because, on some problem instances, it is easier to learn in policy-space (as PEDEL does), while on other instances, it is easier to learn near-optimal actions on individual states directly, and then synthesize these actions into a near-optimal policy (the approach Moca takes). This difference arises because, in the former instance, the minimum policy gap is large  $(V_0^{\star} - V_0^{\pi} = \Omega(1))$  for every deterministic policy  $\pi \neq \pi^{\star}$ , while in the latter instance, the minimum policy gap is small, but all value-function gaps are large, satisfying  $\Delta_h(s,a) = \Omega(1)$  for all  $a \neq \arg\max_{a \in \mathcal{A}} Q_h^{\star}(s,a)$  and all s and h. Thus, on the former instance, it is much easier to learn over the space of policies, while on the latter it is much easier to learn optimal actions in individual states. Resolving this discrepency with an algorithm able to achieve the "best-of-both-worlds" is an interesting direction for future work.

**Deterministic MDPs.** Finally, we turn to the simplified setting of tabular, deterministic MDPs. Here, for each (s, a, h), there exists some s' such that  $P_h(s'|s, a) = 1$ . We still allow the rewards to be random, however, so the agent must still learn in order to find a near-optimal policy. Following the same notation as the recent work of Tirinzoni et al. (2022), let  $\Pi_{sah} = \{\pi \text{ deterministic } : s_h^{\pi} = \{\pi \text{ deterministic } : s$ 

 $s, a_h^{\pi} = a$ }, where  $s_h^{\pi}$  and  $a_h^{\pi}$  are the state and action policy  $\pi$  will be in at step h (note that these quantities are well-defined quantities for deterministic policies). Also define the deterministic return gap as  $\bar{\Delta}_h(s,a) := V_0^{\star} - \max_{\pi \in \Pi_{sah}} V_0^{\pi}$ , and let  $\bar{\Delta}_{\min} := \min_{s,a,h:\bar{\Delta}_h(s,a)>0} \bar{\Delta}_h(s,a)$  in the case when there exists a unique optimal deterministic policy, and  $\bar{\Delta}_{\min} := 0$  otherwise. We obtain the following.

Corollary 4. In the setting of tabular, deterministic MDPs, PEDEL outputs an  $\epsilon$ -optimal policy with probability at least  $1 - \delta$ , and has sample complexity bounded as

$$\widetilde{\mathcal{O}}\left(H^4 \cdot \sum_{h=1}^{H} \sum_{s,a} \frac{1}{\max\{\bar{\Delta}_h(s,a), \bar{\Delta}_{\min}, \epsilon\}^2} \cdot (H + \log \frac{1}{\delta}) + \operatorname{poly}\left(S, A, H, \log \frac{1}{\delta}, \log \frac{1}{\epsilon}\right)\right).$$

Up to H and log factors and lower-order terms, the rate given in Corollary 4 matches the instance-dependent lower bound given in Tirinzoni et al.  $(2022)^3$ . Thus, we conclude that, in the setting of tabular, deterministic MDPs, PEDEL is (nearly) instance-optimal. While Tirinzoni et al. (2022) also obtain instance-optimality in this setting, their algorithm and analysis are specialized to tabular, deterministic MDPs—in contrast, PEDEL requires no modification from its standard operation.

# 5 Online Experiment Design in Linear MDPs

As described in Section 4, to reduce our uncertainty and explore in a way that only targets the relevant feature directions, we must solve an XY-experiment design problem of the form:

$$\inf_{\mathbf{\Lambda}_{\exp} \in \Omega_h} \max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{\mathbf{\Lambda}_{\exp}}^2, \tag{5.1}$$

where here  $\Phi$  will be some set of estimated feature-visitations. Recall that  $\Omega_h$  denotes the set of covariance matrices realizable on our MDP, and therefore without knowledge of the MDP dynamics we cannot specify this set. If follows that it is not in general possible to solve (5.1) without knowledge of the MDP dynamics. In this section we describe our approach to solving (5.1) without this knowledge by relying on a low-regret algorithm as an optimization primitive.

Approximating Frank-Wolfe via Regret Minimization. Given knowledge of the MDP dynamics, we could compute  $\Omega_h$  directly, and apply the celebrated Frank-Wolfe coordinate-descent algorithm (Frank & Wolfe, 1956) to solve (5.1). In this setting the Frank-Wolfe update for (5.1) is:

$$\mathbf{\Gamma}_{t} = \arg\min_{\mathbf{\Gamma} \in \mathbf{\Omega}_{h}} \langle \nabla_{\mathbf{\Lambda}} (\max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{\mathbf{\Lambda}^{-1}}^{2}) |_{\mathbf{\Lambda} = \mathbf{\Lambda}_{t}}, \mathbf{\Gamma} \rangle, \quad \mathbf{\Lambda}_{t+1} = (1 - \gamma_{t}) \mathbf{\Lambda}_{t} + \gamma_{t} \mathbf{\Gamma}_{t}$$
 (5.2)

for step size  $\gamma_t$ . Standard Frank-Wolfe analysis shows that this update converges to a near-optimal solution to (5.1) at a polynomial rate. However, without knowledge of  $\Omega_h$ , we are unable to solve for  $\Gamma_t$  and run the Frank-Wolfe update.

Our critical observation is that the minimization over  $\Omega_h$  in (5.2) can be approximated without knowledge of  $\Omega_h$  by running a low-regret algorithm on a particular objective. Some calculation shows that (except on a measure-zero set, assuming  $\Phi$  is finite)  $\nabla_{\mathbf{\Lambda}}(\max_{\phi \in \Phi} \|\phi\|_{\mathbf{\Lambda}^{-1}}^2)|_{\mathbf{\Lambda} = \mathbf{\Lambda}_t} = -\mathbf{\Lambda}_t^{-1}\widetilde{\phi}_t\widetilde{\phi}_t^{\top}\mathbf{\Lambda}_t^{-1}$  for  $\widetilde{\phi}_t = \arg\max_{\phi \in \Phi} \|\phi\|_{\mathbf{\Lambda}_t^{-1}}^2$ . If  $\mathbf{\Gamma} = \mathbf{\Lambda}_{\pi,h} = \mathbb{E}_{\pi}[\phi_h\phi_h^{\top}]$  for some  $\pi$ , we have

$$\langle \nabla_{\boldsymbol{\Lambda}} (\max_{\boldsymbol{\phi} \in \boldsymbol{\Phi}} \|\boldsymbol{\phi}\|_{\boldsymbol{\Lambda}^{-1}}^2) |_{\boldsymbol{\Lambda} = \boldsymbol{\Lambda}_t}, \boldsymbol{\Gamma} \rangle = -\mathrm{tr}(\boldsymbol{\Lambda}_t^{-1} \widetilde{\boldsymbol{\phi}}_t \widetilde{\boldsymbol{\phi}}_t^\top \boldsymbol{\Lambda}_t^{-1} \boldsymbol{\Lambda}_{\pi,h}) = -\mathbb{E}_{\pi}[(\boldsymbol{\phi}_h^\top \boldsymbol{\Lambda}_t^{-1} \widetilde{\boldsymbol{\phi}}_t)^2].$$

<sup>&</sup>lt;sup>3</sup>The lower bound of Tirinzoni et al. (2022) depends on a slightly different (but nearly equivalent) minimum gap term,  $\bar{\Delta}_{\min}^h$ . Similar to our upper bound, the upper bound of Tirinzoni et al. (2022) scales with  $\bar{\Delta}_{\min}$  instead of  $\bar{\Delta}_{\min}^h$ . We offer a more in-depth discussion of this point in Appendix B.3.

Now, if we run a low-regret algorithm on the (deterministic) reward  $\nu_h^t(s, a) = (\phi(s, a)^{\top} \mathbf{\Lambda}_t^{-1} \widetilde{\phi}_t)^2$  for a sufficiently large number of episodes K, we will be guaranteed to collect reward at a rate close to that of the optimal policy, which implies we will collect some data  $\{\phi_{h,\tau}\}_{\tau=1}^K$  such that

$$K^{-1} \cdot \widetilde{\boldsymbol{\phi}}_t^{\top} \widehat{\boldsymbol{\Gamma}}_K \widetilde{\boldsymbol{\phi}}_t := K^{-1} \sum_{\tau=1}^K (\boldsymbol{\phi}_{h,\tau}^{\top} \boldsymbol{\Lambda}_t^{-1} \widetilde{\boldsymbol{\phi}}_t)^2 \approx \sup_{\pi} \mathbb{E}_{\pi} [(\boldsymbol{\phi}_h^{\top} \boldsymbol{\Lambda}_t^{-1} \widetilde{\boldsymbol{\phi}}_t)^2]. \tag{5.3}$$

However, this implies the covariates we have collected,  $\widehat{\Gamma}_K$ , approximately minimize (5.2). In other words, running a low-regret algorithm on  $\nu_h^t$  allows us to obtain covariates which are approximately the solution to the minimization in the Frank-Wolfe update—without knowledge of  $\Omega_h$ , we can solve the Frank-Wolfe update by running a low-regret algorithm, and therefore solve (5.1). This motivates Algorithm 2.

# Algorithm 2 Online Frank-Wolfe via Regret Minimization (informal)

- 1: **input:** uncertain feature directions  $\Phi$ , step h, regularization  $\Lambda_0 \succ 0$
- 2:  $K_0 \leftarrow$  sufficiently large number of episodes to guarantee (5.3) holds
- 3: Run any policy for  $K_0$  episodes, collect data  $\{\phi_{h,\tau}\}_{\tau=1}^{K_0}$ , set  $\Lambda_1 \leftarrow K_0^{-1} \sum_{\tau=1}^{K_0} \phi_{h,\tau} \phi_{h,\tau}^{\top}$
- 4: **for** t = 1, ..., T 1 **do**
- 5:  $\widetilde{\boldsymbol{\phi}}_t \leftarrow \underset{\boldsymbol{\phi} \in \Phi}{\operatorname{rg} \max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{(\boldsymbol{\Lambda}_t + \boldsymbol{\Lambda}_0)^{-1}}^2}, \ \nu_h^t(s, a) \leftarrow (\boldsymbol{\phi}(s, a)^\top (\boldsymbol{\Lambda}_t + \boldsymbol{\Lambda}_0)^{-1} \widetilde{\boldsymbol{\phi}}_t)^2$
- 6: Run low-regret algorithm on  $\nu_h^t$  for  $K_0$  episodes, collect covariates  $\widehat{\Gamma}_{K_0}^t$
- 7: Set  $\Lambda_{t+1} \leftarrow (1 \gamma_t) \Lambda_t + \gamma_t K_0^{-1} \widehat{\Gamma}_{K_0}^t$  for  $\gamma_t = \frac{1}{t+1}$
- 8: **return:** covariates  $TK_0\Lambda_T = \sum_{t=1}^{T-1} \widehat{\Gamma}_{K_0}^t + \Lambda_1$

**Theorem 5** (informal). Consider running Algorithm 2 with some  $\Lambda_0 \succ 0$ . Then with properly chosen settings of  $K_0$  and T, we can guarantee that, with probability at least  $1 - \delta$ , we will run for at most

$$N \leq 20 \cdot \frac{\inf_{\mathbf{\Lambda}_{\exp} \in \mathbf{\Omega}_h} \max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{(\mathbf{\Lambda}_{\exp} + \mathbf{\Lambda}_0)^{-1}}^2}{\epsilon_{\exp}} + \operatorname{poly}\left(d, H, \|\mathbf{\Lambda}_0^{-1}\|_{\operatorname{op}}, \log |\Phi|, \log 1/\delta\right)$$

episodes, and return covariance  $\widehat{\mathbf{\Lambda}}_N$  satisfying  $\max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{(\widehat{\mathbf{\Lambda}}_N + N\mathbf{\Lambda}_0)^{-1}}^2 \leq \epsilon_{\exp}$ .

Note that this rate is essentially optimal, up to constants and lower-order terms. If we let  $\omega_{\exp}^{\star}$  denote the distribution over policies which minimize (5.1), then to collect covariance  $\widehat{\Lambda}_N$  such that  $\max_{\phi \in \Phi} \|\phi\|_{(\widehat{\Lambda}_N + N\Lambda_0)^{-1}}^2 \le \epsilon_{\exp}$ , in expectation, we would need to play  $\pi \sim \omega_{\exp}^{\star}$  for at least

$$\frac{\inf_{\mathbf{\Lambda}_{\exp} \in \mathbf{\Omega}_h} \max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{(\mathbf{\Lambda}_{\exp} + \mathbf{\Lambda}_0)^{-1}}^2}{\epsilon_{\exp}}$$

episodes, which is the same scaling as obtained in Theorem 5.

In practice, we instead run Algorithm 2 on a smoothed version of the objective in (5.1). We provide a full definition of Algorithm 2 with exact setting of T and  $K_0$  in Appendix  $\mathbb{C}$ .

#### 5.1 Experiment Design in MDPs with General Objective Functions

While the experiment design in (5.1) is the natural design if our goal is to identify a near-optimal policy, in general we may be interested in collecting data to minimize some other objective; that is,

solving an experiment design of the form:

$$\inf_{\mathbf{\Lambda}_{\exp} \in \mathbf{\Omega}_h} f(\mathbf{\Lambda}_{\exp})$$

for some function f defined over the space of PSD matrices. For example, we could take  $f(\mathbf{\Lambda}_{\text{exp}}) = \|\mathbf{\Lambda}_{\text{exp}}^{-1}\|_{\text{op}} = \frac{1}{\lambda_{\min}(\mathbf{\Lambda}_{\text{exp}})}$ , and the above experiment design would correspond to maximizing the minimum eigenvalue of the collected covariates, or E-optimal design (Pukelsheim, 2006).

Motivated by this, in Appendix C we generalize Theorem 5 and Algorithm 2 to handle a much broader class of experiment design problems. In particular, we consider all *smooth experiment design objectives*, which we define as follows.

**Definition 5.1** (Smooth Experiment Design Objectives). We say that  $f(\Lambda): \mathbb{S}^d_+ \to \mathbb{R}$  is a *smooth* experiment design objective if it satisfies the following conditions:

- f is convex, differentiable, and  $\beta$  smooth in the norm  $\|\cdot\|: \|\nabla f(\mathbf{\Lambda}) \nabla f(\mathbf{\Lambda}')\|_* \leq \beta \|\mathbf{\Lambda} \mathbf{\Lambda}'\|.$
- f is L-lipschitz in the operator norm:  $|f(\Lambda) f(\Lambda')| \leq L ||\Lambda \Lambda'||_{\text{op}}$ .
- Let  $\Xi_{\Lambda_0} := -\nabla_{\Lambda} f(\Lambda)|_{\Lambda = \Lambda_0}$ . Then  $\Xi_{\Lambda_0} \succeq 0$  and  $\operatorname{tr}(\Xi_{\Lambda_0}) \leq M$  for all  $\Lambda_0 \succeq 0$  satisfying  $\|\Lambda_0\|_{\operatorname{op}} \leq 1$ .

Our generalization of Algorithm 2 to handle all smooth experiment design objectives—OptCov, defined in Appendix C.3—enjoys the following guarantee.

**Theorem 6.** Fix  $h \in [H]$ , consider some f satisfying Definition 5.1, and let  $f_{\min}$  be some value such that  $\inf_{\mathbf{\Lambda}_{\exp} \in \mathbf{\Omega}_h} f(\mathbf{\Lambda}_{\exp}) \ge f_{\min}$ . Then with probability at least  $1 - \delta$ , given any  $\epsilon > 0$ , OPTCOV runs for at most

$$N \leq 5 \cdot \frac{\inf_{\mathbf{\Lambda}_{\exp} \in \mathbf{\Omega}_h} f(\mathbf{\Lambda}_{\exp})}{\epsilon} + \text{poly}\left(d, H, M, \beta, L, f_{\min}^{-1}, \log 1/\delta\right)$$

episodes, and collects covariates  $\hat{\Sigma}_N = \sum_{\tau=1}^N \phi_{h,\tau} \phi_{h,\tau}^{\top}$  such that

$$f(N^{-1}\widehat{\Sigma}_N) \le N\epsilon.$$

We will often be interested in objectives f that satisfy  $f(a\mathbf{\Lambda}) = a^{-1}f(\mathbf{\Lambda})$  for a scalar a, in which case the guarantee  $f(N^{-1}\widehat{\Sigma}_N) \leq N\epsilon$  reduces to  $f(\widehat{\Sigma}_N) \leq \epsilon$ . We note also that many typical experiment design objectives are non-smooth. As we show in Appendix D, however, it is often possible to derive smoothed versions of such objectives with negligible approximation error.

# 6 Conclusion

In this work, we have shown that it is possible to obtain instance-dependent guarantees in RL with function approximation, and that our algorithm, PEDEL, yields provable gains over low-regret algorithms. As the first result of its kind in this setting, it opens several directions for future work.

The computational complexity of PEDEL scales as  $\operatorname{poly}(d, H, \frac{1}{\epsilon}, |\Pi|, |\mathcal{A}|, \log \frac{1}{\delta})$ . In general, to ensure  $\Pi$  contains an  $\epsilon$ -optimal policy,  $|\Pi|$  must be exponential in problem parameters, rendering PEDEL computationally inefficient. Furthermore, the sample complexity of PEDEL scales with  $\lambda_{\min}^{\star}$ , the "hardest-to-reach" direction. While this is not uncommon in the literature, we might hope that

if a direction is very difficult to reach, learning in that direction should not be necessary, as we are unlikely to ever encounter it. Obtaining an algorithm with a similar instance-dependence but that is computationally efficient and does not depend on  $\lambda_{\min}^{\star}$  is an interesting direction for future work.

Extending our results to the setting of general function approximation is also an exciting direction. While our results do rely on the linear structure of the MDP, we believe the online experiment-design approach we propose could be generally applicable in more complex settings. As a first step, it could be interesting to extend our approach to the setting of Bilinear classes (Du et al., 2021), which also exhibits a certain linear structure.

# Acknowledgements

The work of AW was supported by an NSF GFRP Fellowship DGE-1762114. The work of KJ was funded in part by the AFRL and NSF TRIPODS 2023166.

# References

- Agarwal, N., Chaudhuri, S., Jain, P., Nagaraj, D., and Netrapalli, P. Online target q-learning with reverse experience replay: Efficiently finding the optimal policy for linear mdps. arXiv preprint arXiv:2110.08440, 2021.
- Ayoub, A., Jia, Z., Szepesvari, C., Wang, M., and Yang, L. Model-based reinforcement learning with value-targeted regression. In *International Conference on Machine Learning*, pp. 463–474. PMLR, 2020.
- Baird, L. Residual algorithms: Reinforcement learning with function approximation. In *Machine Learning Proceedings* 1995, pp. 30–37. Elsevier, 1995.
- Bradtke, S. J. and Barto, A. G. Linear least-squares algorithms for temporal difference learning. *Machine learning*, 22(1):33–57, 1996.
- Brafman, R. I. and Tennenholtz, M. R-max-a general polynomial time algorithm for near-optimal reinforcement learning. *Journal of Machine Learning Research*, 3(Oct):213–231, 2002.
- Dann, C. and Brunskill, E. Sample complexity of episodic fixed-horizon reinforcement learning. arXiv preprint arXiv:1510.08906, 2015.
- Dann, C., Li, L., Wei, W., and Brunskill, E. Policy certificates: Towards accountable reinforcement learning. In *International Conference on Machine Learning*, pp. 1507–1516. PMLR, 2019.
- Dann, C., Marinov, T. V., Mohri, M., and Zimmert, J. Beyond value-function gaps: Improved instance-dependent regret bounds for episodic reinforcement learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- Du, S. S., Kakade, S. M., Wang, R., and Yang, L. F. Is a good representation sufficient for sample efficient reinforcement learning? arXiv preprint arXiv:1910.03016, 2019.
- Du, S. S., Kakade, S. M., Lee, J. D., Lovett, S., Mahajan, G., Sun, W., and Wang, R. Bilinear classes: A structural framework for provable generalization in rl. arXiv preprint arXiv:2103.10897, 2021.
- Epasto, A., Mahdian, M., Mirrokni, V., and Zampetakis, E. Optimal approximation-smoothness tradeoffs for soft-max functions. *Advances in Neural Information Processing Systems*, 33:2651–2660, 2020.
- Fiez, T., Jain, L., Jamieson, K. G., and Ratliff, L. Sequential experimental design for transductive linear bandits. *Advances in neural information processing systems*, 32, 2019.
- Foster, D. J., Kakade, S. M., Qian, J., and Rakhlin, A. The statistical complexity of interactive decision making. arXiv preprint arXiv:2112.13487, 2021.
- Frank, M. and Wolfe, P. An algorithm for quadratic programming. *Naval research logistics quarterly*, 3(1-2):95–110, 1956.
- Freedman, D. A. On tail probabilities for martingales. the Annals of Probability, pp. 100–118, 1975.
- Hao, B., Lattimore, T., Szepesvári, C., and Wang, M. Online sparse reinforcement learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 316–324. PMLR, 2021.

- He, J., Zhou, D., and Gu, Q. Logarithmic regret for reinforcement learning with linear function approximation. arXiv preprint arXiv:2011.11566, 2020.
- Jia, Z., Yang, L., Szepesvari, C., and Wang, M. Model-based reinforcement learning with value-targeted regression. In *Learning for Dynamics and Control*, pp. 666–686. PMLR, 2020.
- Jiang, N., Krishnamurthy, A., Agarwal, A., Langford, J., and Schapire, R. E. Contextual decision processes with low bellman rank are pac-learnable. In *International Conference on Machine Learning*, pp. 1704–1713. PMLR, 2017.
- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. Is q-learning provably efficient? In *Proceedings* of the 32nd International Conference on Neural Information Processing Systems, pp. 4868–4878, 2018.
- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pp. 2137–2143. PMLR, 2020.
- Jin, C., Liu, Q., and Miryoosefi, S. Bellman eluder dimension: New rich classes of rl problems, and sample-efficient algorithms. arXiv preprint arXiv:2102.00815, 2021.
- Jonsson, A., Kaufmann, E., Ménard, P., Domingues, O. D., Leurent, E., and Valko, M. Planning in markov decision processes with gap-dependent sample complexity. arXiv preprint arXiv:2006.05879, 2020.
- Kakade, S. M. On the sample complexity of reinforcement learning. PhD thesis, UCL (University College London), 2003.
- Kaufmann, E., Cappé, O., and Garivier, A. On the complexity of best-arm identification in multi-armed bandit models. *The Journal of Machine Learning Research*, 17(1):1–42, 2016.
- Kearns, M. and Singh, S. Finite-sample convergence rates for q-learning and indirect algorithms.

  Advances in neural information processing systems, 11, 1998.
- Lattimore, T. and Szepesvári, C. Bandit algorithms. Cambridge University Press, 2020.
- Marjani, A. A. and Proutiere, A. Best policy identification in discounted mdps: Problem-specific sample complexity. arXiv preprint arXiv:2009.13405, 2020.
- Marjani, A. A., Garivier, A., and Proutiere, A. Navigating to the best policy in markov decision processes. arXiv preprint arXiv:2106.02847, 2021.
- McSherry, F. and Talwar, K. Mechanism design via differential privacy. In 48th Annual IEEE Symposium on Foundations of Computer Science (FOCS'07), pp. 94–103. IEEE, 2007.
- Melo, F. S. and Ribeiro, M. I. Q-learning with linear function approximation. In *International Conference on Computational Learning Theory*, pp. 308–322. Springer, 2007.
- Ménard, P., Domingues, O. D., Jonsson, A., Kaufmann, E., Leurent, E., and Valko, M. Fast active learning for pure exploration in reinforcement learning. arXiv preprint arXiv:2007.13442, 2020.
- Mutny, M., Janik, T., and Krause, A. Active exploration via experiment design in markov chains. arXiv preprint arXiv:2206.14332, 2022.

- Ok, J., Proutiere, A., and Tranos, D. Exploration in structured reinforcement learning. arXiv preprint arXiv:1806.00775, 2018.
- Pukelsheim, F. Optimal design of experiments. SIAM, 2006.
- Simchowitz, M. and Jamieson, K. Non-asymptotic gap-dependent regret bounds for tabular mdps. arXiv preprint arXiv:1905.03814, 2019.
- Soare, M., Lazaric, A., and Munos, R. Best-arm identification in linear bandits. *Advances in Neural Information Processing Systems*, 27, 2014.
- Sutton, R. S., McAllester, D., Singh, S., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.
- Tirinzoni, A., Al-Marjani, A., and Kaufmann, E. Near instance-optimal pac reinforcement learning for deterministic mdps. arXiv preprint arXiv:2203.09251, 2022.
- Vershynin, R. Introduction to the non-asymptotic analysis of random matrices. arXiv preprint arXiv:1011.3027, 2010.
- Wagenmaker, A., Chen, Y., Simchowitz, M., Du, S. S., and Jamieson, K. First-order regret in reinforcement learning with linear function approximation: A robust estimation approach. arXiv preprint arXiv:2112.03432, 2021a.
- Wagenmaker, A., Simchowitz, M., and Jamieson, K. Beyond no regret: Instance-dependent pac reinforcement learning. arXiv preprint arXiv:2108.02717, 2021b.
- Wagenmaker, A., Chen, Y., Simchowitz, M., Du, S. S., and Jamieson, K. Reward-free rl is no harder than reward-aware rl in linear markov decision processes. arXiv preprint arXiv:2201.11206, 2022.
- Wagenmaker, A. J., Simchowitz, M., and Jamieson, K. Task-optimal exploration in linear dynamical systems. In *International Conference on Machine Learning*, pp. 10641–10652. PMLR, 2021c.
- Wang, Y., Wang, R., Du, S. S., and Krishnamurthy, A. Optimism in reinforcement learning with generalized linear function approximation. arXiv preprint arXiv:1912.04136, 2019.
- Wang, Y., Wang, R., and Kakade, S. M. An exponential lower bound for linearly-realizable mdps with constant suboptimality gap. arXiv preprint arXiv:2103.12690, 2021.
- Weisz, G., Amortila, P., and Szepesvári, C. Exponential lower bounds for planning in mdps with linearly-realizable optimal action-value functions. In *Algorithmic Learning Theory*, pp. 1237–1264. PMLR, 2021.
- Xu, H., Ma, T., and Du, S. S. Fine-grained gap-dependent bounds for tabular mdps via adaptive multi-step bootstrap. arXiv preprint arXiv:2102.04692, 2021.
- Yang, L. and Wang, M. Sample-optimal parametric q-learning using linearly additive features. In *International Conference on Machine Learning*, pp. 6995–7004. PMLR, 2019.
- Zanette, A., Kochenderfer, M. J., and Brunskill, E. Almost horizon-free structure-aware best policy identification with a generative model. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

- Zanette, A., Brandfonbrener, D., Brunskill, E., Pirotta, M., and Lazaric, A. Frequentist regret bounds for randomized least-squares value iteration. In *International Conference on Artificial Intelligence and Statistics*, pp. 1954–1964. PMLR, 2020a.
- Zanette, A., Lazaric, A., Kochenderfer, M., and Brunskill, E. Learning near optimal policies with low inherent bellman error. In *International Conference on Machine Learning*, pp. 10978–10989. PMLR, 2020b.
- Zanette, A., Lazaric, A., Kochenderfer, M. J., and Brunskill, E. Provably efficient reward-agnostic navigation with linear value iteration. *Advances in Neural Information Processing Systems*, 33: 11756–11766, 2020c.
- Zhang, Z., Ji, X., and Du, S. S. Is reinforcement learning more difficult than bandits? a near-optimal algorithm escaping the curse of horizon. arXiv preprint arXiv:2009.13503, 2020.
- Zhang, Z., Yang, J., Ji, X., and Du, S. S. Variance-aware confidence set: Variance-dependent bound for linear bandits and horizon-free bound for linear mixture mdp. arXiv preprint arXiv:2101.12745, 2021.
- Zhou, D., Gu, Q., and Szepesvari, C. Nearly minimax optimal reinforcement learning for linear mixture markov decision processes. arXiv preprint arXiv:2012.08507, 2020.
- Zhou, D., He, J., and Gu, Q. Provably efficient reinforcement learning for discounted mdps with feature mapping. In *International Conference on Machine Learning*, pp. 12793–12802. PMLR, 2021.

# Contents

1	Introduction	1
	1.1 Contributions	2
2	Related Work	3
3	Preliminaries	4
	3.1 Linear MDPs	Ę
4	Near-Optimal Policy Identification in Linear MDPs	6
	4.1 Main Results	7
	4.2 Low-Regret Algorithms are Suboptimal for PAC RL in Large State-Spaces	Ĝ
	4.3 Tabular and Deterministic MDPs	10
5	Online Experiment Design in Linear MDPs	11
	5.1 Experiment Design in MDPs with General Objective Functions	12
6	Conclusion	13
$\mathbf{A}$	Technical Results	20
	A.1 Properties of Linear MDPs	20
	A.2 Feature-Visitations in Linear MDPs	21
	A.3 Constructing the Policy Class	23
В	Policy Elimination	27
	B.1 Estimating Feature-Visitations and Rewards	
	B.2 Correctness and Sample Complexity of Pedel	31
	B.3 Interpreting the Complexity	37
$\mathbf{C}$	. 0	42
	C.1 Approximate Frank-Wolfe	
	C.2 Online Frank-Wolfe via Regret Minimization	44
	C.3 Data Collection via Online Frank-Wolfe	46
D	XY-Optimal Design	<b>52</b>
	D.1 Approximating Non-Smooth Optimal Design with Smooth Optimal Design	52
	D.2 Bounding the Smoothness	
	D.3 Obtaining Well-Conditioned Covariates	
	D.4 Online XY-Optimal Design	59
${f E}$	Suboptimality of Optimistic Algorithms	61
	E.1 Linear Bandit Construction	61
	E.2 Mapping to Linear MDPs	66

# A Technical Results

**Lemma A.1** (Vershynin (2010)). For any  $\epsilon > 0$ , the  $\epsilon$ -covering number of the Euclidean ball  $\mathcal{B}^d(R) := \{ \boldsymbol{x} \in \mathbb{R}^d : \|\boldsymbol{x}\|_2 \leq R \}$  with radius R > 0 in the Euclidean metric is upper bounded by  $(1 + 2R/\epsilon)^d$ .

**Lemma A.2** (Lemma A.4 of Wagenmaker et al. (2022)). If  $x \ge C(2n)^n \log^n(2nCB)$  for  $n, C, B \ge 1$ , then  $x \ge C \log^n(Bx)$ .

**Lemma A.3** (McSherry & Talwar (2007); Epasto et al. (2020)). Consider some  $(x_i)_{i=1}^n$ . Then if  $\eta \ge \log(n)/\delta$ , we have

$$\frac{\sum_{i=1}^{n} e^{\eta x_i} x_i}{\sum_{i=1}^{n} e^{\eta x_i}} \ge \max_{i \in [n]} x_i - \delta.$$

**Lemma A.4** (Azuma-Hoeffding). et  $\mathcal{F}_0 \subset \mathcal{F}_1 \subset \ldots \subset \mathcal{F}_T$  be a filtration and let  $X_1, X_2, \ldots, X_T$  be real random variables such that  $X_t$  is  $\mathcal{F}_t$ -measurable,  $\mathbb{E}[X_t|\mathcal{F}_{t-1}] = 0$ , and  $|X_t| \leq b$  almost surely. Then for any  $\delta \in (0,1)$ , we have with probability at least  $1-\delta$ ,

$$\left| \sum_{t=1}^{T} X_t \right| \le \sqrt{8b^2 \log 2/\delta}.$$

**Lemma A.5** (Freedman's Inequality (Freedman, 1975)). Let  $\mathcal{F}_0 \subset \mathcal{F}_1 \subset \ldots \subset \mathcal{F}_T$  be a filtration and let  $X_1, X_2, \ldots, X_T$  be real random variables such that  $X_t$  is  $\mathcal{F}_t$ -measurable,  $\mathbb{E}[X_t|\mathcal{F}_{t-1}] = 0$ ,  $|X_t| \leq b$  almost surely, and  $\sum_{t=1}^T \mathbb{E}[X_t^2|\mathcal{F}_{t-1}] \leq V$  for some fixed V > 0 and b > 0. Then for any  $\delta \in (0,1)$ , we have with probability at least  $1-\delta$ ,

$$\sum_{t=1}^{T} X_t \le 2\sqrt{V \log 1/\delta} + b \log 1/\delta.$$

# A.1 Properties of Linear MDPs

**Lemma A.6.** For any linear MDP satisfying Definition 3.1, we must have that  $\|\phi(s, a)\|_2 \ge 1/\sqrt{d}$  for all s and a, and  $\|\phi_{\pi,h}\|_2 \ge 1/\sqrt{d}$  for all  $\pi$  and h.

*Proof.* By Definition 3.1, we know that  $P_h(\cdot|s,a) = \langle \phi(s,a), \mu_h(\cdot) \rangle$  forms a valid probability distribution, and that  $\|\int_{\mathcal{S}} |\mathrm{d}\mu_h(s)|\|_2 \leq \sqrt{d}$ . It follows that

$$1 = \int_{\mathcal{S}} \langle \boldsymbol{\phi}(s, a), d\boldsymbol{\mu}_h(s) \rangle \le \|\boldsymbol{\phi}(s, a)\|_2 \|\int_{\mathcal{S}} |d\boldsymbol{\mu}_h(s)| \|_2 \le \sqrt{d} \|\boldsymbol{\phi}(s, a)\|_2$$

from which the first result follows.

For the second result, using that  $1 = \int_{\mathcal{S}} \langle \phi(s, a), d\mu_h(s) \rangle$ , we get

$$\int_{\mathcal{S}} \langle \phi_{\pi,h}, d\boldsymbol{\mu}_{h}(s) \rangle = \int_{\mathcal{S}} \langle \mathbb{E}_{\pi}[\phi_{h}], d\boldsymbol{\mu}_{h}(s) \rangle 
= \mathbb{E}_{\pi} \left[ \int_{\mathcal{S}} \langle \phi_{h}, d\boldsymbol{\mu}_{h}(s) \rangle \right] 
= \mathbb{E}_{\pi}[1]$$

where we can exchange the order of integration by Fubini's Theorem since the integrand is absolutely integrable, by Definition 3.1. As above, we then have

$$1 = \int_{\mathcal{S}} \langle \phi_{\pi,h}, d\mu_h(s) \rangle \le \sqrt{d} \|\phi_{\pi,h}\|_2$$

so the second result follows.

#### A.2 Feature-Visitations in Linear MDPs

Define

$$\phi_{\pi,h} = \mathbb{E}_{\pi}[\phi(s_h, a_h)], \quad \phi_{\pi,h}(s) = \sum_{a \in \mathcal{A}} \phi(s, a) \pi_h(a|s)$$

and

$$\mathcal{T}_{\pi,h} := \int \phi_{\pi,h}(s) \mathrm{d}\boldsymbol{\mu}_{h-1}(s)^{\top}.$$

Lemma A.7.  $\phi_{\pi,h} = \mathcal{T}_{\pi,h}\phi_{\pi,h-1} = \ldots = \mathcal{T}_{\pi,h}\ldots\mathcal{T}_{\pi,1}\phi_{\pi,0}$ .

*Proof.* By the linear MDP assumption, we have:

$$\phi_{\pi,h} = \mathbb{E}_{\pi}[\phi(s_h, a_h)]$$

$$= \mathbb{E}_{\pi}[\mathbb{E}[\phi(s_h, a_h) | \mathcal{F}_{h-1}]]$$

$$= \mathbb{E}_{\pi}[\int \int \phi(s, a) d\pi_h(a|s) d\boldsymbol{\mu}_{h-1}(s)^{\top} \phi(s_{h-1}, a_{h-1})]$$

$$= \mathbb{E}_{\pi}[\int \phi_{\pi,h}(s) d\boldsymbol{\mu}_{h-1}(s)^{\top} \phi(s_{h-1}, a_{h-1})]$$

$$= \int \phi_{\pi,h}(s) d\boldsymbol{\mu}_{h-1}(s)^{\top} \mathbb{E}_{\pi}[\phi(s_{h-1}, a_{h-1})]$$

$$= \mathcal{T}_{\pi,h} \phi_{\pi,h-1}.$$

This yields the first equality. Repeating this calculation h-1 more times yields the final equality.  $\Box$ 

**Lemma A.8.** Fix some h and i < h, and consider the vector

$$oldsymbol{v} := \mathcal{T}_{\pi,i+1}^ op \mathcal{T}_{\pi,i+2}^ op \ldots \mathcal{T}_{\pi,h-1}^ op \mathcal{T}_{\pi,h}^ op oldsymbol{u}.$$

Assume that either  $\mathbf{u} = \boldsymbol{\theta}_h$  for some  $\boldsymbol{\theta}_h$  which is a valid reward vector as defined in Definition 3.1, or  $\mathbf{u} \in \mathcal{S}^{d-1}$ . In either case, we have that, for any  $s, a, |\mathbf{v}^{\top} \boldsymbol{\phi}(s, a)| \leq 1$ , and  $||\mathbf{v}||_2 \leq \sqrt{d}$ .

*Proof.* By the linear MDP structure (see Proposition 2.3 of Jin et al. (2020)), for any j,

$$Q_{j}^{\pi}(s, a) = \langle \boldsymbol{\phi}(s, a), \boldsymbol{w}_{j}^{\pi} \rangle$$

$$= \langle \boldsymbol{\phi}(s, a), \boldsymbol{\theta}_{j} \rangle + \int V_{j+1}^{\pi}(s') d\boldsymbol{\mu}_{j}(s')^{\top} \boldsymbol{\phi}(s, a)$$

$$= \langle \boldsymbol{\phi}(s, a), \boldsymbol{\theta}_{j} \rangle + \int \langle \boldsymbol{w}_{j+1}^{\pi}, \boldsymbol{\phi}_{j+1, \pi}(s') \rangle d\boldsymbol{\mu}_{j}(s')^{\top} \boldsymbol{\phi}(s, a)$$

$$= \langle \boldsymbol{\phi}(s, a), \boldsymbol{\theta}_j + \mathcal{T}_{\pi, j+1}^{\top} \boldsymbol{w}_{j+1}^{\pi} \rangle$$

so in general,

$$oldsymbol{w}_i^{\pi} = \sum_{h'=i}^H (\prod_{j=i+1}^{h'} \mathcal{T}_{\pi,j}^{ op}) oldsymbol{ heta}_{h'}$$

where we order the product  $\prod_{j=i+1}^{h'} \mathcal{T}_{\pi,j}^{\top} = \mathcal{T}_{\pi,i+1}^{\top} \mathcal{T}_{\pi,i+1}^{\top} \dots \mathcal{T}_{\pi,h'}^{\top}$ .

Case 1:  $u = \theta_h$ . We first consider the case where  $u = \theta_h$  for some  $\theta_h$  which is a valid reward satisfying Definition 3.1. Assume that the reward in our MDP is set such that for  $h' \neq h$ ,  $\theta_{h'} = 0$ . In this case, we then have that

$$oldsymbol{w}_i^\pi = \mathcal{T}_{\pi.i+1}^ op \mathcal{T}_{\pi.i+2}^ op \ldots \mathcal{T}_{\pi.h}^ op oldsymbol{ heta}_h = oldsymbol{v}.$$

In this case, we know that the trajectory rewards are always bounded by 1, so it follows that  $Q_i^{\pi}(s, a) \leq 1$ . Thus,

$$1 \ge Q_i^{\pi}(s, a) = \langle \boldsymbol{\phi}(s, a), \boldsymbol{w}_i^{\pi} \rangle = \langle \boldsymbol{\phi}(s, a), \boldsymbol{v} \rangle$$

and this holds for any s, a. Since Q-values are always positive, it also holds that  $\langle \phi(s, a), v \rangle \geq 0$ . To bound the norm of v, we note that by the Bellman equation and the calculation above,

$$\|\mathbf{v}\|_{2} = \|\mathbf{w}_{i}^{\pi}\|_{2} = \|\boldsymbol{\theta}_{i} + \int V_{i+1}^{\pi}(s') d\boldsymbol{\mu}_{i}(s')\|_{2}$$

$$\leq \|\boldsymbol{\theta}_{i}\|_{2} + \|\int |V_{i+1}^{\pi}(s')| d\boldsymbol{\mu}_{i}(s')\|_{2}$$

$$\leq \|\int |d\boldsymbol{\mu}_{i}(s')|\|_{2}$$

$$< \sqrt{d}$$

where we have used that  $|V_{i+1}^{\pi}(s')| \leq 1$  since the total episode return is at most 1 on our augmented reward function, and the linear MDP assumption.

Case 2:  $u \in S^{d-1}$ . We can repeat the argument above in the case where we only assume  $u \in S^{d-1}$ . Since  $\|\phi(s,a)\|_2 \leq 1$ , it follows that with the reward vector at level h set to u, the reward will still be bounded in [-1,1]. Thus, essentially the same argument can be used, with the slight modification to handle Q-values that are negative.

**Lemma A.9.** The set  $\Omega_h$  is convex and compact.

Proof. Take  $\Lambda_1, \Lambda_2 \in \Omega_h$ . By definition,  $\Lambda_1 = \mathbb{E}_{\pi \sim \omega_1}[\Lambda_{\pi,h}], \Lambda_2 = \mathbb{E}_{\pi \sim \omega_2}[\Lambda_{\pi,h}]$ . It follows that, for any  $t \in [0,1]$ ,  $t\Lambda_1 + (1-t)\Lambda_2 = \mathbb{E}_{\pi \sim t\omega_1 + (1-t)\omega_2}[\Lambda_{\pi,h}]$ . For  $t\omega_1 + (1-t)\omega_2$  the mixture of  $\omega_1$  and  $\omega_2$ . As  $t\omega_1 + (1-t)\omega_2$  is a valid mixture over policies, it follows that  $t\Lambda_1 + (1-t)\Lambda_2 \in \Omega_h$ , which proves convexity.

Compactness follows since  $\|\phi(s, a)\|_2 \le 1$  for all s, a, so  $\|\mathbf{\Lambda}_{\pi, h}\|_{\text{op}} \le 1$ , which implies  $\|\mathbf{\Lambda}\|_{\text{op}} \le 1$  for any  $\mathbf{\Lambda} \in \mathbf{\Omega}_h$ . Furthermore, the set  $\mathbf{\Omega}_h$  is clearly closed, which proves compactness.

# A.3 Constructing the Policy Class

**Lemma A.10** (Lemma B.1 of Jin et al. (2020)). Let  $\mathbf{w}_h^{\pi}$  denote the set of weights such that  $Q_h^{\pi}(s,a) = \langle \boldsymbol{\phi}(s,a), \mathbf{w}_h^{\pi} \rangle$ . Then  $\|\mathbf{w}_h^{\pi}\|_2 \leq 2H\sqrt{d}$ .

**Lemma A.11.** For any  $\delta > 0$  there exists sets of actions  $(\widetilde{\mathcal{A}}_s)_{s \in \mathcal{S}}$ ,  $\widetilde{\mathcal{A}}_s \subseteq \mathcal{A}$ , such that  $|\widetilde{\mathcal{A}}_s| \leq (1 + 8H\sqrt{d}/\delta)^d$  for all s and, for all  $a \in \mathcal{A}$ , s, h, and any  $\pi$ , there exists some  $\widetilde{a} \in \widetilde{\mathcal{A}}_s$  such that

$$|Q_h^{\pi}(s, a) - Q_h^{\pi}(s, \widetilde{a})| \le \delta, \quad |r_h(s, a) - r_h(s, \widetilde{a})| \le \delta.$$

Proof. Let  $\mathcal{N}$  be a  $\delta/(4H\sqrt{d})$  cover of the unit ball. By Lemma A.1 we can bound  $|\mathcal{N}| \leq (1+8H\sqrt{d}/\delta)^d$ . Take any s and let  $\widetilde{\mathcal{A}}_s = \emptyset$ . Then for each  $\phi \in \mathcal{N}$ , choose any a at random from the set  $\{a \in \mathcal{A} : \|\phi(s,a) - \phi\|_2 \leq \delta/2\}$  and set  $\widetilde{\mathcal{A}}_s \leftarrow \widetilde{\mathcal{A}}_s \cup \{a\}$ . With this construction, we claim that for all  $a \in \mathcal{A}$ , there exists some  $\widetilde{a} \in \widetilde{\mathcal{A}}_s$  such that  $\|\phi(s,a) - \phi(s,\widetilde{a})\|_2 \leq \delta/(2H\sqrt{d})$ . To see why this is, note that by construction of  $\mathcal{N}$ , there always exists some  $\phi \in \mathcal{N}$  such that  $\|\phi(s,a) - \phi\|_2 \leq \delta/(4H\sqrt{d})$ . Since  $\widetilde{\mathcal{A}}_s$  will contain some  $\widetilde{a}$  such that  $\|\phi(s,\widetilde{a}) - \phi\|_2 \leq \delta/(4H\sqrt{d})$ , the claim follows by the triangle inequality.

By Lemma A.10, we have that for any  $\pi$ ,  $\|\boldsymbol{w}_h^{\pi}\|_2 \leq 2H\sqrt{d}$ . Take  $a \in \mathcal{A}$  and let  $\widetilde{a} \in \widetilde{\mathcal{A}}_s$  be the action such that  $\|\boldsymbol{\phi}(s,a) - \boldsymbol{\phi}(s,\widetilde{a})\|_2 \leq \delta/(2H\sqrt{d})$ . Then

$$|Q_h^{\pi}(s,a) - Q_h^{\pi}(s,\widetilde{a})| = |\langle \phi(s,a) - \phi(s,\widetilde{a}), w_h^{\pi} \rangle| \le 2H\sqrt{d} \|\phi(s,a) - \phi(s,\widetilde{a})\|_2 \le \delta.$$

The bound on  $|r_h(s, a) - r_h(s, \widetilde{a})|$  follows analogously, since we assume our rewards are linear, and that  $\|\boldsymbol{\theta}_h\|_2 \leq \sqrt{d}$ .

**Definition A.1** (Linear Softmax Policy). We say a policy is a *linear softmax policy* with parameters  $\eta$  and  $\{\boldsymbol{w}_h\}_{h=1}^H$  if it can be written as

$$\pi_h(a|s) = \frac{e^{\eta \langle \phi(s,a), \mathbf{w}_h \rangle}}{\sum_{a' \in A} e^{\eta \langle \phi(s,a'), \mathbf{w}_h \rangle}}$$

for some  $\boldsymbol{w} = \{\boldsymbol{w}_h\}_{h=1}^H$ . We will denote such a policy as  $\pi^{\boldsymbol{w}}$ .

**Definition A.2** (Restricted-Action Linear Softmax Policy). We say a policy is a restricted-action linear softmax policy with parameters  $\eta$ ,  $\{\boldsymbol{w}_h\}_{h=1}^H$ , and  $(\widetilde{\mathcal{A}}_s)_{s\in\mathcal{S}}$  if it can be written as

$$\widetilde{\pi}_h(a|s) = \frac{e^{\eta \langle \phi(s,a), \mathbf{w}_h \rangle} \cdot \mathbb{I}\{a \in \widetilde{\mathcal{A}}_s\}}{\sum_{a' \in \widetilde{\mathcal{A}}_s} e^{\eta \langle \phi(s,a'), \mathbf{w}_h \rangle}}$$

for some  $\boldsymbol{w} = \{\boldsymbol{w}_h\}_{h=1}^H$ . We will denote such a policy as  $\widetilde{\pi}^{\boldsymbol{w}}$ .

**Lemma A.12.** For any restricted-action linear softmax policies  $\pi^{\boldsymbol{w}}$  and  $\pi^{\boldsymbol{u}}$  with identical restricted sets  $(\widetilde{\mathcal{A}}_s)_{s\in\mathcal{S}}$ , we can bound

$$|V_0^{\pi^{\boldsymbol{w}}}(s_1) - V_0^{\pi^{\boldsymbol{u}}}(s_1)| \le 2dH\eta \sum_{h=1}^H \|\boldsymbol{w}_h - \boldsymbol{u}_h\|_2.$$

*Proof.* Note that for any policy  $\pi$ , the value of the policy can be expressed as

$$V_0^{\pi}(s_1) = \sum_{h=1}^H \langle \boldsymbol{\theta}_h, \boldsymbol{\phi}_{\pi,h} \rangle.$$

Thus,

$$|V_0^{\pi^{\boldsymbol{w}}}(s_1) - V_0^{\pi^{\boldsymbol{u}}}(s_1)| \le \sum_{h=1}^H |\langle \boldsymbol{\theta}_h, \boldsymbol{\phi}_{\pi^{\boldsymbol{w}},h} - \boldsymbol{\phi}_{\pi^{\boldsymbol{u}},h} \rangle|.$$

So it suffices to bound  $|\langle \boldsymbol{\theta}_h, \boldsymbol{\phi}_{\pi^{\boldsymbol{w}},h} - \boldsymbol{\phi}_{\pi^{\boldsymbol{u}},h} \rangle|$ . Using the same decomposition as in the proof of Lemma B.2, we have

$$m{\phi}_{\pi^{m{w}},h} - m{\phi}_{\pi^{m{u}},h} = \sum_{i=0}^{h-1} \left(\prod_{j=h-i+1}^{h} \mathcal{T}_{\pi^{m{w}},j}
ight) (\mathcal{T}_{\pi^{m{w}},h-i} - \mathcal{T}_{\pi^{m{u}},h-i}) m{\phi}_{\pi^{m{u}},h-i-1}.$$

By definition,

$$\mathcal{T}_{\pi^{\boldsymbol{w}},h-i} - \mathcal{T}_{\pi^{\boldsymbol{u}},h-i} = \int (\boldsymbol{\phi}_{\pi^{\boldsymbol{w}},h-i}(s) - \boldsymbol{\phi}_{\pi^{\boldsymbol{u}},h-i}(s)) \mathrm{d}\boldsymbol{\mu}_{h-i-1}(s)^{\top}$$

where

$$\phi_{\pi^{\boldsymbol{w}},h-i}(s) = \sum_{a \in \widetilde{\mathcal{A}}_s} \phi(s,a) \pi_{h-i}^{\boldsymbol{w}}(a|s).$$

Now, for  $a \in \widetilde{\mathcal{A}}_s$ ,

$$\nabla_{\boldsymbol{w}_h} \pi_h^{\boldsymbol{w}}(a|s) = \frac{\eta \boldsymbol{\phi}(s,a) e^{\eta \langle \boldsymbol{\phi}(s,a), \boldsymbol{w}_h \rangle} \cdot \sum_{a' \in \widetilde{\mathcal{A}}_s} e^{\eta \langle \boldsymbol{\phi}(s,a'), \boldsymbol{w}_h \rangle} - e^{\eta \langle \boldsymbol{\phi}(s,a), \boldsymbol{w}_h \rangle} \cdot \sum_{a' \in \widetilde{\mathcal{A}}_s} \eta \boldsymbol{\phi}(s,a') e^{\eta \langle \boldsymbol{\phi}(s,a'), \boldsymbol{w}_h \rangle}}{(\sum_{a' \in \widetilde{\mathcal{A}}_s} e^{\eta \langle \boldsymbol{\phi}(s,a'), \boldsymbol{w}_h \rangle})^2}$$

so

$$\|\nabla_{\boldsymbol{w}_h} \pi_h^{\boldsymbol{w}}(a|s)\|_2 \le \frac{2\eta e^{\eta \langle \boldsymbol{\phi}(s,a), \boldsymbol{w}_h \rangle}}{\sum_{a' \in \widetilde{\mathcal{A}}_s} e^{\eta \langle \boldsymbol{\phi}(s,a'), \boldsymbol{w}_h \rangle}}$$

Thus, by the Mean Value Theorem,

$$|\pi_h^{\boldsymbol{w}}(a|s) - \pi_h^{\boldsymbol{u}}(a|s)| \le \frac{2\eta e^{\eta\langle \phi(s,a), \boldsymbol{w}_h \rangle}}{\sum_{a' \in \widetilde{A}} e^{\eta\langle \phi(s,a'), \boldsymbol{w}_h \rangle}} \cdot \|\boldsymbol{w}_h - \boldsymbol{u}_h\|_2$$

SO

$$\|\phi_{\pi^{\boldsymbol{w}},h-i}(s) - \phi_{\pi^{\boldsymbol{u}},h-i}(s)\|_{2} \leq \sum_{a \in \widetilde{\mathcal{A}}_{s}} |\pi_{h-i}^{\boldsymbol{w}}(a|s) - \pi_{h-i}^{\boldsymbol{u}}(a|s)|$$

$$\leq \sum_{a \in \widetilde{\mathcal{A}}_{s}} \frac{2\eta e^{\eta\langle\phi(s,a),\boldsymbol{w}_{h}\rangle}}{\sum_{a' \in \widetilde{\mathcal{A}}_{s}} e^{\eta\langle\phi(s,a'),\boldsymbol{w}_{h}\rangle}} \cdot \|\boldsymbol{w}_{h-i} - \boldsymbol{u}_{h-i}\|_{2}$$

$$\leq 2\eta \|\boldsymbol{w}_{h-1} - \boldsymbol{u}_{h-1}\|_{2}$$

which, with Definition 3.1, implies that

$$\|\mathcal{T}_{\pi^{\boldsymbol{w}},h-i} - \mathcal{T}_{\pi^{\boldsymbol{u}},h-i}\|_{\text{op}} \leq \int \|\boldsymbol{\phi}_{\pi^{\boldsymbol{w}},h-i}(s) - \boldsymbol{\phi}_{\pi^{\boldsymbol{u}},h-i}(s)\|_{2} \|\mathrm{d}\boldsymbol{\mu}_{h-i-1}(s)\|_{2} \leq 2\sqrt{d}\eta \|\boldsymbol{w}_{h-i} - \boldsymbol{u}_{h-i}\|_{2}.$$

By Lemma A.8, we can bound  $\|\boldsymbol{\theta}_h^{\top} \left(\prod_{j=h-i+1}^h \mathcal{T}_{\pi^{\boldsymbol{w}},j}\right)\|_2$ . Thus, returning to the error decomposition given above, we have

$$|V_0^{\pi^{\boldsymbol{w}}}(s_1) - V_0^{\pi^{\boldsymbol{u}}}(s_1)| \leq \sum_{h=1}^{H} \sum_{i=0}^{h-1} \left| \boldsymbol{\theta}_h^{\top} \left( \prod_{j=h-i+1}^{h} \mathcal{T}_{\pi^{\boldsymbol{w}},j} \right) (\mathcal{T}_{\pi^{\boldsymbol{w}},h-i} - \mathcal{T}_{\pi^{\boldsymbol{u}},h-i}) \boldsymbol{\phi}_{\pi^{\boldsymbol{u}},h-i-1} \right|$$

$$\leq \sqrt{d} \sum_{h=1}^{H} \sum_{i=0}^{h-1} \|\mathcal{T}_{\pi^{\boldsymbol{w}},h-i} - \mathcal{T}_{\pi^{\boldsymbol{u}},h-i}\|_{\text{op}} \|\boldsymbol{\phi}_{\pi^{\boldsymbol{u}},h-i-1}\|_{2}$$

$$\leq 2d\eta \sum_{h=1}^{H} \sum_{i=0}^{h-1} \|\boldsymbol{w}_{h-i} - \boldsymbol{u}_{h-i}\|_{2}$$

$$\leq 2dH\eta \sum_{h=1}^{H} \|\boldsymbol{w}_h - \boldsymbol{u}_h\|_{2}.$$

**Lemma A.13.** Let  $\mathbf{w}^*$  denote the weights such that  $Q_h^*(s,a) = \langle \phi(s,a), \mathbf{w}_h^* \rangle$ , and  $\pi^{\mathbf{w}^*}$  the restricted action linear softmax policy with action sets  $(\widetilde{A}_s)_{s \in \mathcal{S}}$  as defined in Lemma A.11 with  $\delta = \frac{\epsilon}{3(3\sqrt{d})^H}$ . Then

$$|V_0^{\pi^{w^*}}(s_1) - V_0^*(s_1)| \le \epsilon$$

as long as  $\eta \ge 2dH \log(1 + 16Hd/\epsilon) \cdot \frac{(3\sqrt{d})^H}{\epsilon}$ .

*Proof.* We prove this by induction. Assume that at step h, for all s, we have  $|V_h^{\star}(s) - V_h^{\pi^{w^{\star}}}(s)| \leq \delta_h$  for some  $\delta_h$ . Then,

$$|Q_{h-1}^{\pi^{w^{\star}}}(s, a) - Q_{h-1}^{\star}(s, a)| = \left| \int (V_h^{\pi^{w^{\star}}}(s') - V_h^{\star}(s')) d\boldsymbol{\mu}_{h-1}(s')^{\top} \boldsymbol{\phi}(s, a) \right|$$

$$\leq \int |V_h^{\pi^{w^{\star}}}(s') - V_h^{\star}(s')| ||d\boldsymbol{\mu}_{h-1}(s')||_2 ||\boldsymbol{\phi}(s, a)||_2$$

$$\leq \sqrt{d}\delta_h$$

where we use the linear MDP assumption in the last inequality. Thus,

$$V_{h-1}^{\pi^{\mathbf{w}^{\star}}}(s) = \frac{\sum_{a \in \widetilde{\mathcal{A}}_s} e^{\eta \langle \phi(s,a), \mathbf{w}_{h-1}^{\star} \rangle} Q_{h-1}^{\pi^{\mathbf{w}^{\star}}}(s,a)}{\sum_{a \in \widetilde{\mathcal{A}}_s} e^{\eta \langle \phi(s,a), \mathbf{w}_{h-1}^{\star} \rangle}}$$

$$= \frac{\sum_{a \in \widetilde{\mathcal{A}}_s} e^{\eta Q_{h-1}^{\star}(s,a)} Q_{h-1}^{\pi^{\mathbf{w}^{\star}}}(s,a)}{\sum_{a \in \widetilde{\mathcal{A}}_s} e^{\eta Q_{h-1}^{\star}(s,a)}}$$

$$\geq \frac{\sum_{a \in \widetilde{\mathcal{A}}_s} e^{\eta Q_{h-1}^{\star}(s,a)} Q_{h-1}^{\star}(s,a)}{\sum_{a \in \widetilde{\mathcal{A}}_s} e^{\eta Q_{h-1}^{\star}(s,a)}} - \sqrt{d}\delta_h.$$

By Lemma A.3, as long as  $\eta \geq \log |\widetilde{\mathcal{A}}_s|/(\sqrt{d}\delta_h)$ , we can lower bound

$$\frac{\sum_{a \in \widetilde{\mathcal{A}}_s} e^{\eta Q_{h-1}^{\star}(s,a)} Q_{h-1}^{\star}(s,a)}{\sum_{a \in \widetilde{\mathcal{A}}_s} e^{\eta Q_{h-1}^{\star}(s,a)}} - \sqrt{d}\delta_h \ge \max_{a \in \widetilde{\mathcal{A}}_s} Q_{h-1}^{\star}(s,a) - 2\sqrt{d}\delta_h.$$

Furthermore, by Lemma A.11 and our choice of  $\widetilde{\mathcal{A}}_s$ , we have

$$\max_{a \in \widetilde{\mathcal{A}}_s} Q_{h-1}^{\star}(s,a) - 2\sqrt{d}\delta_h \ge \max_{a \in \mathcal{A}} Q_{h-1}^{\star}(s,a) - 2\sqrt{d}\delta_h - \frac{\epsilon}{3(3\sqrt{d})^H} = V_{h-1}^{\star}(s) - 2\sqrt{d}\delta_h - \frac{\epsilon}{3(3\sqrt{d})^H}.$$

Define recursively  $\delta_{h-1} = 3\sqrt{d}\delta_h$  and  $\delta_H = \frac{\epsilon}{(3\sqrt{d})^H}$ . Then  $\delta_{h-1} = \frac{\epsilon}{(3\sqrt{d})^{h-1}} \ge \frac{\epsilon}{(3\sqrt{d})^H}$ , so

$$V_{h-1}^{\star}(s) - 2\sqrt{d}\delta_h - \frac{\epsilon}{3(3\sqrt{d})^H} \ge V_{h-1}^{\star}(s) - 2\sqrt{d}\delta_h - \delta_{h-1}/3 = V_{h-1}^{\star}(s) - \delta_{h-1}.$$

So  $|V_h^{\star}(s) - V_h^{\pi^{w^{\star}}}(s)| \leq \delta_{h-1}$  for all s, which proves the inductive step. For the base case, we have

$$V_H^{\pi^{w^*}}(s) - V_H^*(s) = \frac{\sum_{a \in \widetilde{\mathcal{A}}_s} e^{\eta Q_H^*(s,a)} \nu_H(s,a)}{\sum_{a \in \widetilde{\mathcal{A}}_s} e^{\eta Q_H^*(s,a)}} - \max_a \nu_H(s,a)$$

$$\geq \max_{a \in \widetilde{\mathcal{A}}_s} \nu_H(s,a) - \max_a \nu_H(s,a) - \delta_H/2$$

$$\geq -\delta_H$$

where the first inequality holds by Lemma A.3 as long as  $\eta \geq 2 \log |\widetilde{\mathcal{A}}_s|/\delta_H$ , and the second inequality holds by Lemma A.11 and our choice of  $\widetilde{\mathcal{A}}_s$  and  $\delta_H$ . This proves the base case, since  $V_H^{\pi^{w^*}}(s) \leq V_H^{\star}(s)$ .

Recursing this all the way back, we conclude that

$$V_0^{\pi^{w^*}}(s_1) \ge V_0^*(s_1) - \delta_0$$

for  $\delta_0 = (3\sqrt{d})^H \delta_H = \epsilon$ .

For this argument to hold, we must choose  $\eta \geq 2 \log |\widetilde{\mathcal{A}}_s|/\delta_H$  and  $\eta \geq \log |\widetilde{\mathcal{A}}_s|/(\sqrt{d}\delta_h)$  for all s and h. By Lemma A.11 and our choice of  $\widetilde{\mathcal{A}}_s$ , we can bound

$$|\widetilde{\mathcal{A}}_s| \le (1 + 8H\sqrt{d}(2\sqrt{d})^H/\epsilon)^d \le (1 + 16Hd/\epsilon)^{dH}$$

so it suffices that we take  $\eta \geq 2dH \log(1 + 16Hd/\epsilon) \cdot \frac{(3\sqrt{d})^H}{\epsilon}$ .

**Lemma A.14.** Let  $\eta = 2dH \log(1 + 16Hd/\epsilon) \cdot \frac{(3\sqrt{d})^H}{\epsilon}$  and W an  $\frac{\epsilon}{4dH^2\eta}$ -net of  $\mathcal{B}^d(2H\sqrt{d})$ . Let  $\Pi$  denote the set of restricted-action linear softmax policy with vectors  $\mathbf{w} \in W^H$ , parameter  $\eta$ , and action sets  $(\widetilde{\mathcal{A}}_s)_{s \in \mathcal{S}}$  as defined in Lemma A.11 with  $\delta = \frac{\epsilon}{3(3\sqrt{d})^H}$ . Then for any MDP and reward function, there exists some  $\pi \in \Pi$  such that  $|V_0^{\pi} - V_0^{\star}| \leq \epsilon$ , and

$$|\Pi| \le \left(1 + \frac{32H^4d^{5/2}\log(1 + 16Hd/\epsilon)}{\epsilon^2}\right)^{dH^2}.$$

*Proof.* Consider some MDP and reward function, and let  $\{\boldsymbol{w}_h^{\star}\}_{h=1}^{H}$  denote the optimal Q-function linear representation:  $Q_h^{\star}(s,a) = \langle \boldsymbol{\phi}(s,a), \boldsymbol{w}_h^{\star} \rangle$ . Let  $\widetilde{\boldsymbol{w}}$  denote the vector in  $\mathcal{W}^H$  such that  $\sum_{h=1}^{H} \|\boldsymbol{w}_h^{\star} - \widetilde{\boldsymbol{w}}_h\|_2$  is minimized. Then by Lemma A.12 and Lemma A.13, as long as  $\eta \geq 2dH \log(1 + 16Hd/\epsilon) \cdot \frac{(3\sqrt{d})^H}{\epsilon}$ , we have

$$|V_0^{\pi^{\widetilde{w}}}(s_1) - V_0^{\star}(s_1)| \le |V_0^{\pi^{\widetilde{w}}}(s_1) - V_0^{\pi^{w^{\star}}}(s_1)| + |V_0^{\pi^{w^{\star}}}(s_1) - V_0^{\star}(s_1)|$$

26

$$\leq 2dH\eta\sum_{h=1}^{H}\|\boldsymbol{w}_{h}^{\star}-\widetilde{\boldsymbol{w}}_{h}\|_{2}+\epsilon/2.$$

The first conclusion then follows as long as we can find some  $\widetilde{\boldsymbol{w}}$  such that

$$2dH\eta \sum_{h=1}^{H} \|\boldsymbol{w}_h^{\star} - \widetilde{\boldsymbol{w}}_h\|_2 \le \epsilon/2.$$

However, by Lemma A.10, we can bound  $\|\boldsymbol{w}_h^{\star}\|_2 \leq 2H\sqrt{d}$ . Therefore, since  $\mathcal{W}$  is a  $\frac{\epsilon}{4dH^2\eta}$ -net of  $\mathcal{B}^d(2H\sqrt{d})$ , for each h there will exist some  $\widetilde{\boldsymbol{w}}_h \in \mathcal{W}$  such that  $\|\boldsymbol{w}_h^{\star} - \widetilde{\boldsymbol{w}}_h\|_2 \leq \frac{\epsilon}{4dH^2\eta}$ , which implies that we can find  $\widetilde{\boldsymbol{w}} \in \mathcal{W}^d$  such that

$$2dH\eta \sum_{h=1}^{H} \|\boldsymbol{w}_h^{\star} - \widetilde{\boldsymbol{w}}_h\|_2 \le \epsilon/2,$$

which gives the first conclusion.

To bound the size of  $\Pi$ , we apply Lemma A.1 and our choice of  $\eta$  to bound

$$|\mathcal{W}| \le (1 + \frac{16H^3d^{3/2}\eta}{\epsilon})^d \le (1 + \frac{32H^4d^{5/2}\log(1 + 16Hd/\epsilon)}{\epsilon^2})^{dH}.$$

The bound on  $|\Pi|$  follows since  $|\Pi| = |\mathcal{W}|^H$ .

# **B** Policy Elimination

Throughout this section, assuming we have run for some number of episodes K, we let  $(\mathcal{F}_{\tau})_{\tau=1}^{K}$  the filtration on this, with  $\mathcal{F}_{\tau}$  the filtration up to and including episode  $\tau$ . We also let  $\mathcal{F}_{\tau,h}$  denote the filtration on all episodes  $\tau' < \tau$ , and on steps  $h' = 1, \ldots, h$  of episode  $\tau$ .

# **B.1** Estimating Feature-Visitations and Rewards

**Lemma B.1.** Assume that we have collected some data  $\{(s_{h-1,\tau}, a_{h-1,\tau}, s_{h,\tau})\}_{\tau=1}^K$ , where, for each  $\tau'$ ,  $s_{h,\tau'}|\mathcal{F}_{h-1,\tau'}$  is independent of  $\{(s_{h-1,\tau}, a_{h-1,\tau}, s_{h,\tau})\}_{\tau\neq\tau'}$ . Denote  $\phi_{h-1,\tau} = \phi(s_{h-1,\tau}, a_{h-1,\tau})$  and  $\Lambda_{h-1} = \sum_{\tau=1}^K \phi_{h-1,\tau} \phi_{h-1,\tau}^{\top} + \lambda I$ . Fix  $\pi$  and let

$$\widehat{\mathcal{T}}_{\pi,h} = \left(\sum_{ au=1}^K oldsymbol{\phi}_{\pi,h}(s_{h, au})oldsymbol{\phi}_{h-1, au}^ op oldsymbol{\Lambda}_{h-1}^{-1}.$$

Fix  $\mathbf{v} \in \mathbb{R}^d$  satisfying  $|\mathbf{v}^{\top} \boldsymbol{\phi}_{\pi,h}(s)| \leq 1$  for all s and  $\mathbf{u} \in \mathbb{R}^d$ . Then with probability at least  $1 - \delta$ , we can bound

$$|\boldsymbol{v}^{\top}(\mathcal{T}_{\pi,h} - \widehat{\mathcal{T}}_{\pi,h})\boldsymbol{u}| \leq \left(2\sqrt{\log 2/\delta} + \frac{\log 2/\delta}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{h-1})}} + \sqrt{\lambda}\|\mathcal{T}_{\pi,h}^{\top}\boldsymbol{v}\|_{2}\right) \cdot \|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_{h-1}^{-1}}.$$

*Proof.* Let  $\mathfrak{D} = \{(s_{h-1,\tau}, a_{h-1,\tau})\}_{\tau=1}^K$ , our data collected at step h-1. Then by our assumption on the independence of  $s_{h,\tau}$ , we have that  $s_{h,\tau}|\mathcal{F}_{h-1,\tau}$  has the same distribution as  $s_{h,\tau}|(\mathcal{F}_{h-1,\tau},\mathfrak{D})$ . Conditioning on  $\mathfrak{D}$ , the  $\phi_{h-1,\tau}$  vectors are fixed, so  $\Lambda_{h-1}$  is also fixed. Note that

$$\mathcal{T}_{\pi,h} = \int \boldsymbol{\phi}_{\pi,h}(s) \mathrm{d}\boldsymbol{\mu}_{h-1}(s)^{\top}$$

$$= \int \boldsymbol{\phi}_{\pi,h}(s) \mathrm{d}\boldsymbol{\mu}_{h-1}(s)^{\top} \left( \sum_{\tau=1}^{K} \boldsymbol{\phi}_{h-1,\tau} \boldsymbol{\phi}_{h-1,\tau}^{\top} \right) \boldsymbol{\Lambda}_{h-1}^{-1} + \lambda \int \boldsymbol{\phi}_{\pi,h}(s) \mathrm{d}\boldsymbol{\mu}_{h-1}(s)^{\top} \boldsymbol{\Lambda}_{h-1}^{-1}$$

$$= \sum_{\tau=1}^{K} \left( \int \boldsymbol{\phi}_{\pi,h}(s) \mathrm{d}\boldsymbol{\mu}_{h-1}(s)^{\top} \boldsymbol{\phi}_{h-1,\tau} \right) \boldsymbol{\phi}_{h-1,\tau}^{\top} \boldsymbol{\Lambda}_{h-1}^{-1} + \lambda \int \boldsymbol{\phi}_{\pi,h}(s) \mathrm{d}\boldsymbol{\mu}_{h-1}(s)^{\top} \boldsymbol{\Lambda}_{h-1}^{-1}$$

$$= \sum_{\tau=1}^{K} \mathbb{E}[\boldsymbol{\phi}_{\pi,h}(s_{h,\tau}) | \mathcal{F}_{h-1,\tau}] \boldsymbol{\phi}_{h-1,\tau}^{\top} \boldsymbol{\Lambda}_{h-1}^{-1} + \lambda \int \boldsymbol{\phi}_{\pi,h}(s) \mathrm{d}\boldsymbol{\mu}_{h-1}(s)^{\top} \boldsymbol{\Lambda}_{h-1}^{-1}$$

$$= \sum_{\tau=1}^{K} \mathbb{E}[\boldsymbol{\phi}_{\pi,h}(s_{h,\tau}) | \mathcal{F}_{h-1,\tau}] \boldsymbol{\phi}_{h-1,\tau}^{\top} \boldsymbol{\Lambda}_{h-1}^{-1} + \lambda \mathcal{T}_{\pi,h} \boldsymbol{\Lambda}_{h-1}^{-1}$$

SO

$$|\boldsymbol{v}^{\top}(\mathcal{T}_{\pi,h} - \widehat{\mathcal{T}}_{\pi,h})\boldsymbol{u}| \leq \underbrace{\left|\sum_{\tau=1}^{K} \boldsymbol{v}^{\top} \left(\mathbb{E}[\boldsymbol{\phi}_{\pi,h}(s_{h,\tau})|\mathcal{F}_{h-1,\tau}] - \boldsymbol{\phi}_{\pi,h}(s_{h,\tau})\right) \boldsymbol{\phi}_{h-1,\tau}^{\top} \boldsymbol{\Lambda}_{h-1}^{-1} \boldsymbol{u}\right|}_{(a)} + \underbrace{\left|\lambda \boldsymbol{v}^{\top} \mathcal{T}_{\pi,h} \boldsymbol{\Lambda}_{h-1}^{-1} \boldsymbol{u}\right|}_{(b)}.$$

Conditioned on  $\mathfrak{D}$ , (a) is simply the sum of mean 0 random variables, where the  $\tau$ th random variable has magnitude bounded as

$$\begin{aligned} |\boldsymbol{v}^{\top} \left( \mathbb{E}[\phi_{\pi,h}(s_{h,\tau})|\mathcal{F}_{h-1,\tau}] - \phi_{\pi,h}(s_{h,\tau}) \right) \phi_{h-1,\tau}^{\top} \boldsymbol{\Lambda}_{h-1}^{-1} \boldsymbol{u}| &\leq 2|\phi_{h-1,\tau}^{\top} \boldsymbol{\Lambda}_{h-1}^{-1} \boldsymbol{u}| \\ &\leq 2\|\phi_{h-1,\tau}\|_{\boldsymbol{\Lambda}_{h-1}^{-1}} \|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_{h-1}^{-1}} \\ &\leq 2\|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_{h-1}^{-1}} / \sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{h-1})} \end{aligned}$$

Furthermore, the variance of each term in (a) is bounded as

$$\operatorname{Var}\left[\boldsymbol{v}^{\top}\left(\mathbb{E}[\boldsymbol{\phi}_{\pi,h}(s_{h,\tau})|\mathcal{F}_{h-1,\tau}] - \boldsymbol{\phi}_{\pi,h}(s_{h,\tau})\right)\boldsymbol{\phi}_{h-1,\tau}^{\top}\boldsymbol{\Lambda}_{h-1}^{-1}\boldsymbol{u}|\mathcal{F}_{h-1}\right]$$

$$= \mathbb{E}\left[\left(\boldsymbol{v}^{\top}\left(\mathbb{E}[\boldsymbol{\phi}_{\pi,h}(s_{h,\tau})|\mathcal{F}_{h-1,\tau}] - \boldsymbol{\phi}_{\pi,h}(s_{h,\tau})\right)\boldsymbol{\phi}_{h-1,\tau}^{\top}\boldsymbol{\Lambda}_{h-1}^{-1}\boldsymbol{u}\right)^{2}|\mathcal{F}_{h-1}\right]$$

$$\leq \boldsymbol{u}^{\top}\boldsymbol{\Lambda}_{h-1}^{-1}\boldsymbol{\phi}_{h-1,\tau}\boldsymbol{\phi}_{h-1,\tau}^{\top}\boldsymbol{\Lambda}_{h-1}^{-1}\boldsymbol{u}.$$

It follows that, by Bernstein's Inequality, we can bound, with probability at least  $1 - \delta$  conditioned on  $\mathfrak{D}$ :

$$(a) \leq 2\sqrt{\sum_{\tau=1}^{K} \boldsymbol{u}^{\top} \boldsymbol{\Lambda}_{h-1}^{-1} \boldsymbol{\phi}_{h-1,\tau} \boldsymbol{\phi}_{h-1,\tau}^{\top} \boldsymbol{\Lambda}_{h-1}^{-1} \boldsymbol{u} \cdot \log \frac{2}{\delta}} + \frac{2\|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_{h-1}^{-1}}}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{h-1})}} \cdot \log \frac{2}{\delta}$$
$$\leq 2(\sqrt{\log 2/\delta} + \frac{\log 2/\delta}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{h-1})}}) \cdot \|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_{h-1}^{-1}}.$$

In other words,

$$\mathbb{P}\left[(a) \geq 2(\sqrt{\log 2/\delta} + \frac{\log 2/\delta}{\sqrt{\lambda_{\min}(\mathbf{\Lambda}_{h-1})}}) \cdot \|\boldsymbol{u}\|_{\mathbf{\Lambda}_{h-1}^{-1}} |\mathfrak{D}\right] \leq \delta$$

so, by the law of total probability, for any distribution F over  $\mathfrak{D}$ ,

$$\begin{split} & \mathbb{P}\left[(a) \geq 2(\sqrt{\log 2/\delta} + \min\{1, \lambda^{-1}\} \log 2/\delta) \cdot \|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_{h-1}^{-1}}\right] \\ & = \int \mathbb{P}\left[(a) \geq 2(\sqrt{\log 2/\delta} + \min\{1, \lambda^{-1}\} \log 2/\delta) \cdot \|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_{h-1}^{-1}} |\mathfrak{D}\right] \mathrm{d}F(\mathfrak{D}) \\ & \leq \delta \int \mathrm{d}F(\mathfrak{D}) \\ & = \delta. \end{split}$$

We can also bound

$$(b) \leq \sqrt{\lambda} \|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_{h-1}^{-1}} \|\mathcal{T}_{\pi,h}^{\top} \boldsymbol{v}\|_{2}.$$

Combining these gives the result.

**Lemma B.2.** Fix  $\pi$  and let

$$\widehat{\phi}_{\pi,h} = \widehat{\mathcal{T}}_{\pi,h} \widehat{\mathcal{T}}_{\pi,h-1} \dots \widehat{\mathcal{T}}_{\pi,2} \widehat{\mathcal{T}}_{\pi,1} \phi_{\pi,0}.$$

Fix  $\mathbf{u} \in \mathcal{S}^{d-1}$  or  $\mathbf{u}$  a valid reward vector as defined by Definition 3.1. Then with probability at least  $1 - \delta$ :

$$|\langle \boldsymbol{u}, \boldsymbol{\phi}_{\pi,h} - \widehat{\boldsymbol{\phi}}_{\pi,h} \rangle| \leq \sum_{i=1}^{h-1} \left( 2\sqrt{\log \frac{2H}{\delta}} + \frac{\log \frac{2H}{\delta}}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_i)}} + \sqrt{d\lambda} \right) \cdot \|\widehat{\boldsymbol{\phi}}_{\pi,i}\|_{\boldsymbol{\Lambda}_i^{-1}}.$$

Proof. Note that

$$\begin{split} \phi_{\pi,h} - \widehat{\phi}_{\pi,h} &= \mathcal{T}_{\pi,h} \phi_{\pi,h-1} - \widehat{\mathcal{T}}_{\pi,h} \widehat{\phi}_{\pi,h-1} \\ &= \mathcal{T}_{\pi,h} (\phi_{\pi,h-1} - \widehat{\phi}_{\pi,h-1}) + (\mathcal{T}_{\pi,h} - \widehat{\mathcal{T}}_{\pi,h}) \widehat{\phi}_{\pi,h-1}. \end{split}$$

Thus, unrolling this all the way back, we get

$$\phi_{\pi,h}-\widehat{\phi}_{\pi,h}=\sum_{i=1}^{h-1}\left(\prod_{j=h-i+1}^{h}\mathcal{T}_{\pi,j}
ight)(\mathcal{T}_{\pi,h-i}-\widehat{\mathcal{T}}_{\pi,h-i})\widehat{\phi}_{\pi,h-i-1}$$

where we order the product  $\prod_{j=h-i+1}^h \mathcal{T}_{\pi,j} = \mathcal{T}_{\pi,h} \mathcal{T}_{\pi,h-1} \dots \mathcal{T}_{\pi,h-i+1}$ . It follows that

$$|\langle oldsymbol{u}, oldsymbol{\phi}_{\pi,h} - \widehat{oldsymbol{\phi}}_{\pi,h} 
angle| \leq \sum_{i=1}^{h-1} \left| oldsymbol{u}^ op \left( \prod_{j=h-i+1}^h \mathcal{T}_{\pi,j} 
ight) (\mathcal{T}_{\pi,h-i} - \widehat{\mathcal{T}}_{\pi,h-i}) \widehat{oldsymbol{\phi}}_{\pi,h-i-1} 
ight|.$$

Denote  $v_i := u^{\top} \left( \prod_{j=h-i+1}^{h} \mathcal{T}_{\pi,j} \right)$ . By Lemma A.8 and our assumption on u, we can bound  $\|v_i\|_2 \leq \sqrt{d}$  and also have that for all  $s, a, |v_i^{\top} \phi(s, a)| \leq 1$ , which implies

$$|\boldsymbol{v}_i^{\top} \boldsymbol{\phi}_{\pi,j}(s)| = \left| \sum_{a \in \mathcal{A}} \boldsymbol{v}_i^{\top} \boldsymbol{\phi}(s, a) \pi_h(a|s) \right| \leq \sum_{a \in \mathcal{A}} \pi_h(a|s) = 1.$$

We can therefore apply Lemma B.1 to get that, with probability at least  $1 - \delta$ , for all i,

$$\left| \boldsymbol{v}_i^\top (\mathcal{T}_{\pi,h-i} - \widehat{\mathcal{T}}_{\pi,h-i}) \widehat{\boldsymbol{\phi}}_{\pi,h-i-1} \right| \leq \left( 2 \sqrt{\log \frac{2H}{\delta}} + \frac{\log \frac{2H}{\delta}}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{h-i-1})}} + \sqrt{\lambda} \| \mathcal{T}_{\pi,h}^\top \boldsymbol{v}_i \|_2 \right) \cdot \| \widehat{\boldsymbol{\phi}}_{\pi,h-i-1} \|_{\boldsymbol{\Lambda}_{h-i-1}^{-1}} \|_{\boldsymbol{\Lambda}_{h-$$

By Lemma A.8, the definition of  $v_i$ , and our assumption on u, we can bound  $\|\mathcal{T}_{\pi,h}^{\top}v_i\|_2 \leq \sqrt{d}$ . Summing over i proves the result.

**Lemma B.3.** With probability at least  $1 - \delta$ :

$$\|\widehat{\boldsymbol{\phi}}_{\pi,h} - \boldsymbol{\phi}_{\pi,h}\|_{2} \leq d \sum_{h'=1}^{h-1} \left( 2\sqrt{\log \frac{2Hd}{\delta}} + \frac{\log \frac{2Hd}{\delta}}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{h'})}} + \sqrt{d\lambda} \right) \cdot \|\widehat{\boldsymbol{\phi}}_{\pi,h'}\|_{\boldsymbol{\Lambda}_{h'}^{-1}}.$$

*Proof.* We have:

$$\|\widehat{\phi}_{\pi,h} - \phi_{\pi,h}\|_2 \leq \|\widehat{\phi}_{\pi,h} - \phi_{\pi,h}\|_1 = \sum_{i=1}^d |[\widehat{\phi}_{\pi,h}]_i - [\phi_{\pi,h}]_i| = \sum_{i=1}^d |\langle e_i, \widehat{\phi}_{\pi,h} - \phi_{\pi,h} \rangle|.$$

Since  $e_i \in \mathcal{S}^{d-1}$ , we can apply Lemma B.2 to bound, with probability  $1 - \delta/d$ ,

$$|\langle \boldsymbol{e}_i, \widehat{\boldsymbol{\phi}}_{\pi,h} - \boldsymbol{\phi}_{\pi,h} \rangle| \leq \sum_{h'=0}^{h-1} \left( 2\sqrt{\log \frac{2Hd}{\delta}} + \frac{\log \frac{2Hd}{\delta}}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{h'})}} + \sqrt{d\lambda} \right) \cdot \|\widehat{\boldsymbol{\phi}}_{\pi,h'}\|_{\boldsymbol{\Lambda}_{h'}^{-1}}.$$

Summing over i gives the result.

**Lemma B.4.** Assume we have collected data  $\{\phi(s_{h,\tau}, a_{h,\tau}), r_h(s_{h,\tau}, a_{h,\tau})\}_{\tau=1}^K$  and that for each  $\tau'$ ,  $r_h(s_{h,\tau'}, a_{h,\tau'})|(s_{h,\tau'}, a_{h,\tau'})|$  is independent of  $\{(s_{h,\tau}, a_{h,\tau})\}_{\tau\neq\tau'}$ . Let

$$\widehat{\boldsymbol{\theta}}_h = \operatorname*{arg\,min}_{\boldsymbol{\theta}} \sum_{\tau=1}^K (r_{h,\tau} - \langle \boldsymbol{\phi}_{h,\tau}, \boldsymbol{\theta} \rangle)^2 + \lambda \|\boldsymbol{\theta}\|_2^2$$

and fix  $\mathbf{u} \in \mathbb{R}^d$  that is independent of  $\{\phi(s_{h,\tau}, a_{h,\tau}), r_h(s_{h,\tau}, a_{h,\tau})\}_{\tau=1}^K$ . Then with probability at least  $1 - \delta$ :

$$|\langle \boldsymbol{u}, \widehat{\boldsymbol{\theta}}_h - \boldsymbol{\theta}_h \rangle| \le \left(\sqrt{\log 2/\delta} + \frac{\log 2/\delta}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_h)}} + \sqrt{d\lambda}\right) \cdot \|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_h^{-1}}.$$

*Proof.* Let  $\mathfrak{D} = \{(s_{h,\tau}, a_{h,\tau})\}_{\tau=1}^K$ . Then by our assumption on the independence of  $r_{h,\tau}$ , we have that  $r_{h,\tau}|(s_{h,\tau}, a_{h,\tau})$  has the same distribution as  $r_{h,\tau}|\mathfrak{D}$ . Conditioning on  $\mathfrak{D}$ , the  $\phi_{h,\tau}$  vectors are fixed, so  $\Lambda_h$  is also fixed.

By construction we have

$$\widehat{oldsymbol{ heta}}_h = oldsymbol{\Lambda}_h^{-1} \sum_{ au=1}^K oldsymbol{\phi}_{h, au} r_{h, au}.$$

Furthermore:

$$oldsymbol{ heta}_h = oldsymbol{\Lambda}_h^{-1} oldsymbol{\Lambda}_h oldsymbol{ heta}_h = oldsymbol{\Lambda}_h^{-1} \sum_{ au=1}^K oldsymbol{\phi}_{h, au} \mathbb{E}[r_{h, au}|\mathcal{F}_{h-1, au}] + \lambda oldsymbol{\Lambda}_h^{-1} oldsymbol{ heta}_h.$$

Thus,

$$|\langle \boldsymbol{u}, \widehat{\boldsymbol{\theta}}_h - \boldsymbol{\theta}_h \rangle| \leq \underbrace{\left| \sum_{\tau=1}^K \boldsymbol{u}^\top \boldsymbol{\Lambda}_h^{-1} \boldsymbol{\phi}_{h,\tau} (r_{h,\tau} - \mathbb{E}[r_{h,\tau} | \mathcal{F}_{h-1,\tau}]) \right|}_{(b)} + \underbrace{\left| \lambda \boldsymbol{u}^\top \boldsymbol{\Lambda}_h^{-1} \boldsymbol{\theta}_h \right|}_{(b)}.$$

Since  $R_{h,\tau} \in [0,1]$  almost surely, we can bound

$$|\boldsymbol{u}^{\top}\boldsymbol{\Lambda}_{h}^{-1}\boldsymbol{\phi}_{h,\tau}(r_{h,\tau} - \mathbb{E}[r_{h,\tau}|\mathcal{F}_{h-1,\tau}])| \leq \|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_{h}^{-1}}\|\boldsymbol{\phi}_{h,\tau}\|_{\boldsymbol{\Lambda}_{h}^{-1}} \leq \|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_{h}^{-1}}/\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{h})}.$$

Furthermore, we can bound

$$\operatorname{Var}\left[\boldsymbol{u}^{\top}\boldsymbol{\Lambda}_{h}^{-1}\boldsymbol{\phi}_{h,\tau}(r_{h,\tau} - \mathbb{E}[r_{h,\tau}|\mathcal{F}_{h-1,\tau}])|\mathfrak{D}\right]$$

$$= \mathbb{E}\left[(\boldsymbol{u}^{\top}\boldsymbol{\Lambda}_{h}^{-1}\boldsymbol{\phi}_{h,\tau}(r_{h,\tau} - \mathbb{E}[r_{h,\tau}|\mathcal{F}_{h-1,\tau}]))^{2}|\mathfrak{D}\right]$$

$$\leq \boldsymbol{u}^{\top}\boldsymbol{\Lambda}_{h}^{-1}\boldsymbol{\phi}_{h,\tau}\boldsymbol{\phi}_{h,\tau}^{\top}\boldsymbol{\Lambda}_{h}^{-1}\boldsymbol{u}.$$

By Bernstein's inequality, we then have, with probability at least  $1 - \delta$  conditioned on  $\mathfrak{D}$ :

$$(a) \leq \sqrt{\sum_{\tau=1}^{K} \boldsymbol{u}^{\top} \boldsymbol{\Lambda}_{h}^{-1} \boldsymbol{\phi}_{h,\tau} \boldsymbol{\phi}_{h,\tau}^{\top} \boldsymbol{\Lambda}_{h}^{-1} \boldsymbol{u} \cdot \log 2/\delta} + \frac{\|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_{h}^{-1}} \cdot \log 2/\delta}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{h})}} \\ \leq (\sqrt{\log 2/\delta} + \frac{\log 2/\delta}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{h})}}) \cdot \|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_{h}^{-1}}.$$

Applying the Law of Total Probability as in Lemma B.1, we obtain

$$\mathbb{P}\left[(a) \geq (\sqrt{\log 2/\delta} + \frac{\log 2/\delta}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_h)}}) \cdot \|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_h^{-1}}\right] \leq \delta.$$

By Definition 3.1, we can also bound

$$(b) \leq \sqrt{\lambda} \|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_h^{-1}} \|\boldsymbol{\theta}_h\|_2 \leq \sqrt{d\lambda} \|\boldsymbol{u}\|_{\boldsymbol{\Lambda}_h^{-1}}.$$

Combining these proves the result.

# B.2 Correctness and Sample Complexity of Pedel

**Lemma B.5.** Let  $\mathcal{E}_{est}^{\ell,h}$  denote the event on which, for all  $\pi \in \Pi_{\ell}$ :

$$\begin{split} |\langle \boldsymbol{\theta}_{h+1}, \widehat{\boldsymbol{\phi}}_{\pi,h+1}^{\ell} - \boldsymbol{\phi}_{\pi,h+1} \rangle| &\leq \sum_{i=1}^{h} \left( 3\sqrt{\log \frac{4H^{2}|\Pi_{\ell}|\ell^{2}}{\delta}} + \frac{\log \frac{4H^{2}|\Pi_{\ell}|\ell^{2}}{\delta}}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{i,\ell})}} \right) \cdot \|\widehat{\boldsymbol{\phi}}_{\pi,i}^{\ell}\|_{\boldsymbol{\Lambda}_{i,\ell}^{-1}}, \\ \|\widehat{\boldsymbol{\phi}}_{\pi,h+1}^{\ell} - \boldsymbol{\phi}_{\pi,h+1}\|_{2} &\leq d\sum_{i=1}^{h} \left( 3\sqrt{\log \frac{4H^{2}d|\Pi_{\ell}|\ell^{2}}{\delta}} + \frac{\log \frac{4H^{2}d|\Pi_{\ell}|\ell^{2}}{\delta}}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{i,\ell})}} \right) \cdot \|\widehat{\boldsymbol{\phi}}_{\pi,i}^{\ell}\|_{\boldsymbol{\Lambda}_{i,\ell}^{-1}}, \\ |\langle \widehat{\boldsymbol{\phi}}_{\pi,h}^{\ell}, \widehat{\boldsymbol{\theta}}_{h} - \boldsymbol{\theta}_{h} \rangle| &\leq \left( 2\sqrt{\log \frac{4H^{2}|\Pi_{\ell}|\ell^{2}}{\delta}} + \frac{\log \frac{4H^{2}|\Pi_{\ell}|\ell^{2}}{\delta}}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{h,\ell})}} \right) \cdot \|\widehat{\boldsymbol{\phi}}_{\pi,h}^{\ell}\|_{\boldsymbol{\Lambda}_{h,\ell}^{-1}}. \end{split}$$

Then  $\mathbb{P}[(\mathcal{E}_{\mathrm{est}}^{\ell,h})^c] \leq \frac{\delta}{2H\ell^2}$ .

## Algorithm 3 Policy Learning via Experiment Design in Linear MDPs (Pedel, full version)

- 1: **input:** tolerance  $\epsilon$ , confidence  $\delta$ , policy set  $\Pi$
- 2:  $\ell_0 \leftarrow \lceil \log_2 \frac{d^{3/2}}{H} \rceil$ ,  $\Pi_{\ell_0} \leftarrow \Pi$ ,  $\widehat{\phi}_{\pi,1}^1 \leftarrow \mathbb{E}_{a \sim \pi_1(\cdot|s_1)}[\phi(s_1, a)], \forall \pi \in \Pi$
- 3: **for**  $\ell = \ell_0, \ell_0 + 1, \dots, \lceil \log \frac{4}{\epsilon} \rceil$  **do**
- 4:  $\epsilon_{\ell} \leftarrow 2^{-\ell}, \ \beta_{\ell} \leftarrow 64H^4 \log \frac{4H^2|\Pi_{\ell}|\ell^2}{\delta}$
- 5: **for** h = 1, 2, ..., H **do**
- 6: Run procedure described in Theorem 9 with parameters

$$\epsilon_{\mathrm{exp}} \leftarrow \frac{\epsilon_{\ell}^2}{\beta_{\ell}}, \quad \delta \leftarrow \frac{\delta}{2H\ell^2}, \quad \underline{\lambda} \leftarrow \log \frac{4H^2|\Pi_{\ell}|\ell^2}{\delta}, \quad \Phi \leftarrow \Phi_{h,\ell} := \{\widehat{\phi}_{\pi,h}^{\ell} : \pi \in \Pi_{\ell}\}$$

and denote returned data as  $\{(s_{h,\tau},a_{h,\tau},r_{h,\tau},s_{h+1,\tau})\}_{\tau=1}^{K_{h,\ell}}$ , for  $K_{h,\ell}$  total number of episodes run , and covariates

$$\boldsymbol{\Lambda}_{h,\ell} \leftarrow \sum_{\tau=1}^{K_{h,\ell}} \boldsymbol{\phi}(s_{h,\tau}, a_{h,\tau}) \boldsymbol{\phi}(s_{h,\tau}, a_{h,\tau})^\top + 1/d \cdot I$$

- 7:  $\mathbf{for} \ \pi \in \Pi_\ell \ \mathbf{do}$  // Estimate feature-visitations for active policies
- 8:  $\widehat{\boldsymbol{\phi}}_{\pi,h+1}^{\ell} \leftarrow \left(\sum_{\tau=1}^{K_{h,\ell}} \boldsymbol{\phi}_{\pi,h+1}(s_{h+1,\tau}) \boldsymbol{\phi}_{h,\tau}^{\top} \boldsymbol{\Lambda}_{h,\ell}^{-1}\right) \widehat{\boldsymbol{\phi}}_{\pi,h}^{\ell}$
- 9:  $\widehat{m{ heta}}_h^\ell \leftarrow m{\Lambda}_{h,\ell}^{-1} \sum_{ au=1}^{\hat{K_{h,\ell}}} m{\phi}_{h, au} r_{h, au}$  // Estimate reward vectors
- 10: // Remove provably suboptimal policies from active policy set

$$\Pi_{\ell+1} \leftarrow \Pi_{\ell} \setminus \left\{ \pi \in \Pi_{\ell} : \widehat{V}_0^{\pi} < \sup_{\pi' \in \Pi_{\ell}} \widehat{V}_0^{\pi'} - 2\epsilon_{\ell} \right\} \quad \text{for} \quad \widehat{V}_0^{\pi} := \sum_{h=1}^{H} \langle \widehat{\boldsymbol{\phi}}_{\pi,h}^{\ell}, \widehat{\boldsymbol{\theta}}_{h}^{\ell} \rangle$$

- 11: if  $|\Pi_{\ell+1}| = 1$  then return  $\pi \in \Pi_{\ell+1}$
- 12: **return** any  $\pi \in \Pi_{\ell+1}$

Proof. Note that the data collection procedure outlined in Theorem 9 collects data that satisfies the independence requirement of Lemma B.1 and Lemma B.4, since Theorem 9 operates on the h-truncated-horizon MDP defined with respect to our original MDP (see Definition C.1 and following discussion), so by construction the data obtained at step h is independent of  $s_{h+1}$  and  $r_h(s_h, a_h)$ . Note also that  $\hat{\phi}_{\pi,h}^{\ell}$  is independent of  $\{r_{h,\tau}^{\ell}\}_{\tau=1}^{K_{h,\ell}}|\{(s_{h,\tau}, a_{h,\tau})\}_{\tau=1}^{K_{h,\ell}}$ , since we construct  $\hat{\phi}_{\pi,h}^{\ell}$  using only observations taken at step h-1.

The result follows by Lemma B.2, Lemma B.3, and Lemma B.4, and setting  $\lambda = 1/d$ .

**Lemma B.6.** Let  $\mathcal{E}_{\exp}^{\ell,h}$  denote the event on which:

• The exploration procedure on Line 6 terminates after running for at most

$$C \cdot \frac{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\boldsymbol{\phi} \in \Phi_{\ell,h}} \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2}{\epsilon_{\ell}^2 / \beta_{\ell}} + \operatorname{poly}\left(d, H, \log \frac{\ell^2}{\delta}, \frac{1}{\lambda_{\min}^{\star}}, \log |\Pi_{\ell}|\right)$$

episodes.

• The covariates returned by Line 6 for any  $(h, \ell)$ ,  $\Lambda_{h,\ell}$ , satisfy

$$\max_{\boldsymbol{\phi} \in \Phi_{\ell,h}} \|\boldsymbol{\phi}\|_{\boldsymbol{\Lambda}_{h,\ell}^{-1}}^2 \leq \frac{\epsilon_\ell^2}{\beta_\ell}, \qquad \lambda_{\min}(\boldsymbol{\Lambda}_{h,\ell}) \geq \log \frac{4H^2 |\Pi_\ell| \ell^2}{\delta}.$$

Then 
$$\mathbb{P}[(\mathcal{E}_{\exp}^{\ell,h})^c \cap \mathcal{E}_{\operatorname{est}}^{\ell,h-1} \cap (\cap_{i=1}^{h-1} \mathcal{E}_{\exp}^{\ell,i})] \leq \frac{\delta}{2H\ell^2}$$

*Proof.* By Lemma B.7, on the event  $\mathcal{E}_{\text{est}}^{\ell,h-1} \cap (\bigcap_{i=1}^{h-1} \mathcal{E}_{\exp}^{\ell,i})$  we can bound  $\|\widehat{\phi}_{\pi,h}^{\ell} - \phi_{\pi,h}\|_2 \leq d\epsilon_{\ell}/2H$ . By Lemma A.6, we can lower bound  $\|\phi_{\pi,h}\|_2 \geq 1/\sqrt{d}$ . By the reverse triangle inequality,

$$\|\widehat{\phi}_{\pi,h}^{\ell}\|_{2} \ge \|\phi_{\pi,h}\|_{2} - \|\widehat{\phi}_{\pi,h}^{\ell} - \phi_{\pi,h}\|_{2} \ge 1/\sqrt{d} - d\epsilon_{\ell}/2H.$$

It follows that as long as  $\epsilon_{\ell} \leq H/d^{3/2}$ , that we can lower bound  $\|\widehat{\phi}_{\pi,h}^{\ell}\|_2 \geq 1/(2\sqrt{d})$ . Since we start  $\ell$  at  $\ell = \lceil \log_2 \frac{d^{3/2}}{H} \rceil$ , we will have that  $\epsilon_{\ell} = 2^{-\ell} \leq H/d^{3/2}$ .

The result then follows by applying Theorem 9 with our chosen parameters and  $\gamma_{\Phi} \leftarrow 1/(2\sqrt{d})$ .

**Lemma B.7.** On the event  $\mathcal{E}_{\mathrm{est}}^{\ell,h} \cap (\cap_{i=1}^h \mathcal{E}_{\mathrm{exp}}^{\ell,i})$ , for all  $\pi \in \Pi_{\ell}$ :

$$\begin{aligned} |\langle \boldsymbol{\theta}_{h+1}, \widehat{\boldsymbol{\phi}}_{\pi,h+1}^{\ell} - \boldsymbol{\phi}_{\pi,h+1} \rangle| &\leq \epsilon_{\ell}/2H, \\ \|\widehat{\boldsymbol{\phi}}_{\pi,h+1}^{\ell} - \boldsymbol{\phi}_{\pi,h+1}\|_{2} &\leq d\epsilon_{\ell}/2H, \\ |\langle \widehat{\boldsymbol{\phi}}_{\pi,h}^{\ell}, \widehat{\boldsymbol{\theta}}_{h} - \boldsymbol{\theta}_{h} \rangle| &\leq \epsilon_{\ell}/2H. \end{aligned}$$

*Proof.* On  $\mathcal{E}_{\exp}^{\ell,i}$ , we can lower bound

$$\lambda_{\min}(\mathbf{\Lambda}_{i,\ell}) \ge \log \frac{4H^2|\Pi_{\ell}|\ell^2}{\delta}$$

which implies

$$3\sqrt{\log\frac{4H^2|\Pi_\ell|\ell^2}{\delta}} + \frac{\log\frac{4H^2|\Pi_\ell|\ell^2}{\delta}}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{i,\ell})}} \leq 4\sqrt{\log\frac{4H^2|\Pi_\ell|\ell^2}{\delta}}.$$

Furthermore, on  $\mathcal{E}_{\exp}^{\ell,i}$ ,  $\|\widehat{\phi}_{\pi,i}^{\ell}\|_{\Lambda_{i,\ell}^{-1}} \leq \frac{\epsilon_{\ell}}{\sqrt{\beta_{\ell}}}$ . Since  $\beta_{\ell} = 64H^4 \log \frac{4H^2|\Pi_{\ell}|\ell^2}{\delta}$ , on  $\mathcal{E}_{\mathrm{est}}^{\ell,h}$ , we can then upper bound

$$\begin{aligned} |\langle \boldsymbol{\theta}_{h+1}, \widehat{\boldsymbol{\phi}}_{\pi,h+1}^{\ell} - \boldsymbol{\phi}_{\pi,h+1} \rangle| &\leq \sum_{i=1}^{h} \left( 3\sqrt{\log \frac{4H^{2}|\Pi_{\ell}|\ell^{2}}{\delta}} + \frac{\log \frac{4H^{2}|\Pi_{\ell}|\ell^{2}}{\delta}}{\sqrt{\lambda_{\min}(\boldsymbol{\Lambda}_{i,\ell})}} \right) \cdot \|\widehat{\boldsymbol{\phi}}_{\pi,i}^{\ell}\|_{\boldsymbol{\Lambda}_{i,\ell}^{-1}} \\ &\leq H4\sqrt{\log \frac{4H^{2}|\Pi_{\ell}|\ell^{2}}{\delta}} \frac{\epsilon_{\ell}}{\sqrt{\beta_{\ell}}} \\ &\leq \epsilon_{\ell}/2H. \end{aligned}$$

The same calculation gives the bounds on  $\|\widehat{\phi}_{\pi,h}^{\ell} - \phi_{\pi,h}\|_2$  and  $|\langle \widehat{\phi}_{\pi,h}^{\ell}, \widehat{\theta}_h - \theta_h \rangle|$ .

**Lemma B.8.** Define  $\mathcal{E}_{\exp} = \cap_{\ell} \cap_{h} \mathcal{E}_{\exp}^{\ell,h}$  and  $\mathcal{E}_{est} = \cap_{\ell} \cap_{h} \mathcal{E}_{est}^{\ell,h}$ . Then  $\mathbb{P}[\mathcal{E}_{est} \cap \mathcal{E}_{\exp}] \geq 1 - 2\delta$  and on  $\mathcal{E}_{est} \cap \mathcal{E}_{exp}$ , for all  $h, \ell$ , and  $\pi \in \Pi_{\ell}$ ,

$$\begin{aligned} |\langle \boldsymbol{\theta}_{h+1}, \widehat{\boldsymbol{\phi}}_{\pi,h+1}^{\ell} - \boldsymbol{\phi}_{\pi,h+1} \rangle| &\leq \epsilon_{\ell}/2H, \\ \|\widehat{\boldsymbol{\phi}}_{\pi,h+1}^{\ell} - \boldsymbol{\phi}_{\pi,h+1}\|_{2} &\leq d\epsilon_{\ell}/2H, \\ |\langle \widehat{\boldsymbol{\phi}}_{\pi,h}^{\ell}, \widehat{\boldsymbol{\theta}}_{h} - \boldsymbol{\theta}_{h} \rangle| &\leq \epsilon_{\ell}/2H. \end{aligned}$$

Proof. Clearly,

$$\begin{split} \mathcal{E}_{\mathrm{est}}^{c} \cup \mathcal{E}_{\mathrm{exp}}^{c} &= \bigcup_{\ell=\ell_{0}}^{\lceil \log 4/\epsilon \rceil} \bigcup_{h=1}^{H} ((\mathcal{E}_{\mathrm{est}}^{\ell,h})^{c} \cup (\mathcal{E}_{\mathrm{exp}}^{\ell,h})^{c}) \\ &= \bigcup_{\ell=\ell_{0}}^{\lceil \log 4/\epsilon \rceil} \bigcup_{h=1}^{H} (\mathcal{E}_{\mathrm{est}}^{\ell,h})^{c} \backslash \left( (\mathcal{E}_{\mathrm{est}}^{\ell,h-1})^{c} \cup (\cup_{i=1}^{h-1} (\mathcal{E}_{\mathrm{exp}}^{\ell,i})^{c}) \right) \cup \bigcup_{\ell=\ell_{0}}^{\lceil \log 4/\epsilon \rceil} \bigcup_{h=1}^{H} (\mathcal{E}_{\mathrm{exp}}^{\ell,h})^{c} \\ &= \bigcup_{\ell=\ell_{0}}^{\lceil \log 4/\epsilon \rceil} \bigcup_{h=1}^{H} (\mathcal{E}_{\mathrm{est}}^{\ell,h})^{c} \cap \left( \mathcal{E}_{\mathrm{est}}^{\ell,h-1} \cap (\cup_{i=1}^{h-1} \mathcal{E}_{\mathrm{exp}}^{\ell,i}) \right) \cup \bigcup_{\ell=\ell_{0}}^{\lceil \log 4/\epsilon \rceil} \bigcup_{h=1}^{H} (\mathcal{E}_{\mathrm{exp}}^{\ell,h})^{c}. \end{split}$$

The first conclusion follows by Lemma B.5, Lemma B.5, and since we can bound

$$\sum_{\ell} \sum_{h=1}^{H} 2 \cdot \frac{\delta}{2H\ell^2} \le \frac{\pi^2}{6} \delta \le 2\delta.$$

The second conclusion follows by Lemma B.7.

**Lemma B.9.** On the event  $\mathcal{E}_{est} \cap \mathcal{E}_{exp}$ , for all  $\ell > \ell_0$ , every policy  $\pi \in \Pi_{\ell}$  satisfies  $V_0^{\star}(\Pi) - V_0^{\pi} \leq 4\epsilon_{\ell}$  and  $\widetilde{\pi}^{\star} \in \Pi_{\ell}$ , for  $\widetilde{\pi}^{\star} = \arg \max_{\pi \in \Pi} V_0^{\pi}$ .

*Proof.* The value of a policy  $\pi$  is given by

$$\sum_{h=1}^{H} \langle oldsymbol{ heta}_h, oldsymbol{\phi}_{\pi,h} 
angle.$$

By Lemma B.8, for all  $\pi \in \Pi_{\ell}$  we can bound

$$|\langle \widehat{\boldsymbol{\theta}}_h, \widehat{\phi}_{\pi,h}^\ell \rangle - \langle \boldsymbol{\theta}_h, \boldsymbol{\phi}_{\pi,h} \rangle| \leq |\langle \widehat{\boldsymbol{\theta}}_h - \boldsymbol{\theta}_h, \widehat{\phi}_{\pi,h}^\ell \rangle| + |\langle \boldsymbol{\theta}_h, \widehat{\phi}_{\pi,h}^\ell - \boldsymbol{\phi}_{\pi,h} \rangle| \leq \epsilon_\ell / 2H + \epsilon_\ell / 2H = \epsilon_\ell / H.$$

Thus,

$$\left| \sum_{h=1}^{H} \langle \widehat{\boldsymbol{\theta}}_h^{\ell}, \widehat{\boldsymbol{\phi}}_{\pi,h}^{\ell} \rangle - \sum_{h=1}^{H} \langle \boldsymbol{\theta}_h, \boldsymbol{\phi}_{\pi,h} \rangle \right| \leq \epsilon_{\ell}.$$

We will only include  $\pi \in \Pi_{\ell+1}$  if  $\pi \in \Pi_{\ell}$  and

$$\sum_{h=1}^{H} \langle \widehat{\phi}_{\pi,h}^{\ell}, \widehat{\theta}_{h}^{\ell} \rangle \ge \sup_{\pi' \in \Pi_{\ell}} \sum_{h=1}^{H} \langle \widehat{\phi}_{\pi',h}^{\ell}, \widehat{\theta}_{h}^{\ell} \rangle - 2\epsilon_{\ell}.$$

Using the estimation error given above, this implies that for any  $\pi \in \Pi_{\ell}$ ,

$$V_0^{\pi} = \sum_{h=1}^{H} \langle \boldsymbol{\theta}_h, \boldsymbol{\phi}_{\pi,h} \rangle \ge \sup_{\pi' \in \Pi_{\ell}} \sum_{h=1}^{H} \langle \boldsymbol{\theta}_h, \boldsymbol{\phi}_{\pi',h} \rangle - 4\epsilon_{\ell} = \sup_{\pi' \in \Pi_{\ell}} V_0^{\pi'} - 4\epsilon_{\ell}.$$

Both claims then follow if we can show  $\widetilde{\pi}^*$  is always contained in the active set. Assume that  $\widetilde{\pi}^* \in \Pi_{\ell}$ . Then

$$\sum_{h=1}^{H} \langle \widehat{\phi}_{\widetilde{\pi}^{\star},h}^{\ell}, \widehat{\theta}_{h}^{\ell} \rangle \geq V_{0}^{\widetilde{\pi}^{\star}} - \epsilon_{\ell}, \quad \sup_{\pi' \in \Pi_{\ell}} \sum_{h=1}^{H} \langle \widehat{\phi}_{\pi',h}^{\ell}, \widehat{\theta}_{h}^{\ell} \rangle \leq \sup_{\pi' \in \Pi_{\ell}} \sum_{h=1}^{H} \langle \phi_{\pi',h}, \boldsymbol{\theta}_{h} \rangle + \epsilon_{\ell} = V_{0}^{\widetilde{\pi}^{\star}} + \epsilon_{\ell}.$$

Rearranging this gives

$$\sum_{h=1}^{H} \langle \widehat{\phi}_{\widetilde{\pi}^{\star},h}^{\ell}, \widehat{\theta}_{h}^{\ell} \rangle \geq \sup_{\pi' \in \Pi_{\ell}} \sum_{h=1}^{H} \langle \widehat{\phi}_{\pi',h}^{\ell}, \widehat{\theta}_{h}^{\ell} \rangle - 2\epsilon_{\ell}$$

so  $\widetilde{\pi}^{\star} \in \Pi_{\ell+1}$ .

**Theorem 7.** With probability at least  $1-2\delta$ , Algorithm 1 will terminate after collecting at most

$$CH^{4} \sum_{h=1}^{H} \sum_{\ell=\ell_{0}+1}^{\iota_{0}} \frac{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_{h}} \max_{\pi \in \Pi(4\epsilon_{\ell})} \|\phi_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^{2}}{\epsilon_{\ell}^{2}} \cdot \log \frac{H|\Pi(4\epsilon_{\ell})| \log \frac{1}{\epsilon}}{\delta} + \operatorname{poly}\left(d, H, \frac{1}{\lambda_{\min}^{\star}}, \log \frac{1}{\delta}, \log |\Pi|, \log \frac{1}{\epsilon}\right) \\ + CH^{4} \sum_{h=1}^{H} \frac{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_{h}} \max_{\pi \in \Pi} \|\phi_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^{2}}{\epsilon_{\ell_{0}}^{2}} \cdot \log \frac{H|\Pi| \log(1/\epsilon)}{\delta}$$

episodes for  $\iota_0 := \min\{\lceil \log \frac{4}{\epsilon} \rceil, \log \frac{4}{\Delta_{\min}(\Pi)} \}$ , and will output a policy  $\widehat{\pi}$  such that

$$V_0^{\widehat{\pi}} \ge \max_{\pi \in \Pi} V_0^{\pi} - \epsilon,$$

where here  $\Pi(4\epsilon_{\ell}) = \{\pi \in \Pi : V_0^{\pi} \ge \max_{\pi \in \Pi} V_0^{\pi} - 4\epsilon_{\ell}\}.$ 

*Proof.* By Lemma B.8 the event  $\mathcal{E}_{est} \cap \mathcal{E}_{exp}$  occurs with probability at least  $1 - 2\delta$ . Henceforth we assume we are on this event.

Correctness follows by Lemma B.9, since upon termination,  $\Pi_{\ell}$  will only contain policies  $\pi$  satisfying  $V_0^{\pi} \geq \max_{\pi \in \Pi} V_0^{\pi} - \epsilon$  (and will contain at least 1 policy since  $\widetilde{\pi}^{\star} \in \Pi_{\ell}$  for all  $\ell$ ). Furthermore, by Lemma B.9, if  $4\epsilon_{\ell} < \Delta_{\min}(\Pi)$ , we must have that  $\Pi_{\ell} = \{\widetilde{\pi}^{\star}\}$ , and will therefore terminate on Line 11 since  $|\Pi_{\ell}| = 1$ . Thus, we can bound the number of number of epochs by

$$\iota_0 := \min\{\lceil \log \frac{4}{\epsilon} \rceil, \log \frac{4}{\Delta_{\min}(\Pi)} \}.$$

By Lemma B.6, the total number of episodes collected is bounded by

$$\begin{split} &\sum_{h=1}^{H} \sum_{\ell=1}^{\iota_0} C \cdot \frac{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\boldsymbol{\phi} \in \Phi_{\ell,h}} \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2}{\epsilon_\ell^2 / \beta_\ell} + \operatorname{poly}\left(d, H, \log \frac{1}{\delta}, \frac{1}{\lambda_{\min}^*}, \log |\Pi|, \log \frac{1}{\epsilon}\right) \\ &\leq \sum_{h=1}^{H} \sum_{\ell=1}^{\iota_0} C \cdot \frac{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\boldsymbol{\phi} \in \Phi_{\ell,h}} \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2}{\epsilon_\ell^2} \cdot H^4 \log \frac{H|\Pi_\ell| \log(1/\epsilon)}{\delta} + \operatorname{poly}\left(d, H, \log \frac{1}{\delta}, \frac{1}{\lambda_{\min}^*}, \log |\Pi|, \log \frac{1}{\epsilon}\right). \end{split}$$

On  $\mathcal{E}_{\text{est}} \cap \mathcal{E}_{\text{exp}}$ , by Lemma B.8, for each  $\pi \in \Pi_{\ell}$ , we have  $\|\widehat{\phi}_{\pi,h}^{\ell} - \phi_{\pi,h}\|_{2} \leq d\epsilon_{\ell}/2H$ . As  $\Phi_{\ell,h} = \{\widehat{\phi}_{\pi,h}^{\ell} : \pi \in \Pi_{\ell}\}$ , it follows that we can upper bound

$$\begin{split} \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\boldsymbol{\phi} \in \Phi_{\ell,h}} \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2 &= \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi_{\ell}} \|\widehat{\boldsymbol{\phi}}_{\pi,h}^{\ell}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2 \\ &\leq \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi_{\ell}} (2\|\boldsymbol{\phi}_{\pi,h}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2 + 2\|\widehat{\boldsymbol{\phi}}_{\pi,h}^{\ell} - \boldsymbol{\phi}_{\pi,h}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2) \\ &\leq \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi_{\ell}} (2\|\boldsymbol{\phi}_{\pi,h}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2 + \frac{d^2 \epsilon_{\ell}^2}{2H^2 \lambda_{\min}(\mathbf{A}(\mathbf{\Lambda}))}) \end{split}$$

$$\leq \inf_{\mathbf{\Lambda} \in \Omega_h} \max_{\pi \in \Pi_{\ell}} 4 \|\phi_{\pi,h}\|_{\mathbf{\Lambda}(\mathbf{\Lambda})^{-1}}^2 + \inf_{\pi} \frac{d^2 \epsilon_{\ell}^2}{H^2 \lambda_{\min}(\mathbf{\Lambda}(\mathbf{\Lambda}))}$$
$$\leq \inf_{\mathbf{\Lambda} \in \Omega_h} \max_{\pi \in \Pi_{\ell}} 4 \|\phi_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2 + \frac{d^2 \epsilon_{\ell}^2}{H^2 \lambda_{\min}^*}$$

SO

$$\frac{\inf_{\mathbf{\Lambda}\in\mathbf{\Omega}_h}\max_{\boldsymbol{\phi}\in\Phi_{\ell,h}}\|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2}{\epsilon_\ell^2}\leq \frac{\inf_{\mathbf{\Lambda}\in\mathbf{\Omega}_h}\max_{\pi\in\Pi_\ell}4\|\boldsymbol{\phi}_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{\epsilon_\ell^2}+\frac{d^2}{H^2\lambda_{\min}^\star}.$$

Note also that, by Lemma B.9, for  $\ell > \ell_0$ , every policy  $\pi \in \Pi_{\ell}$  will be  $4\epsilon_{\ell}$  optimal, so we therefore have

$$\Pi_{\ell} \subseteq \{ \pi \in \Pi : V_0^{\pi} \ge V_0^{\widetilde{\pi}^{\star}} - 4\epsilon_{\ell} \} =: \Pi(4\epsilon_{\ell}).$$

Putting this together, we can upper bound the complexity by

$$\begin{split} \sum_{h=1}^{H} \sum_{\ell=\ell_0+1}^{\iota_0} C \cdot \frac{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi(4\epsilon_\ell)} \|\boldsymbol{\phi}_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{\epsilon_\ell^2} \cdot H^4 \log \frac{H|\Pi(4\epsilon_\ell)| \log(1/\epsilon)}{\delta} + \operatorname{poly}\left(d, H, \log \frac{1}{\delta}, \frac{1}{\lambda_{\min}^{\star}}, \log |\Pi|, \log \frac{1}{\epsilon}\right) \\ + C \cdot \frac{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi} \|\boldsymbol{\phi}_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{\epsilon_{\ell_0}^2} \cdot H^4 \log \frac{H|\Pi| \log(1/\epsilon)}{\delta}. \end{split}$$

Corollary 5 (Full Statement of Theorem 1). With probability at least  $1 - \delta$ , the complexity of Algorithm 1 can be bounded as

$$CH^{4} \log \frac{1}{\epsilon} \cdot \sum_{h=1}^{H} \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_{h}} \max_{\pi \in \Pi} \frac{\|\phi_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^{2}}{(V_{0}^{\star}(\Pi) - V_{0}^{\pi})^{2} \vee \epsilon^{2} \vee \Delta_{\min}(\Pi)^{2}} \cdot \log \frac{H|\Pi| \log \frac{1}{\epsilon}}{\delta} + \operatorname{poly} \left(d, H, \frac{1}{\lambda_{\min}^{\star}}, \log \frac{1}{\delta}, \log |\Pi|, \log \frac{1}{\epsilon}\right)$$

episodes, and Algorithm 1 will output a policy  $\hat{\pi}$  such that

$$V_0^{\widehat{\pi}} \ge \max_{\pi \in \Pi} V_0^{\pi} - \epsilon.$$

*Proof.* By the definition of  $\Pi(4\epsilon_{\ell})$ , for each  $\pi \in \Pi(4\epsilon_{\ell})$  we have

$$\epsilon_{\ell}^2 = \frac{1}{16} \left( (V_0^{\star}(\Pi) - V_0^{\pi})^2 \vee (4\epsilon_{\ell})^2 \right).$$

We can therefore upper bound

$$\begin{split} &\sum_{\ell=\ell_0+1}^{\iota_0} \frac{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi(4\epsilon_\ell)} \|\boldsymbol{\phi}_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{\epsilon_\ell^2} \cdot \log \frac{H|\Pi(4\epsilon_\ell)| \log \frac{1}{\epsilon}}{\delta} \\ &\leq C \sum_{\ell=\ell_0+1}^{\iota_0} \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi(4\epsilon_\ell)} \frac{\|\boldsymbol{\phi}_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{(V_0^{\star}(\Pi) - V_0^{\pi})^2 \vee \epsilon_\ell^2} \cdot \log \frac{H|\Pi(4\epsilon_\ell)| \log \frac{1}{\epsilon}}{\delta} \\ &\leq C \log \frac{1}{\epsilon} \cdot \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi} \frac{\|\boldsymbol{\phi}_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{(V_0^{\star}(\Pi) - V_0^{\pi})^2 \vee \epsilon^2 \vee \Delta_{\min}(\Pi)^2} \cdot \log \frac{H|\Pi| \log \frac{1}{\epsilon}}{\delta}. \end{split}$$

Furthermore, since  $\ell_0 = \lceil \log_2 d^{3/2}/H \rceil$ , using Lemma B.10 we can also bound

$$C \cdot \frac{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi} \|\boldsymbol{\phi}_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{\epsilon_{\ell_0}^2} \cdot H^4 \log \frac{H|\Pi| \log(1/\epsilon)}{\delta} \leq \operatorname{poly}\left(d, H, \log 1/\delta, \log |\Pi|, \log 1/\epsilon\right).$$

The result then follows by Theorem 7.

Proof of Corollary 1. By Lemma A.14, we can choose  $\Pi_{\epsilon}$  to be the restricted-action linear softmax policy set constructed in Lemma A.14. Lemma A.14 shows that  $\Pi_{\epsilon}$  will contain an  $\epsilon$ -optimal policy for any MDP and reward function, and that

$$|\Pi_{\epsilon}| \le \left(1 + \frac{32H^4d^{5/2}\log(1 + 16Hd/\epsilon)}{\epsilon^2}\right)^{dH^2}.$$

Combining this with the guarantee of Corollary 5 shows that  $V_0^{\widehat{\pi}} \geq V_0^{\star} - 2\epsilon$  and that  $V_0^{\star}(\Pi) - V_0^{\pi}$  is within a factor of  $\epsilon$  of  $V_0^{\star} - V_0^{\pi}$ . To bound the complexity of this procedure, we apply the bound given in Corollary 5 with the bound on the cardinality of  $\Pi_{\epsilon}$  given above.

# **B.3** Interpreting the Complexity

**Lemma B.10.** For any set of policies  $\Pi$ , we can bound

$$\inf_{\mathbf{\Lambda}\in\mathbf{\Omega}_h}\sup_{\pi\in\Pi}\|\phi_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2\leq d.$$

*Proof.* By Jensen's inequality, for any  $\boldsymbol{v} \in \mathbb{R}^d$ , we have

$$oldsymbol{v}^ op oldsymbol{\Lambda}_{\pi,h} oldsymbol{v} = \mathbb{E}_{\pi}[(oldsymbol{v}^ op oldsymbol{\phi}_h)^2] \geq (\mathbb{E}_{\pi}[oldsymbol{v}^ op oldsymbol{\phi}_h])^2 = (oldsymbol{v}^ op oldsymbol{\phi}_{\pi,h})^2.$$

It follows that, for any  $\pi$ ,

$$oldsymbol{\Lambda}_{\pi,h}\succeq oldsymbol{\phi}_{\pi,h}oldsymbol{\phi}_{\pi,h}^{ op}$$

Take  $\Lambda \in \Omega_h$ . Then,

$$oldsymbol{\Lambda} = \mathbb{E}_{\pi \sim \omega}[oldsymbol{\Lambda}_{\pi,h}] \succeq \mathbb{E}_{\pi \sim \omega}[oldsymbol{\phi}_{\pi,h}oldsymbol{\phi}_{\pi,h}^{ op}].$$

It follows that we can upper bound

$$\inf_{\mathbf{\Lambda} \in \Omega_h} \sup_{\pi \in \Pi} \|\phi_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2 \leq \inf_{\lambda \in \Delta_\Pi} \sup_{\pi \in \Pi} \|\phi_{\pi,h}\|_{A(\lambda)^{-1}}^2$$

where  $A(\lambda) = \sum_{\pi} \lambda_{\pi} \phi_{\pi,h} \phi_{\pi,h}^{\top}$ . By Kiefer-Wolfowitz (Lattimore & Szepesvári, 2020), this is upper bounded by d.

*Proof of Corollary* 2. This follows directly from Lemma B.10 and Corollary 1, by upper bounding:

$$\inf_{\mathbf{\Lambda} \in \Omega_h} \max_{\pi \in \Pi_{\epsilon}} \frac{\|\phi_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{\max\{V_{\mathbf{\Lambda}}^{\star} - V_{\mathbf{\Lambda}}^{\pi}, \epsilon\}^2} \leq \inf_{\mathbf{\Lambda} \in \Omega_h} \max_{\pi \in \Pi_{\epsilon}} \frac{\|\phi_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{\epsilon^2} \leq \frac{d}{\epsilon^2}.$$

### **B.3.1** Linear Contextual Bandits

Since we always assume the MDP starts in some state  $s_1$ , to encode a linear contextual bandit, the direct mapping of our linear MDP in Definition 3.1 would require considering an H = 2 MDP, where we encode the "context" in the transition to state s at step h = 2. While we could run our algorithm directly on this, in the standard contextual bandit setting, the learner has no control over the context, and so their action before receiving that context has no effect. Thus, there is no need for the learner to explore at stage h = 1. To account for this, we can simply run our algorithm but ignore the exploration at stage h = 1, which will reduce the h = 1 term in the sample complexity.

### B.3.2 Tabular MDPs

**Lemma B.11.** In the tabular MDP setting, assuming that  $\Pi$  contains an optimal policy,

$$\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi} \frac{\|\phi_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{(V_0^{\star} - V_0^{\pi})^2 \vee \epsilon^2 \vee \Delta_{\min}(\Pi)^2} \\
\leq \inf_{\pi_{\exp}} \max_{\pi \in \Pi} \max_{s,a} \frac{1}{w_h^{\pi_{\exp}}(s,a)} \min \left\{ \frac{1}{w_h^{\pi}(s,a)\Delta_h(s,a)^2}, \frac{w_h^{\pi}(s,a)}{\epsilon^2 \vee \Delta_{\min}(\Pi)^2} \right\} \\
\leq \inf_{\pi_{\exp}} \max_{s,a} \frac{1}{w_h^{\pi_{\exp}}(s,a)} \cdot \frac{1}{\epsilon \max\{\Delta_h(s,a),\epsilon,\Delta_{\min}(\Pi)\}}.$$

Proof. We have that  $[\phi_{\pi,h}]_{s,a} = w_h^{\pi}(s,a)$ . Furthermore,  $\phi(s,a) = e_{s,a}$ , so for any  $\Lambda \in \Omega_h$ ,  $\Lambda$  is diagonal with  $[\Lambda]_{sa,sa} = \mathbb{E}_{\pi \sim \omega}[w_h^{\pi}(s,a)]$ . Furthermore, by the Performance-Difference Lemma,  $V_0^{\star} - V_0^{\pi} = \sum_{s,a,h} w_h^{\pi}(s,a) \Delta_h(s,a)$ . Thus,

$$\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi} \frac{\|\phi_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{(V_0^{\star} - V_0^{\pi})^2 \vee \epsilon^2 \vee \Delta_{\min}(\Pi)^2} \leq \inf_{\pi_{\exp}} \max_{\pi \in \Pi} \frac{\sum_{s,a} \frac{w_h^{\pi}(s,a)^2}{w_h^{\exp}(s,a)}}{(\sum_{s',a',h'} w_{h'}^{\pi}(s',a')\Delta_{h'}(s',a'))^2 \vee \epsilon^2 \vee \Delta_{\min}(\Pi)^2}.$$
(B.1)

We have

$$\sum_{s,a} \frac{w_h^{\pi}(s,a)^2}{w_h^{\pi_{\exp}}(s,a)} \leq \left(\sum_{s,a} w_h^{\pi}(s,a)\right) \cdot \max_{s,a} \frac{w_h^{\pi}(s,a)}{w_h^{\pi_{\exp}}(s,a)} = \max_{s,a} \frac{w_h^{\pi}(s,a)}{w_h^{\pi_{\exp}}(s,a)}.$$

Thus,

$$(\mathbf{B}.1) \leq \inf_{\pi_{\exp}} \max_{\pi \in \Pi} \max_{s,a} \frac{w_h^{\pi}(s,a)/w_h^{\pi_{\exp}}(s,a)}{(\sum_{s',a',h'} w_{h'}^{\pi}(s',a')\Delta_{h'}(s',a'))^2 \vee \epsilon^2 \vee \Delta_{\min}(\Pi)^2}$$

$$\leq \inf_{\pi_{\exp}} \max_{\pi \in \Pi} \max_{s,a} \frac{w_h^{\pi}(s,a)/w_h^{\pi_{\exp}}(s,a)}{(w_h^{\pi}(s,a)\Delta_h(s,a))^2 \vee \epsilon^2 \vee \Delta_{\min}(\Pi)^2}$$

$$= \inf_{\pi_{\exp}} \max_{\pi \in \Pi} \max_{s,a} \frac{1}{w_h^{\pi_{\exp}}(s,a)} \min \left\{ \frac{1}{w_h^{\pi}(s,a)\Delta_h(s,a)^2}, \frac{w_h^{\pi}(s,a)}{\epsilon^2 \vee \Delta_{\min}(\Pi)^2} \right\}. \tag{B.2}$$

We can further upper bound

$$\min\left\{\frac{1}{w_h^{\pi}(s,a)\Delta_h(s,a)^2}, \frac{w_h^{\pi}(s,a)}{\epsilon^2 \vee \Delta_{\min}(\Pi)^2}\right\} \leq \frac{1}{\Delta_h(s,a)(\epsilon \vee \Delta_{\min}(\Pi))}$$

so

$$\begin{aligned} &(\mathbf{B.2}) \leq \inf_{\pi_{\text{exp}}} \max_{\pi \in \Pi} \max_{s, a} \frac{1}{w_h^{\pi_{\text{exp}}}(s, a)} \min \left\{ \frac{1}{\Delta_h(s, a) \epsilon}, \frac{w_h^{\pi}(s, a)}{\epsilon^2 \vee \Delta_{\min}(\Pi)^2} \right\} \\ &\leq \inf_{\pi_{\text{exp}}} \max_{s, a} \frac{1}{w_h^{\pi_{\text{exp}}}(s, a)} \frac{1}{\epsilon \max \{ \Delta_h(s, a), \epsilon, \Delta_{\min}(\Pi) \}}. \end{aligned}$$

**Lemma B.12.** If PEDEL is run with a set  $\Pi$  that contains an optimal policy, the complexity of PEDEL is upper bounded as

$$\widetilde{\mathcal{O}}\left(H^4\sum_{h=1}^{H}\sup_{\epsilon'\geq \max\{\epsilon,\Delta_{\min}(\Pi)/4\}}\inf_{\pi_{\exp}}\max_{\pi\in\Pi(\epsilon')}\max_{s,a}\frac{1}{w_h^{\pi_{\exp}}(s,a)}\min\left\{\frac{1}{w_h^{\pi}(s,a)\Delta_h(s,a)^2},\frac{w_h^{\pi}(s,a)}{(\epsilon')^2}\right\}\cdot\log\frac{|\Pi(\epsilon')|}{\delta}+C_0\right)$$

for 
$$\Pi(\epsilon') = \{ \pi \in \Pi : V_0^{\pi} \ge V_0^{\star}(\Pi) - \epsilon \}.$$

*Proof.* Using an argument identical to that in Lemma B.11, we can upper bound

$$\begin{split} \frac{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi(4\epsilon_\ell)} \|\boldsymbol{\phi}_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{\epsilon_\ell^2} &\leq \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi(4\epsilon_\ell)} \frac{c \|\boldsymbol{\phi}_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{(V_0^* - V_0^\pi)^2 \vee \epsilon_\ell^2} \\ &\leq \inf_{\pi_{\exp}} \max_{\pi \in \Pi(4\epsilon_\ell)} \max_{s,a} \frac{1}{w_h^{\pi_{\exp}}(s,a)} \min \left\{ \frac{1}{w_h^\pi(s,a)\Delta_h(s,a)^2}, \frac{w_h^\pi(s,a)}{\epsilon_\ell^2} \right\}. \end{split}$$

The result then follows from Theorem 7, noting that we will never run for  $\epsilon_{\ell} < \Delta_{\min}(\Pi)/4$ .

Proof of Corollary 3. Note that in the tabular MDP setting, we can choose  $\Pi$  to be the set of all deterministic policies, since this set is guaranteed to contain an optimal policy. We can then bound  $|\Pi| \leq A^{SH}$ . The result then follows directly from Lemma B.11 and Theorem 1.

*Proof of Proposition* 4. We begin with an example where PEDEL has complexity smaller than the Gap-Visitation Complexity, and then turn to an example where the reverse is true.

**Pedel Improves on Gap-Visitation Complexity.** Consider the tabular MDP with |S| = |A| = N, and where

$$\begin{split} P_h(s_1|s_1,a_1) &= 1, \quad \nu_h(s_1,a_1) = 1, \forall h \in [H] \\ P_h(s_1|s_1,a_j) &= 0, \quad \nu_h(s_1,a_j) = 0, \forall h \in [H], j \neq 1 \\ P_h(s_1|s_i,a_j) &= 0, \forall h \in [H], j \in [N], i \neq 1 \\ P_h(s_i|s_j,a_i) &= 1, \forall h \in [H], j \in [N], i \neq 1 \\ r_h(s_i,a_1) &= \epsilon, \forall h \in [H], i \neq 1, \quad \nu_h(s_i,a_j) = 0, \forall h \in [H], j \neq 1, i \neq 1. \end{split}$$

In this MDP, the optimal policy simply plays action  $a_1$  H times and is always in state  $s_1$ . The total reward it collects is H. Any deterministic policy that does not play  $a_1$  H consecutive times has optimality gap of at least  $1 - \epsilon$ . Furthermore, every other state can be reached with probability 1. In this case, then, assuming that we take  $\Pi$  to be the set of all deterministic policies, we have  $\Delta_{\min}(\Pi) = 1 - \epsilon$  (note that since there always exists a deterministic policy that is optimal, it suffices to take  $\Pi$  to be the set of all deterministic policies).

By Corollary 3, we can therefore upper bound the complexity of the leading-order term by  $\widetilde{\mathcal{O}}(H^5S^2A)$ , so PEDEL will identify the optimal policy (since  $\Pi$  contains an optimal policy). Thus, the total complexity of PEDEL is  $\mathcal{O}(\text{poly}(S, A, H, \log 1/\delta))$ .

On this example, in every state  $s_i$ ,  $i \neq 1$ , action  $a_1$  still collects a reward of  $\epsilon$ . Thus, we have that  $\Delta_h(s_i, a_j) = \epsilon$  for  $j \neq 1$ . The Gap-Visitation complexity is given by

$$\sum_{h=1}^{H} \inf_{\pi} \max_{s,a} \min \left\{ \frac{1}{w_h^{\pi}(s,a)\Delta_h(s,a)^2}, \frac{W_h(s)^2}{\epsilon^2} \right\}.$$

Since  $W_h(s) = 1$  for each s, we conclude that

$$\sum_{h=1}^{H} \inf_{\pi} \max_{s,a} \min \left\{ \frac{1}{w_h^{\pi}(s,a)\Delta_h(s,a)^2}, \frac{W_h(s)^2}{\epsilon^2} \right\} \ge \sum_{h=1}^{H} \frac{1}{\epsilon^2}.$$

Thus, for small  $\epsilon$ , the Gap-Visitation complexity can be arbitrarily worse than the complexity of PEDEL.

The Gap-Visitation Complexity Improves on Pedel. To show that the Gap-Visitation Complexity improves on the complexity of PEDEL, we consider the example in Instance Class 5.1 of Wagenmaker et al. (2021b). As shown by Proposition 6 of Wagenmaker et al. (2021b), on this example, for any  $\epsilon$ , the Gap-Visitation Complexity is  $\widetilde{\mathcal{O}}(\text{poly}(S))$ .

To bound the complexity of PEDEL on this example, we consider the complexity given in Theorem 7 with  $\Pi$  the set of all deterministic policies, which is slightly tighter than the complexity of Corollary 3. Take  $\epsilon \geq 2^{-S}$ . Then, on this example, it follows that  $\Delta_{\min}(\Pi) \leq \mathcal{O}(\epsilon)$ , since we can find a policy  $\pi$  which is optimal on every state  $s_i$  at step h=2 for  $i=\mathcal{O}(\log 1/\epsilon)$ , which will give it a policy gap of  $\mathcal{O}(\epsilon)$ . Furthermore, any near-optimal policy will have  $[\phi_{\pi,2}]_{s_1,a_1} = w_2^{\pi}(s_1,a_1) = \mathcal{O}(1)$ , so we always have  $\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_2} \max_{\pi \in \Pi(4\epsilon_{\ell})} \|\phi_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2 \geq \Omega(1)$ . It follows that the complexity of PEDEL is lower bounded by  $\Omega(1/\epsilon^2)$ .

# B.3.3 Deterministic, Tabular MDPs

**Lemma B.13.** In the deterministic MDP setting,

$$\inf_{\mathbf{\Lambda}\in\mathbf{\Omega}_h}\max_{\pi\in\Pi}\frac{\|\boldsymbol{\phi}_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{(V_0^{\star}-V_0^{\pi})^2\vee\epsilon^2}\leq\sum_{s,a}\frac{1}{\bar{\Delta}_h(s,a)^2\vee\epsilon^2}.$$

Proof. Note that  $[\phi_{\pi,h}]_{s_h^{\pi},a_h^{\pi}} = 1$ , and otherwise, for  $(s,a) \neq (s_h^{\pi},a_h^{\pi})$ ,  $[\phi_{\pi,h}]_{s,a} = 0$ . Furthermore,  $\Lambda_{\pi_{\text{exp}},h}$  will always be diagonal, with diagonal elements  $w_h^{\pi}(s,a)$ . We then have  $\|\phi_{\pi,h}\|_{\Lambda_{\text{exp}}^{-1},h}^2 = \frac{1}{w_h^{\pi_{\text{exp}}}(s_h^{\pi},a_h^{\pi})}$ , so

$$\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}_h} \max_{\pi \in \Pi} \frac{\|\phi_{\pi,h}\|_{\mathbf{\Lambda}^{-1}}^2}{(V_0^{\star} - V_0^{\pi})^2 \vee \epsilon^2} \leq \inf_{\pi \exp} \max_{\pi \in \Pi} \frac{\|\phi_{\pi,h}\|_{\mathbf{\Lambda}_{\pi \exp,h}}^2}{(V_0^{\star} - V_0^{\pi})^2 \vee \epsilon^2} \\
= \inf_{\pi \exp} \max_{\pi \in \Pi} \frac{w_h^{\pi \exp}(s_h^{\pi}, a_h^{\pi})^{-1}}{(V_0^{\star} - V_0^{\pi})^2 \vee \epsilon^2}$$

$$\begin{array}{l} \overset{(a)}{=} \inf_{\pi_{\exp}} \max_{s,a} \max_{\pi \in \Pi_{sah}} \frac{w_h^{\pi_{\exp}}(s_h^{\pi}, a_h^{\pi})^{-1}}{(V_0^{\star} - V_0^{\pi})^2 \vee \epsilon^2} \\ \overset{(b)}{=} \inf_{\pi_{\exp}} \max_{s,a} \max_{\pi \in \Pi_{sah}} \frac{w_h^{\pi_{\exp}}(s, a)^{-1}}{(V_0^{\star} - V_0^{\pi})^2 \vee \epsilon^2} \\ = \inf_{\pi_{\exp}} \max_{s,a} \frac{w_h^{\pi_{\exp}}(s, a)^{-1}}{(V_0^{\star} - \max_{\pi \in \Pi_{sah}} V_0^{\pi})^2 \vee \epsilon^2} \\ \overset{(c)}{=} \inf_{\pi_{\exp}} \max_{s,a} \frac{w_h^{\pi_{\exp}}(s, a)^{-1}}{\bar{\Delta}_h(s, a)^2 \vee \epsilon^2} \end{array}$$

where (a) follows since  $\Pi = \bigcup_{s,a} \Pi_{sah}$ , (b) follows since by definition, for any  $\pi \in \Pi_{sah}$ ,  $(s_h^{\pi}, a_h^{\pi}) = (s, a)$ , and (c) follows by the definition of  $\bar{\Delta}_h(s, a)$ .

Let  $\pi^{sa}$  denote any policy such that  $(s_h^{\pi}, a_h^{\pi}) = (s, a)$ . Set

$$\lambda_{\pi^{sa}} = \frac{\max\{\bar{\Delta}_h(s,a),\epsilon\}^{-2}}{\sum_{s',a'}\max\{\bar{\Delta}_h(s',a'),\epsilon\}^{-2}}.$$

Note that this is a valid distribution. Let  $\pi_{\exp} = \sum_{s,a} \lambda_{\pi^{sa}} \pi^{sa}$ , then  $w_h^{\pi_{\exp}}(s,a) = \lambda_{\pi^{sa}}$ , so

$$\inf_{\pi_{\exp}} \max_{s,a} \frac{w_h^{\pi_{\exp}}(s,a)^{-1}}{\bar{\Delta}_h(s,a)^2 \vee \epsilon^2} \le \max_{s,a} \frac{\lambda_{\pi^{sa}}^{-1}}{\bar{\Delta}_h(s,a)^2 \vee \epsilon^2}$$
$$\le \sum_{s,a} \frac{1}{\bar{\Delta}_h(s,a)^2 \vee \epsilon^2}$$

which proves the result.

Proof of Corollary 4. As in tabular MDPs, we can set  $\Pi$  to correspond to the set of all deterministic policies. However, since our MDP is also deterministic, at any given h, we only need to specify  $\pi_h(s)$  for a single s—the state we will end up in at step h with probability 1. Thus, we can take  $\Pi$  to be a set of cardinality  $|\Pi| = A^H$ . The result then follows directly from Lemma B.13 and Theorem 1.  $\square$ 

Comparison to Lower Bound of Tirinzoni et al. (2022). The precise definition for  $\bar{\Delta}_{\min}^h$  is  $\bar{\Delta}_{\min}^h := \min_{s,a:\bar{\Delta}_h(s,a)>0} \bar{\Delta}_h(s,a)$  in the setting when every deterministic  $\epsilon$ -optimal policy will reach the same (s,a) at step h, and  $\bar{\Delta}_{\min}^h := 0$  otherwise.

The exact lower bound given in Tirinzoni et al. (2022) scales as  $\varphi^*(\underline{c})$  which does not have an explicit form. However, they show that

$$\max_{h \in [H]} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \frac{\log(1/4\delta)}{4 \max\{\bar{\Delta}_h(s, a), \bar{\Delta}_{\min}^h, \epsilon\}^2} \leq \varphi^{\star}(\underline{c}) \leq \sum_{h \in [H]} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \frac{\log(1/4\delta)}{4 \max\{\bar{\Delta}_h(s, a), \bar{\Delta}_{\min}^h, \epsilon\}^2}.$$

Up to H factors, then, this matches the complexity of our upper bound in every term but the  $\bar{\Delta}_{\min}^h$  term.  $\bar{\Delta}_{\min}^h \geq \bar{\Delta}_{\min}$ , so this lower bound is potentially smaller than our upper bound in this dependence. We remark, however, that the algorithm presented in Tirinzoni et al. (2022) obtains the same scaling as we do, depending on  $\bar{\Delta}_{\min}$  instead of  $\bar{\Delta}_{\min}^h$ . Furthermore, in general we can think of these quantities as scaling in a similar manner, since they each quantify the minimum policy gap.

# C Experiment Design via Online Frank-Wolfe

Through the remainder of Appendix C as well as Appendix D, we will be interested in the problem of data collection in linear MDPs. In general, we will seek to collect data for a particular  $h \in [H]$ . We will therefore consider the following truncation to our MDP.

**Definition C.1** (Truncated Horizon MDPs). Given some MDP  $\mathcal{M}$  with horizon H, we define the h-truncated-horizon MDP  $\mathcal{M}_{tr,h}$  to be the MDP that is identical to  $\mathcal{M}$  for  $h' \leq h$ , but that terminates after reaching state  $s_h$  and playing action  $a_h$ .

We can simulated a truncated-horizon MDP by playing in our standard MDP  $\mathcal{M}$ , and after taking an action at step h,  $a_h$ , taking random actions for h' > h and ignoring all future observations.

The utility of considering truncated-horizon MDPs is that we can therefore guarantee the data we collect,  $\{\{(s_{h',\tau},a_{h',\tau})\}_{h'=1}^{h}\}_{\tau=1}^{K}$  is uncorrelated with the true next state and reward at step h obtained in  $\mathcal{M}$ ,  $\{(s_{h+1,\tau},r_{h,\tau})\}_{\tau=1}^{K}$ . While we do not allow our algorithm to use  $\{(s_{h+1,\tau},r_{h,\tau})\}_{\tau=1}^{K}$  in its operation, it is allowed to store this data and return it.

For the remainder of Appendix C and Appendix D, then, we assume there is some fixed h we are interested in, and that we are running our algorithms in the h-truncated-horizon MDP defined with respect to our original MDP. We will also drop the subscript of h from observations, so  $\Lambda_{\pi} = \Lambda_{\pi,h}$ ,  $\phi_{\tau} = \phi_{\tau,h}$ , and  $\Omega = \Omega_{h}$ .

Our main experiment design algorithm, OPTCOV, relies on a regret-minimization algorithm satisfying the following guarantee.

**Definition C.2** (Regret Minimization Algorithm). We say REGMIN is a regret minimization algorithm if it has regret scaling as, with probability at least  $1 - \delta$ ,

$$\mathcal{R}_K := \sum_{k=1}^K (V_0^{\star} - V_0^{\pi_k}) \le \sqrt{\mathcal{C}_1 K \log^{p_1}(HK/\delta)} + \mathcal{C}_2 \log^{p_2}(HK/\delta)$$

for any deterministic reward function  $r_h(s, a) \in [0, 1]$ .

Throughout this section, we will let  $\Lambda$  refer to covariates normalized by time, and  $\Sigma$  unnormalized covariates. So, for example, we might have  $\Sigma = \sum_{\tau=1}^{T} \phi_{\tau} \phi_{\tau}^{\mathsf{T}}$  and  $\Lambda = \frac{1}{T} \sum_{\tau=1}^{T} \phi_{\tau} \phi_{\tau}^{\mathsf{T}}$ .

The rest of this section is organized as follows. First, in Appendix C.1 we show that a variant of the Frank-Wolfe algorithm that relies on only approximate updates enjoys a convergence rate similar to the standard Frank-Wolfe rate. Next, in Appendix C.2 we show that for a smooth experiment design objective, we can approximately optimize the objective in a linear MDP by approximating the Frank-Wolfe updates via a regret minimization algorithm. Finally, in Appendix C.3 we present our main experiment-design algorithm, OPTCOV, which relies on our online Frank-Wolfe procedure to collect covariates that minimize an online experimental design objective up to an arbitrarily tolerance.

### C.1 Approximate Frank-Wolfe

We will consider the following approximate variant of the Frank-Wolfe algorithm:

**Lemma C.1.** Consider running Algorithm  $\frac{1}{4}$  with some convex function f that is  $\beta$ -smooth with respect to some norm  $\|\cdot\|$ , and let  $R := \sup_{x,y \in \mathcal{X}} \|x - y\|$ . Then for  $T \geq 2$ , we have

$$f(x_{T+1}) - \min_{x \in \mathcal{X}} f(x) \le \frac{\beta R^2 (\log T + 1)}{2(T+1)} + \frac{1}{T+1} \sum_{t=1}^{T} \epsilon_t.$$

# Algorithm 4 Approximate Frank-Wolfe

- 1: **input**: function to optimize f, number of iterations to run T, starting iterate  $x_1$
- 2: **for** t = 1, 2, ..., T **do**
- 3: Set  $\gamma_t \leftarrow \frac{1}{t+1}$
- 4: Choose  $y_t$  to be any point such that

$$abla f(oldsymbol{x}_t)^{ op} oldsymbol{y}_t \leq \min_{oldsymbol{y} \in \mathcal{X}} 
abla f(oldsymbol{x}_t)^{ op} oldsymbol{y} + \epsilon_t$$

- 5:  $x_{t+1} \leftarrow (1 \gamma_t)x_t + \gamma_t y_t$
- 6: return  $x_{T+1}$

*Proof.* Let  $\mathbf{x}^* = \arg\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ . Using that f is  $\beta$ -smooth, the definition of  $\mathbf{y}_s$ , and the convexity of f, we have that for any s,

$$f(\boldsymbol{x}_{s+1}) - f(\boldsymbol{x}_s) \leq \nabla f(\boldsymbol{x}_s)^{\top} (\boldsymbol{x}_{s+1} - \boldsymbol{x}_s) + \frac{\beta}{2} \|\boldsymbol{x}_{s+1} - \boldsymbol{x}_s\|^2$$

$$\leq \gamma_s \nabla f(\boldsymbol{x}_s)^{\top} (\boldsymbol{y}_s - \boldsymbol{x}_s) + \frac{\beta}{2} \gamma_s^2 R^2$$

$$\leq \gamma_s \nabla f(\boldsymbol{x}_s)^{\top} (\boldsymbol{x}^* - \boldsymbol{x}_s) + \gamma_s \epsilon_s + \frac{\beta}{2} \gamma_s^2 R^2$$

$$\leq \gamma_s (f(\boldsymbol{x}^*) - f(\boldsymbol{x}_s)) + \gamma_s \epsilon_s + \frac{\beta}{2} \gamma_s^2 R^2.$$

Letting  $\delta_s = f(\boldsymbol{x}_s) - f(\boldsymbol{x}^*)$ , this implies that

$$\delta_{s+1} \le (1 - \gamma_s)\delta_s + \gamma_s\epsilon_s + \frac{\beta}{2}\gamma_s^2 R^2.$$

Unrolling this backwards gives

$$\delta_{T+1} \leq (1 - \gamma_T)\delta_T + \gamma_T \epsilon_T + \frac{\beta}{2} \gamma_T^2 R^2$$

$$\leq (1 - \gamma_T)(1 - \gamma_{T-1})\delta_{T-1} + (1 - \gamma_T)(\gamma_{T-1}\epsilon_{T-1} + \frac{\beta}{2} \gamma_{T-1}^2 R^2) + \gamma_T \epsilon_T + \frac{\beta}{2} \gamma_T^2 R^2$$

$$\leq \sum_{t=1}^T \left( \prod_{s=t+1}^T (1 - \gamma_s) \right) (\gamma_t \epsilon_t + \frac{\beta}{2} \gamma_t^2 R^2).$$

We can write

$$\prod_{s=t+1}^{T} (1 - \gamma_s) = \prod_{s=t+1}^{T} \frac{s}{s+1} = \frac{t+1}{T+1}$$

so

$$\sum_{t=1}^{T} \left( \prod_{s=t+1}^{T} (1 - \gamma_s) \right) \frac{\beta}{2} \gamma_t^2 R^2 = \sum_{t=1}^{T} \frac{t+1}{T+1} \frac{\beta}{2} \frac{1}{(t+1)^2} R^2$$
$$= \frac{\beta R^2}{2(T+1)} \sum_{t=1}^{T} \frac{1}{t+1}$$

$$\leq \frac{\beta R^2(\log T + 1)}{2(T+1)}$$

and

$$\sum_{t=1}^{T} \left( \prod_{s=t+1}^{T} (1 - \gamma_s) \right) \gamma_t \epsilon_t = \sum_{t=1}^{T} \frac{t+1}{T+1} \frac{1}{t+1} \epsilon_t$$
$$= \frac{1}{T+1} \sum_{t=1}^{T} \epsilon_t$$

which proves the result.

**Lemma C.2.** When running Algorithm 4, we have

$$oldsymbol{x}_{T+1} = rac{1}{T+1} \left( \sum_{t=1}^{T} oldsymbol{y}_t + oldsymbol{x}_1 
ight).$$

*Proof.* We have:

$$\begin{aligned} \boldsymbol{x}_{T+1} &= \sum_{t=1}^{T} \left( \prod_{s=t+1}^{T} (1 - \gamma_s) \right) \gamma_t \boldsymbol{y}_t + \left( \prod_{s=1}^{T} (1 - \gamma_s) \right) \boldsymbol{x}_1 \\ &= \sum_{t=1}^{T} \frac{t+1}{T+1} \frac{1}{t+1} \boldsymbol{y}_t + \frac{1}{T+1} \boldsymbol{x}_1 \\ &= \frac{1}{T+1} \sum_{t=1}^{T} \boldsymbol{y}_t + \frac{1}{T+1} \boldsymbol{x}_1. \end{aligned}$$

# C.2 Online Frank-Wolfe via Regret Minimization

### Algorithm 5 Online Frank-Wolfe via Regret Minimization (FWREGRET)

- 1: **input**: function to optimize f, number of iterates T, episodes per iterate K
- 2: Play any policy for K episodes, denote collected covariates as  $\Gamma_0$ , collected data as  $\mathfrak{D}_0$
- 3:  $\Lambda_1 \leftarrow K^{-1}\Gamma_0$
- 4: **for** t = 1, 2, ..., T **do**
- 5: Set  $\gamma_t \leftarrow \frac{1}{t+1}$
- 6: Run REGMIN on reward  $r_h^t(s, a) = \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \cdot \boldsymbol{\phi}(s, a) \boldsymbol{\phi}(s, a)^{\top})/M$  for K episodes, denote collected covariates as  $\Gamma_t$ , collected data as  $\mathfrak{D}_t$
- 7:  $\mathbf{\Lambda}_{t+1} \leftarrow (1 \gamma_t) \mathbf{\Lambda}_t + \gamma_t K^{-1} \mathbf{\Gamma}_t$
- 8: **return**  $\Lambda_{T+1}$ ,  $\cup_{t=0}^{T} \mathfrak{D}_t$

**Lemma C.3.** Consider running Algorithm 5 with a function f satisfying Definition 5.1 and a regret minimization algorithm satisfying Definition C.2. Denote  $K_0(T, \beta, M, \delta)$  the minimum integer value of K satisfying

$$K \geq \max\left\{\frac{72T^2M^2\log(4T/\delta)}{\beta^2R^4}, \frac{8T^2M^2\mathcal{C}_1\log^{p_1}(2HKT/\delta)}{\beta^2R^4}, \frac{3TM\mathcal{C}_2\log^{p_2}(2HKT/\delta)}{\beta R^2}\right\}.$$

Then as long as  $K \geq K_0(T, \beta, M, \delta)$ , we have that, with probability at least  $1 - \delta$ ,

$$f(\mathbf{\Lambda}_{T+1}) - \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} f(\mathbf{\Lambda}) \le \frac{\beta R^2 (\log T + 3)}{2(T+1)}$$

for  $R = \sup_{\pi,\pi'} \|\mathbf{\Lambda}_{\pi} - \mathbf{\Lambda}_{\pi'}\|$ .

Proof. Note that by Lemma C.2 and since  $\|\phi(s,a)\|_2 \leq 1$ , we can bound  $\|\mathbf{\Lambda}_t\|_{\text{op}} \leq 1$  and  $\|\phi(s,a)\phi(s,a)^{\top}\|_{\text{op}} \leq 1$ . Definition 5.1 it follows that  $r_h^t(s,a) \in [0,1]$  for all s,a, since  $\text{tr}(\Xi_{\mathbf{\Lambda}_t} \cdot \phi(s,a)\phi(s,a)^{\top}) \leq \|\phi(s,a)\phi(s,a)^{\top}\|_{\text{op}} \cdot \text{tr}(\Xi_{\mathbf{\Lambda}_t}) \leq \text{tr}(\Xi_{\mathbf{\Lambda}_t}) \leq M$ , and  $\text{tr}(\Xi_{\mathbf{\Lambda}_t} \cdot \phi(s,a)\phi(s,a)^{\top}) \geq 0$  since  $\Xi_{\mathbf{\Lambda}_t} \succeq 0$ . If we run RegMin for K episodes on reward function  $r_h^t$ , by Definition 5.1 and Definition C.2 we then have that, with probability at least  $1 - \delta/2T$ ,

$$\sqrt{C_1 K \log^{p_1}(2HKT/\delta)} + C_2 \log^{p_2}(2HKT/\delta) \ge K \sup_{\pi} \mathbb{E}_{\pi} [\operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \cdot \boldsymbol{\phi} \boldsymbol{\phi}^{\top})/M] - \sum_{k=1}^{K} \mathbb{E}_{\pi_k} [\operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \cdot \boldsymbol{\phi} \boldsymbol{\phi}^{\top})/M]$$

$$= K \sup_{\pi} \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \mathbf{\Lambda}_{\pi})/M - K \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \cdot K^{-1} \sum_{k=1}^{K} \mathbf{\Lambda}_{\pi_k})/M$$

which implies

$$\sqrt{\frac{M^2 \mathcal{C}_1 \log^{p_1}(2HKT/\delta)}{K}} + \frac{M \mathcal{C}_2 \log^{p_2}(2HKT/\delta)}{K} \ge \sup_{\pi} \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \mathbf{\Lambda}_{\pi}) - \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \cdot K^{-1} \sum_{k=1}^{K} \mathbf{\Lambda}_{\pi_k}).$$

Furthermore, we have that

$$\left| \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \cdot K^{-1} \sum_{k=1}^K \mathbf{\Lambda}_{\pi_k}) - \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \cdot K^{-1} \mathbf{\Gamma}_t) \right| = \left| \frac{1}{K} \sum_{k=1}^K \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \mathbf{\Lambda}_{\pi_k}) - \frac{1}{K} \sum_{k=1}^K \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \boldsymbol{\phi}_k \boldsymbol{\phi}_k^{\top}) \right|.$$

Note that  $\mathbb{E}_{\pi_k}[\operatorname{tr}(\Xi_{\Lambda_t}\phi_k\phi_k^{\top})] = \operatorname{tr}(\Xi_{\Lambda_t}\Lambda_{\pi_k})$ ,  $\operatorname{tr}(\Xi_{\Lambda_t}\phi_k\phi_k^{\top}) \in [0, M]$ , and  $\pi_k$  is  $\mathcal{F}_{k-1}$ -measurable. We can therefore apply Azuma-Hoeffding (Lemma A.4) to get that, with probability at least  $1 - \delta/2T$ ,

$$\left| \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \cdot K^{-1} \sum_{k=1}^K \mathbf{\Lambda}_{\pi_k}) - \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \cdot K^{-1} \mathbf{\Gamma}_t) \right| \leq \sqrt{\frac{8M^2 \log(4T/\delta)}{K}}.$$

Therefore,

$$\sqrt{\frac{8M^2 \log(4T/\delta)}{K}} + \sqrt{\frac{M^2 \mathcal{C}_1 \log^{p_1}(2HKT/\delta)}{K}} + \frac{M \mathcal{C}_2 \log^{p_2}(2HKT/\delta)}{K} \\
\geq \sup_{\pi} \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \mathbf{\Lambda}_{\pi}) - \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \cdot K^{-1} \mathbf{\Gamma}_t).$$

Given our condition on K, we have

$$\sqrt{\frac{8M^2\log(4T/\delta)}{K}} + \sqrt{\frac{M^2\mathcal{C}_1\log^{p_1}(2HKT/\delta)}{K}} + \frac{M\mathcal{C}_2\log^{p_2}(2HKT/\delta)}{K} \leq \frac{\beta R^2}{T}$$

which implies

$$\sup_{\pi} \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \mathbf{\Lambda}_{\pi}) - \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t} \cdot K^{-1} \mathbf{\Gamma}_t) \le \frac{\beta R^2}{T}.$$
 (C.1)

Note that, for any  $\Lambda \in \Omega$ , we have

$$\operatorname{tr}(\Xi_{\mathbf{\Lambda}_t}\mathbf{\Lambda}) = \operatorname{tr}(\Xi_{\mathbf{\Lambda}_t}\mathbb{E}_{\pi\sim\omega}[\mathbf{\Lambda}_{\pi}]) = \mathbb{E}_{\pi\sim\omega}[\operatorname{tr}(\Xi_{\mathbf{\Lambda}_t}\mathbf{\Lambda}_{\pi})]$$

SO

$$\sup_{\boldsymbol{\Lambda} \in \boldsymbol{\Omega}} \operatorname{tr}(\boldsymbol{\Xi}_{\boldsymbol{\Lambda}_t} \boldsymbol{\Lambda}) = \sup_{\boldsymbol{\omega} \in \boldsymbol{\Omega}_{\pi}} \mathbb{E}_{\boldsymbol{\pi} \sim \boldsymbol{\omega}} [\operatorname{tr}(\boldsymbol{\Xi}_{\boldsymbol{\Lambda}_t} \boldsymbol{\Lambda}_{\pi})] = \sup_{\boldsymbol{\pi}} \operatorname{tr}(\boldsymbol{\Xi}_{\boldsymbol{\Lambda}_t} \boldsymbol{\Lambda}_{\pi})$$

By definition,  $\Xi_{\Lambda_t} = -\nabla_{\Lambda} f(\Lambda)|_{\Lambda = \Lambda_t}$ , so it follows that

$$-\sup_{\mathbf{\Lambda}'\in\mathbf{\Omega}}\operatorname{tr}(\Xi_{\mathbf{\Lambda}_t}\mathbf{\Lambda}')=\inf_{\mathbf{\Lambda}'\in\mathbf{\Omega}}\operatorname{tr}(\nabla_{\mathbf{\Lambda}}f(\mathbf{\Lambda})|_{\mathbf{\Lambda}=\mathbf{\Lambda}_t}\cdot\mathbf{\Lambda}').$$

It follows that (C.1) is precisely the guarantee required on  $y_t$  by Algorithm 4 with  $\epsilon_t = \frac{\beta R^2}{T}$ . Since f is  $\beta$ -smooth by Definition 5.1 and since the set  $\Omega$  is convex and compact by Lemma A.9, we can apply Lemma C.1 with a union bound over t to get the result.

### C.3 Data Collection via Online Frank-Wolfe

### Algorithm 6 Collect Optimal Covariates (OPTCOV)

- 1: **input**: functions to optimize  $(f_i)_i$ , constraint tolerance  $\epsilon$ , confidence  $\delta$
- 2: **for**  $i = 1, 2, 3, \dots$  **do**
- 3:  $T_i \leftarrow 2^i, K_i \leftarrow 2^i T_i^2$
- 4: **if**  $K_i \geq \widetilde{K}_0(T_i, \beta_i, M_i, \frac{\delta}{4i^2})T_i^2 + \widetilde{K}_1(T_i, \beta_i, M_i, \frac{\delta}{4i^2})T_i$  for  $\widetilde{K}_0$  and  $\widetilde{K}_1$  as in Lemma C.5 **then**
- 5:  $\widehat{\mathbf{\Lambda}}, \mathfrak{D}_i \leftarrow \text{FWREGRET}(f_i, T_i 1, K_i)$
- 6: if  $f_i(\widehat{\mathbf{\Lambda}}) \leq K_i T_i \epsilon$  and  $f_i(\widehat{\mathbf{\Lambda}}) \geq \frac{\beta_i R^2 (\log T_i + 3)}{T_i}$  then
- 7: return  $\widehat{\mathbf{\Lambda}}$ ,  $K_iT_i$ ,  $\mathfrak{D}_i$

**Theorem 8.** Let  $(f_i)_i$  denote some sequence of functions which satisfy Definition 5.1 with constants  $(\beta_i, L_i, M_i)$  and assume  $\beta_i \geq 1$ . Let  $(\beta, L, M)$  be some values such that  $\beta_i \leq \beta, L_i \leq L, M_i \leq M$  for all i, and let f be some function such that  $f_i(\mathbf{\Lambda}) \leq f(\mathbf{\Lambda})$  for all i and  $\mathbf{\Lambda} \succeq 0$ . Denote  $f_{\min}$  a lower bound on all  $f_i$ :  $\min_i \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} f_i(\mathbf{\Lambda}) \geq f_{\min}$ .

Define

$$N^{\star}(\epsilon; f) := \frac{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} f(\mathbf{\Lambda})}{\epsilon}.$$
 (C.2)

Then, if we run Algorithm 6 on  $(f_i)_i$  with constraint tolerance  $\epsilon$  and confidence  $\delta$ , we have that with probability at least  $1 - \delta$ , it will run for at most

$$5N^{\star}(\epsilon; f) + \text{poly}\left(2^{p_1+p_2}, \mathcal{C}_1, \mathcal{C}_2, M, \beta, R, L, f_{\min}^{-1}, \log 1/\delta\right)$$

episodes, and will return data  $\{\phi_{\tau}\}_{\tau=1}^{N}$  with covariance  $\widehat{\Sigma}_{N} = \sum_{\tau=1}^{N} \phi_{\tau} \phi_{\tau}^{\top}$  such that

$$f_{\widehat{i}}(N^{-1}\widehat{\Sigma}_N) \le N\epsilon,$$

where  $\hat{i}$  is the iteration on which OptCov terminates.

Corollary 6 (Theorem 6). Instantiating REGMIN with the computationally efficient version of the FORCE algorithm of Wagenmaker et al. (2021a), we obtain a complexity of

$$5N^{\star}(\epsilon; f) + \text{poly}\left(d, H, M, \beta, R, L, f_{\min}^{-1}, \log 1/\delta\right).$$

*Proof.* This result is immediate since Force satisfies Definition C.2 with

$$C_1 = c_1 d^4 H^4$$
,  $C_2 = c_2 d^4 H^3$ ,  $p_1 = 3$ ,  $p_2 = 7/2$ 

for universal numerical constants  $c_1$  and  $c_2$ .

Proof of Theorem 8. We first show that the condition  $f_i(\widehat{\Lambda}) \geq \frac{\beta R^2(\log T_i + 3)}{T_i}$  is sufficient to ensure a 2-approximate minimum of  $f_i$ , and then show a sufficient condition on  $K_i$  and  $T_i$  that will guarantee the condition on Line 6 is met.

Guaranteeing 2-optimality. We first show that for a fixed i, the condition  $f_i(\widehat{\Lambda}) \geq \frac{\beta_i R^2(\log T_i + 3)}{T_i}$  will only be met once

$$f_i(\widehat{\boldsymbol{\Lambda}}) \leq 2 \cdot \inf_{\boldsymbol{\Lambda} \in \Omega} f_i(\boldsymbol{\Lambda})$$

and that it will take at most

$$T_i \ge \frac{2\beta R^2(\log T_i + 3)}{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} f_i(\mathbf{\Lambda})}$$

iterations to do so, as long as

$$T_i K_i \ge \frac{L^2}{2(d \log(1 + 8\sqrt{T_i K_i}) + \log(4i^2/\delta)) \cdot (\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} f_i(\mathbf{\Lambda}))^2}.$$

The first part follows by applying Lemma C.3. Note that the if statement on Line 4 will only be met once

$$K_i \ge K_0(T_i, \beta_i, M_i, \delta/4i^2).$$

This follows by Lemma C.5. Thus, the condition on  $K_i$  required by Lemma C.3 will be met, so it follows that with probability at least  $1 - \delta/(4i^2)$ ,

$$f_i(\widehat{\mathbf{\Lambda}}) - \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} f_i(\mathbf{\Lambda}) \le \frac{\beta_i R^2 (\log T_i + 3)}{2T_i}.$$

Therefore, if  $f_i(\widehat{\boldsymbol{\Lambda}}) \geq \frac{\beta_i R^2(\log T_i + 3)}{T_i}$ , we have

$$f_i(\widehat{\mathbf{\Lambda}}) - \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} f_i(\mathbf{\Lambda}) \le \frac{1}{2} f_i(\widehat{\mathbf{\Lambda}}) \implies \frac{1}{2} f_i(\widehat{\mathbf{\Lambda}}) \le \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} f_i(\mathbf{\Lambda})$$

$$\implies f_i(\widehat{\mathbf{\Lambda}}) \leq 2 \cdot \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} f_i(\mathbf{\Lambda}).$$

We will show a sufficient condition for  $f_i(\widehat{\Lambda}) \geq \frac{\beta R^2(\log T_i + 3)}{T_i}$ , which implies that  $f_i(\widehat{\Lambda}) \geq \frac{\beta_i R^2(\log T_i + 3)}{T_i}$  since  $\beta_i \leq \beta$ . By Lemma C.2 and the procedure run by Algorithm 5, we have that  $\widehat{\Lambda} = \frac{1}{T_i K_i} \sum_{\tau=1}^{T_i K_i} \phi_{\tau} \phi_{\tau}^{\top}$  where at episodes  $\tau$  we run some  $\mathcal{F}_{\tau-1}$ -measurable policy  $\pi_{\tau}$  to acquire  $\phi_{\tau}$ . Now if  $\widehat{\Lambda} = \widetilde{\Lambda}$  for some  $\widetilde{\Lambda} \in \Omega$ , then the second part follows trivially since  $\inf_{\Lambda \in \Omega} f_i(\Lambda) \leq f_i(\widetilde{\Lambda})$ , so a sufficient condition for  $f_i(\widehat{\Lambda}) \geq \frac{\beta R^2(\log T_i + 3)}{T_i}$  is that  $\inf_{\Lambda \in \Omega} f_i(\Lambda) \geq \frac{\beta R^2(\log T_i + 3)}{T_i}$ . However, since  $\widehat{\Lambda}$  is stochastic, we may not have that  $\widehat{\Lambda} \in \Omega$ . Let  $\widetilde{\Lambda} := \frac{1}{T_i K_i} \sum_{\tau=1}^{T_i K_i} \Lambda_{\pi_{\tau}}$  and note that  $\widetilde{\Lambda} \in \Omega$ . Applying Lemma C.4, we have that with probability at least  $1 - \delta/(4i^2)$ ,

$$\left\|\widetilde{\mathbf{\Lambda}} - \widehat{\mathbf{\Lambda}}\right\|_{\text{op}} \le \sqrt{\frac{8d\log(1 + 8\sqrt{T_i K_i}) + 8\log(4i^2/\delta)}{T_i K_i}}$$

for  $\widetilde{\pi}$  the uniform mixture of  $\{\pi_{\tau}\}_{\tau=1}^{T_iK_i}$ . By the Lipschitz condition of Definition 5.1, this implies

$$f_{i}(\widehat{\boldsymbol{\Lambda}}) \geq f_{i}(\widetilde{\boldsymbol{\Lambda}}) - L_{i} \|\widehat{\boldsymbol{\Lambda}} - \widetilde{\boldsymbol{\Lambda}}\|_{\text{op}}$$

$$\geq f_{i}(\widetilde{\boldsymbol{\Lambda}}) - L \|\widehat{\boldsymbol{\Lambda}} - \widetilde{\boldsymbol{\Lambda}}\|_{\text{op}}$$

$$\geq f_{i}(\widetilde{\boldsymbol{\Lambda}}) - L \sqrt{\frac{8d \log(1 + 8\sqrt{T_{i}K_{i}}) + 8\log(4i^{2}/\delta)}{T_{i}K_{i}}}$$

$$\geq \inf_{\boldsymbol{\Lambda} \in \Omega} f_{i}(\boldsymbol{\Lambda}) - L \sqrt{\frac{8d \log(1 + 8\sqrt{T_{i}K_{i}}) + 8\log(4i^{2}/\delta)}{T_{i}K_{i}}}.$$

Thus, a sufficient condition for  $f_i(\widehat{\mathbf{\Lambda}}) \geq \frac{\beta R^2(\log T_i + 3)}{T_i}$  is that

$$\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} f_i(\mathbf{\Lambda}) - L \sqrt{\frac{8d \log(1 + 8\sqrt{T_i K_i}) + 8\log(4i^2/\delta)}{T_i K_i}} \ge \frac{\beta R^2(\log T_i + 3)}{T_i}$$

$$\iff T_i \ge \frac{\beta R^2(\log T_i + 3)}{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} f_i(\mathbf{\Lambda}) - L \sqrt{\frac{8d \log(1 + 8\sqrt{T_i K_i}) + 8\log(4i^2/\delta)}{T_i K_i}}}.$$

If

$$T_i K_i \ge \frac{L^2}{2(d \log(1 + 8\sqrt{T_i K_i}) + \log(4i^2/\delta)) \cdot (\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} f_i(\mathbf{\Lambda}))^2}$$

it follows that a sufficient condition is

$$T_i \ge \frac{2\beta R^2(\log T_i + 3)}{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} f_i(\mathbf{\Lambda})}.$$

Union bounding over the events considered above for all i, we have that the total probability of failure is bounded as

$$\sum_{i=1}^{\infty} \left(\frac{\delta}{4i^2} + \frac{\delta}{4i^2}\right) = \frac{\pi^2}{12}\delta \le \delta.$$

**Termination Guarantee.** We next show a sufficient condition to ensure that the if statements on Line 4 and Line 6 are met.

Assume the if statement on Line 4 has been met and that we are in the regime where

$$T_i K_i \ge \frac{L^2}{2(d \log(1 + 8\sqrt{T_i K_i}) + \log(4i^2/\delta)) \cdot f_{\min}^2}, \quad T_i \ge \frac{2\beta R^2(\log T_i + 3)}{f_{\min}}.$$
 (C.3)

By the argument above and since  $\inf_{\Lambda \in \Omega} f_i(\Lambda) \geq f_{\min}$ , these conditions are sufficient to guarantee a 2-optimal solutions has been found, that is,

$$f_i(\widehat{\boldsymbol{\Lambda}}) \leq 2 \cdot \inf_{\boldsymbol{\Lambda} \in \Omega} f_i(\boldsymbol{\Lambda}),$$

and that the condition  $f_i(\widehat{\Lambda}) \geq \frac{\beta R^2(\log T_i + 3)}{T_i}$  has been met. Thus, if (C.3) holds, a sufficient condition for  $f_i(\widehat{\Lambda}) \leq T_i K_i \epsilon$  is

$$2 \cdot \inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} f_i(\mathbf{\Lambda}) \leq T_i K_i \epsilon.$$

It follows that this condition will be met (assuming (C.3) holds) once  $T_iK_i \geq N^*(\frac{\epsilon}{2}; f_i)$ . Since  $f_i \leq f$ ,  $N^*(\frac{\epsilon}{2}; f_i) \leq N^*(\frac{\epsilon}{2}; f)$ , so a sufficient condition is that  $T_iK_i \geq N^*(\frac{\epsilon}{2}; f)$ .

To upper bound the total complexity, it suffices then to guarantee that we run for enough epochs so that

$$K_i = 2^{3i} \ge \widetilde{K}_0(T_i, \beta_i, M_i, \frac{\delta}{4i^2})T_i^2 + \widetilde{K}_1(T_i, \beta_i, M_i, \frac{\delta}{4i^2})T_i$$
 (C.4)

$$T_i K_i = 2^{4i} \ge \frac{L^2}{2(d \log(1 + 8\sqrt{T_i K_i}) + \log(4i^2/\delta)) \cdot f_{\min}^2}$$
 (C.5)

$$T_i = 2^i \ge \frac{2\beta R^2(\log T_i + 3)}{f_{\min}} \tag{C.6}$$

$$T_i K_i = 2^{4i} \ge N^*(\frac{\epsilon}{2}; f). \tag{C.7}$$

Here (C.4) guarantees the if statement on Line 4 is met, and (C.5)-(C.7) guarantee the if statement on line Line 6 is met.

By assumption,  $M_i \leq M$  and  $\beta_i \geq 1$ , and note that  $\widetilde{K}_0(T_i, \beta_i, M_i, \frac{\delta}{4i^2})$  and  $\widetilde{K}_1(T_i, \beta_i, M_i, \frac{\delta}{4i^2})$  are both increasing in  $M_i$  and decreasing in  $\beta_i$ . Thus, a sufficient condition to ensure (C.4) is met is

$$2^{3i} \ge \widetilde{K}_0(2^i, 1, M, \frac{\delta}{4i^2})2^{2i} + \widetilde{K}_1(2^i, 1, M, \frac{\delta}{4i^2})2^i.$$
 (C.8)

Some calculation shows that

$$\widetilde{K}_0(2^i, 1, M, \frac{\delta}{4i^2}) \le (5i)^{p_1} \widetilde{K}_0(2, 1, M, \frac{\delta}{4}), \quad \widetilde{K}_1(2^i, 1, M, \frac{\delta}{4i^2}) \le (4i)^{p_2} \widetilde{K}_1(2, 1, M, \frac{\delta}{4})$$

so a sufficient condition to meet (C.8) is

$$2^{i} \ge 2(5i)^{p_1} \widetilde{K}_0(2, 1, M, \frac{\delta}{4}), \quad 2^{2i} \ge 2(4i)^{p_2} \widetilde{K}_1(2, 1, M, \frac{\delta}{4}).$$

By Lemma A.2 and some calculation, this will be met once

$$i \ge \max \left\{ 4p_1 \log_2(2p_1) + 2\log_2(2(5)^{p_1}\widetilde{K}_0(2,1,M,\frac{\delta}{4})), 2p_2 \log_2(p_2) + 2\log_2(2(4)^{p_2}\widetilde{K}_1(2,1,M,\frac{\delta}{4})) \right\} =: i_0.$$

To meet (C.5) it suffices to take

$$i \ge \frac{1}{4} \log_2 \frac{L^2}{df_{\min}^2} =: i_1$$

By Lemma A.2, a sufficient condition to meet (C.6) is that

$$T_i \ge \max\left\{\frac{6\beta R^2}{f_{\min}}, \frac{4\beta R^2}{f_{\min}}\log\frac{4\beta R^2}{f_{\min}}\right\}$$

so it suffices that

$$i \ge \log_2\left(\frac{6\beta R^2}{f_{\min}}\log\frac{4\beta R^2}{f_{\min}}\right) =: i_2.$$

Finally, to meet (C.7), it suffices that

$$i \ge \frac{1}{4} \log_2 N^*(\epsilon/2; f) =: i_3.$$

If we terminate at epoch  $\hat{i}$ , the total sample complexity will be bounded by

$$\sum_{i=1}^{\hat{i}} T_i K_i = \sum_{i=1}^{\hat{i}} 2^{4i} \le \frac{16}{15} \cdot 2^{4\hat{i}}.$$

By the above argument, we can bound  $\hat{i} \leq \lceil \max\{i_0, i_1, i_2, i_3\} \rceil$ . Furthermore, we see that

$$\begin{split} 2^{4\lceil i_0 \rceil} &= \operatorname{poly}\left(2^{p_1}, 2^{p_2}, M, \mathcal{C}_1, \mathcal{C}_2, \log 1/\delta\right) \\ 2^{4\lceil i_1 \rceil} &= \operatorname{poly}(L, f_{\min}^{-1}) \\ 2^{4\lceil i_2 \rceil} &= \operatorname{poly}(\beta, R, f_{\min}^{-1}) \\ 2^{4\lceil i_3 \rceil} &\leq 2N^{\star}(\epsilon/2; f) \end{split}$$

so we can bound the total sample complexity by

$$\frac{16}{15} \cdot 2^{4\lceil \max\{i_0, i_1, i_2, i_3\} \rceil} \le \frac{32}{15} N^{\star}(\epsilon/2; f) + \text{poly}\left(2^{p_1}, 2^{p_2}, \beta, R, L, f_{\min}^{-1}, M, \mathcal{C}_1, \mathcal{C}_2, \log 1/\delta\right).$$

This completes the proof since  $N^{\star}(\frac{\epsilon}{2};f) = 2N^{\star}(\epsilon;f)$  and since, by Lemma C.2,  $\widehat{\mathbf{\Lambda}}$  is simply the average of the observed feature vectors:  $\widehat{\mathbf{\Lambda}} = \frac{1}{T_i K_i} \sum_{\tau=1}^{T_i K_i} \phi_{\tau} \phi_{\tau}^{\top}$ .

**Lemma C.4.** Let  $\Lambda_K$  denote the time-normalized covariates obtained by playing policies  $\{\pi_k\}_{k=1}^K$ , where  $\pi_k$  is  $\mathcal{F}_{k-1}$ -measurable. Then, with probability at least  $1-\delta$ ,

$$\left\| \frac{1}{K} \sum_{k=1}^{K} \mathbf{\Lambda}_{\pi_k} - \mathbf{\Lambda}_K \right\|_{\text{op}} \leq \sqrt{\frac{8d \log(1 + 8\sqrt{K}) + 8 \log 1/\delta}{K}}.$$

*Proof.* Let  $\mathcal{V}$  denote an  $\epsilon$ -net of  $\mathcal{S}^{d-1}$ , for some  $\epsilon$  to be chosen. Then,

$$\left\| \frac{1}{K} \sum_{k=1}^{K} \mathbf{\Lambda}_{\pi_{k}} - \mathbf{\Lambda}_{K} \right\|_{\text{op}} = \sup_{\mathbf{v} \in \mathcal{S}^{d-1}} \left| \mathbf{v}^{\top} \left( \frac{1}{K} \sum_{k=1}^{K} \mathbf{\Lambda}_{\pi_{k}} - \mathbf{\Lambda}_{K} \right) \mathbf{v} \right|$$

$$\leq \sup_{\widetilde{\mathbf{v}} \in \mathcal{V}} \left| \widetilde{\mathbf{v}}^{\top} \left( \frac{1}{K} \sum_{k=1}^{K} \mathbf{\Lambda}_{\pi_{k}} - \mathbf{\Lambda}_{K} \right) \widetilde{\mathbf{v}} \right|$$

$$+ \sup_{\mathbf{v} \in \mathcal{S}^{d-1}} \inf_{\widetilde{\mathbf{v}} \in \mathcal{V}} \left| \mathbf{v}^{\top} \left( \frac{1}{K} \sum_{k=1}^{K} \mathbf{\Lambda}_{\pi_{k}} - \mathbf{\Lambda}_{K} \right) \mathbf{v} - \widetilde{\mathbf{v}}^{\top} \left( \frac{1}{K} \sum_{k=1}^{K} \mathbf{\Lambda}_{\pi_{k}} - \mathbf{\Lambda}_{K} \right) \widetilde{\mathbf{v}} \right|.$$

$$(b)$$

Via a union bound over V and application of Azuma-Hoeffding, we can bound, with probability at least  $1 - \delta$ ,

$$(a) \le \sqrt{\frac{2\log |\mathcal{V}|/\delta}{K}}.$$

We can bound (b) as

$$(b) \leq \sup_{\boldsymbol{v} \in \mathcal{S}^{d-1}} \inf_{\widetilde{\boldsymbol{v}} \in \mathcal{V}} 2 \left| \boldsymbol{v}^{\top} \left( \frac{1}{K} \sum_{k=1}^{K} \boldsymbol{\Lambda}_{\pi_{k}} - \boldsymbol{\Lambda}_{K} \right) (\boldsymbol{v} - \widetilde{\boldsymbol{v}}) \right|$$

$$\leq \sup_{\boldsymbol{v} \in \mathcal{S}^{d-1}} \inf_{\widetilde{\boldsymbol{v}} \in \mathcal{V}} 2 \|\boldsymbol{v} - \widetilde{\boldsymbol{v}}\|_{2} \left\| \frac{1}{K} \sum_{k=1}^{K} \boldsymbol{\Lambda}_{\pi_{k}} - \boldsymbol{\Lambda}_{K} \right\|_{\text{op}}$$

$$\leq 4\epsilon$$

where the last inequality follows since  $\|\frac{1}{K}\sum_{k=1}^K \mathbf{\Lambda}_{\pi_k}\|_{\text{op}} \leq 1$ , and  $\|\frac{1}{K}\mathbf{\Lambda}_K\|_{\text{op}} \leq 1$ , and since  $\mathcal{V}$  is an  $\epsilon$ -net. Setting  $\epsilon = 1/(4\sqrt{K})$ , Lemma A.1 gives that  $|\mathcal{V}| \leq (1 + 8\sqrt{K})^d$ , and we conclude that with probability at least  $1 - \delta$ :

$$\left\| \frac{1}{K} \sum_{k=1}^{K} \mathbf{\Lambda}_{\pi_k} - \mathbf{\Lambda}_K \right\|_{\text{op}} \leq \sqrt{\frac{2 \log |\mathcal{V}|/\delta}{K}} + 4\epsilon$$

$$\leq \sqrt{\frac{2d \log(1 + 8\sqrt{K}) + 2 \log 1/\delta}{K}} + \frac{1}{\sqrt{K}}$$

$$\leq 2\sqrt{\frac{2d \log(1 + 8\sqrt{K}) + 2 \log 1/\delta}{K}}.$$

Lemma C.5. We can bound

$$K_0(T, \beta, M, \delta) \le \widetilde{K}_0(T, \beta, M, \delta)T^2 + \widetilde{K}_1(T, \beta, M, \delta)T$$

for

$$\widetilde{K}_0(T,\beta,M,\delta) := \max \left\{ \frac{72M^2 \log(4T/\delta)}{\beta^2 R^4}, \frac{8M^2 \mathcal{C}_1}{\beta^2 R^4} \cdot (2p_1)^{p_1} \log^{p_1} \left( \frac{32p_1 H T^3 M^2 \mathcal{C}_1}{\beta^2 R^4 \delta} \right) \right\}$$

$$\widetilde{K}_1(T,\beta,M,\delta) := \frac{3M\mathcal{C}_2}{\beta R^2} \cdot (2p_2)^{p_2} \log^{p_2} \left( \frac{12p_2HT^2M\mathcal{C}_2}{\beta R^2 \delta} \right),$$

*Proof.* By definition  $K_0(T, \beta, M, \delta)$  is the smallest integer value of K that satisfies:

$$K \ge \max\left\{\frac{72T^2M^2\log(4T/\delta)}{\beta^2R^4}, \frac{8T^2M^2C_1\log^{p_1}(2HKT/\delta)}{\beta^2R^4}, \frac{3TMC_2\log^{p_2}(2HKT/\delta)}{\beta R^2}\right\}.$$
 (C.9)

By Lemma A.2, we have that if

$$K \ge \frac{8T^2M^2\mathcal{C}_1}{\beta^2R^4} \cdot (2p_1)^{p_1} \log^{p_1} \left( \frac{8T^2M^2\mathcal{C}_1}{\beta^2R^4} \cdot \frac{4p_1HT}{\delta} \right), \quad K \ge \frac{3TM\mathcal{C}_2}{\beta R^2} \cdot (2p_2)^{p_2} \log^{p_2} \left( \frac{3TM\mathcal{C}_2}{\beta R^2} \cdot \frac{4p_2HT}{\delta} \right)$$

and

$$K \ge \frac{72T^2M^2\log(4T/\delta)}{\beta^2R^4}$$

then Equation (C.9) will be satisfied. Some algebra gives the result.

# D XY-Optimal Design

We are interested in optimizing the function

$$\mathsf{XY}_{\mathrm{opt}}(\mathbf{\Lambda}) = \max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2 \quad \text{for} \quad \mathbf{A}(\mathbf{\Lambda}) = \mathbf{\Lambda} + \mathbf{\Lambda}_0$$

with  $\Lambda_0 \succ 0$  some fixed regularizer. This objective, however, is not smooth, so we relax it to the following:

$$\widetilde{\mathsf{XY}}_{\mathrm{opt}}(\boldsymbol{\Lambda}) := \mathrm{LogSumExp}\left(\left\{e^{\eta\|\boldsymbol{\phi}\|_{\mathbf{A}(\boldsymbol{\Lambda})^{-1}}^{2}}\right\}_{\boldsymbol{\phi} \in \Phi}; \eta\right) = \frac{1}{\eta} \log \left(\sum_{\boldsymbol{\phi} \in \Phi} e^{\eta\|\boldsymbol{\phi}\|_{\mathbf{A}(\boldsymbol{\Lambda})^{-1}}^{2}}\right). \tag{D.1}$$

We first offer some properties on how well  $\widetilde{\mathsf{XY}}_{\mathrm{opt}}(\boldsymbol{\Lambda})$  approximates  $\mathsf{XY}_{\mathrm{opt}}(\boldsymbol{\Lambda})$ , and then show that we can bound the smoothness constant of  $\widetilde{\mathsf{XY}}_{\mathrm{opt}}(\boldsymbol{\Lambda})$ . Throughout this section, we will denote  $\gamma_{\Phi} := \max_{\phi \in \Phi} \|\phi\|_2$  and let  $f(\boldsymbol{\Lambda}) := \widetilde{\mathsf{XY}}_{\mathrm{opt}}(\boldsymbol{\Lambda})$ .

# D.1 Approximating Non-Smooth Optimal Design with Smooth Optimal Design Lemma D.1.

$$|\mathsf{XY}_{\mathrm{opt}}(\mathbf{\Lambda}) - \widetilde{\mathsf{XY}}_{\mathrm{opt}}(\mathbf{\Lambda})| \leq \frac{\log |\Phi|}{\eta}, \qquad \mathsf{XY}_{\mathrm{opt}}(\mathbf{\Lambda}) \leq \widetilde{\mathsf{XY}}_{\mathrm{opt}}(\mathbf{\Lambda}).$$

*Proof.* This result is standard but we include the proof for completeness. We prove it for some generic sequence  $(a_i)_{i=1}^n$ . Take  $\eta > 0$ . Clearly,

$$\exp(\max_{i} \eta a_i) \le \sum_{i=1}^{n} \exp(\eta a_i) \le n \exp(\max_{i} \eta a_i)$$

so

$$\max_{i} \eta a_{i} \leq \log \left( \sum_{i=1}^{n} \exp(\eta a_{i}) \right) \leq \log n + \max_{i} \eta a_{i}.$$

The result follows by rearranging and dividing by  $\eta$ .

Lemma D.2. If  $\eta \geq \widetilde{\eta} \geq 0$ , then  $\widetilde{\mathsf{XY}}_{\mathrm{opt}}(\boldsymbol{\Lambda}; \eta) \leq \widetilde{\mathsf{XY}}_{\mathrm{opt}}(\boldsymbol{\Lambda}; \widetilde{\eta})$ .

*Proof.* We will prove this for some generic sequence  $(a_i)_{i=1}^n$ ,  $a_i \geq 0$ . Note that,

$$\frac{\mathrm{d}}{\mathrm{d}\eta} \frac{1}{\eta} \log \left( \sum_{i} e^{\eta a_{i}} \right) = -\frac{1}{\eta^{2}} \log \left( \sum_{i} e^{\eta a_{i}} \right) + \frac{1}{\eta} \frac{1}{\sum_{i} e^{\eta a_{i}}} \cdot \sum_{i} a_{i} e^{\eta a_{i}}.$$

We are done if we can show this is non-positive. Note that,

$$\log\left(\sum_{i} e^{\eta a_i}\right) \ge \max_{i} \log\left(e^{\eta a_i}\right) = \max_{i} \eta a_i$$

so

$$-\frac{1}{\eta^2}\log\left(\sum_{i}e^{\eta a_i}\right) + \frac{1}{\eta}\frac{1}{\sum_{i}e^{\eta a_i}} \cdot \sum_{i}a_ie^{\eta a_i} \le -\frac{1}{\eta}\max_{i}a_i + \frac{1}{\eta}\frac{1}{\sum_{i}e^{\eta a_i}} \cdot \sum_{i}a_ie^{\eta a_i}$$
$$\le -\frac{1}{\eta}\max_{i}a_i + \frac{1}{\eta}\max_{i}a_i$$
$$= 0$$

The result follows since  $\widetilde{XY}_{\mathrm{opt}}$  has this form.

Lemma D.3. We have,

$$\inf_{\boldsymbol{\Lambda}\succeq 0, \|\boldsymbol{\Lambda}\|_{\mathrm{op}}\leq 1} \mathsf{XY}_{\mathrm{opt}}(\boldsymbol{\Lambda}) \geq \frac{\gamma_{\Phi}}{1+\|\boldsymbol{\Lambda}_{0}\|_{\mathrm{op}}}.$$

*Proof.* Note that  $\|\mathbf{A}(\mathbf{\Lambda})\|_{\mathrm{op}} \leq 1 + \|\mathbf{\Lambda}_0\|_{\mathrm{op}}$ , so

$$\inf_{\mathbf{\Lambda}\succeq 0, \|\mathbf{\Lambda}\|_{\mathrm{op}}\leq 1} \max_{\boldsymbol{\phi}\in\Phi} \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2 \geq \inf_{\mathbf{\Lambda}\succeq 0, \|\mathbf{\Lambda}\|_{\mathrm{op}}\leq 1+\|\mathbf{\Lambda}_0\|_{\mathrm{op}}} \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2 \geq \frac{\max_{\boldsymbol{\phi}\in\Phi} \|\boldsymbol{\phi}\|_2}{1+\|\mathbf{\Lambda}_0\|_{\mathrm{op}}}.$$

**Lemma D.4.** Assume that we set  $\eta \geq \frac{2}{\gamma_{\Phi}}(1 + \|\mathbf{\Lambda}_0\|_{\text{op}}) \cdot \log |\Phi|$ . Then

$$N^{\star}(\epsilon; \widetilde{\mathsf{XY}}_{\mathrm{opt}}(\mathbf{\Lambda})) \leq 2N^{\star}(\epsilon; \mathsf{XY}_{\mathrm{opt}}(\mathbf{\Lambda})).$$

*Proof.* Denote  $f(\mathbf{\Lambda}) \leftarrow \text{LogSumExp}\left(\left\{e^{\eta\|\phi\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2}\right\}_{\phi \in \Phi}; \eta\right)$ . By Lemma D.1 and Lemma D.3, we have

$$\begin{split} |\max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{\mathbf{A}(\boldsymbol{\Lambda})^{-1}}^2 - f(\boldsymbol{\Lambda})| &\leq \frac{\log |\Phi|}{\eta} \leq \frac{\gamma_{\Phi}}{2(1 + \|\boldsymbol{\Lambda}_0\|_{\mathrm{op}})} \leq \min_{\boldsymbol{\Lambda} \succeq 0, \|\boldsymbol{\Lambda}\|_{\mathrm{op}} \leq 1} \frac{1}{2} f(\boldsymbol{\Lambda}) \\ &\Longrightarrow f(\boldsymbol{\Lambda}) \leq 2 \max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{\mathbf{A}(\boldsymbol{\Lambda})^{-1}}^2. \end{split}$$

Let  $\Lambda^*$  denote the matrix that minimizes  $\max_{\phi \in \Phi} \|\phi\|_{\mathbf{A}(\Lambda)^{-1}}^2$  over the constraint set:  $\max_{\phi \in \Phi} \|\phi\|_{\mathbf{A}(\Lambda^*)^{-1}}^2 = \inf_{\Lambda \in \Omega} \max_{\phi \in \Phi} \|\phi\|_{\mathbf{A}(\Lambda)^{-1}}^2$ . Then it follows that, by definition of  $N^*(\epsilon; \max_{\phi \in \Phi} \|\phi\|_{\mathbf{A}(\Lambda)^{-1}}^2)$ :

$$\max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{\mathbf{A}(\boldsymbol{\Lambda}^{\star})^{-1}}^2 \leq \epsilon \cdot N^{\star}(\epsilon; \max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{\mathbf{A}(\boldsymbol{\Lambda})^{-1}}^2).$$

However, this implies

$$\frac{1}{2}f(\mathbf{\Lambda}^{\star}) \leq \epsilon \cdot N^{\star}(\epsilon; \max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2),$$

so  $(\Lambda^*, 2N^*(\epsilon; \max_{\phi \in \Phi} \|\phi\|_{\mathbf{A}(\Lambda)^{-1}}^2))$  is a feasible solution to the optimization (C.2) for f. As  $N^*(\epsilon; f)$  is the minimum solution, it follows that  $N^*(\epsilon; f) \leq 2N^*(\epsilon; \max_{\phi \in \Phi} \|\phi\|_{\mathbf{A}(\Lambda)^{-1}}^2)$ .

# D.2 Bounding the Smoothness

**Lemma D.5.**  $f(\Lambda) = \widetilde{\mathsf{XY}}_{\mathrm{opt}}(\Lambda)$  satisfies all conditions of Definition 5.1 with

$$L = \|\boldsymbol{\Lambda}_0^{-1}\|_{\mathrm{op}}^2, \quad \beta = 2\|\boldsymbol{\Lambda}_0^{-1}\|_{\mathrm{op}}^3 (1 + \eta \|\boldsymbol{\Lambda}_0^{-1}\|_{\mathrm{op}}), \quad M = \|\boldsymbol{\Lambda}_0^{-1}\|_{\mathrm{op}}^2$$

$$\nabla_{\mathbf{\Lambda}} f(\mathbf{\Lambda}) = \left(\sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2}\right)^{-1} \cdot \sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2} \mathbf{A}(\mathbf{\Lambda})^{-1} \boldsymbol{\phi} \boldsymbol{\phi}^{\top} \mathbf{A}(\mathbf{\Lambda})^{-1} =: \Xi_{\mathbf{\Lambda}}.$$

*Proof.* Using Lemma D.6, the gradient of  $f(\Lambda)$  with respect to  $\Lambda_{ij}$  is

$$\nabla_{\mathbf{\Lambda}_{ij}} f(\mathbf{\Lambda}) = -\left(\sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2}\right)^{-1} \cdot \sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2} \boldsymbol{\phi}^{\top} \mathbf{A}(\mathbf{\Lambda})^{-1} \boldsymbol{e}_i \boldsymbol{e}_j^{\top} \mathbf{A}(\mathbf{\Lambda})^{-1} \boldsymbol{\phi}$$

from which the expression for  $\nabla_{\mathbf{\Lambda}} f(\mathbf{\Lambda})$  follows directly.

To bound the Lipschitz constant of f, by the Mean Value Theorem it suffices to bound

$$\sup_{\boldsymbol{\Lambda}, \widetilde{\boldsymbol{\Lambda}} \succeq 0, \|\boldsymbol{\Lambda}\|_{\text{op}} \le 1, \|\widetilde{\boldsymbol{\Lambda}}\|_{\text{op}} \le 1} |\text{tr}(\nabla f(\boldsymbol{\Lambda})^{\top} \widetilde{\boldsymbol{\Lambda}})| \le \left(\sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\boldsymbol{\Lambda})^{-1}}^{2}}\right)^{-1} \cdot \sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\boldsymbol{\Lambda})^{-1}}^{2} \|\mathbf{A}(\boldsymbol{\Lambda})^{-1}\|_{\text{op}}^{2} \|\widetilde{\boldsymbol{\Lambda}}\|_{\text{op}}$$
$$\le \|\boldsymbol{\Lambda}_{0}^{-1}\|_{\text{op}}^{2}$$

where the last inequality follows since  $\mathbf{A}(\mathbf{\Lambda}) \succeq \mathbf{\Lambda}_0$  for all  $\mathbf{\Lambda}$ . This also suffices as a bound on M.

To bound the smoothness, again by the Mean Value Theorem it suffices to bound the operator norm of the Hessian. Standard calculus gives that, using  $\nabla^2 f(\mathbf{\Lambda})[\widetilde{\mathbf{\Lambda}}, \overline{\mathbf{\Lambda}}]$  to denote the Hessian of f in direction  $(\widetilde{\mathbf{\Lambda}}, \overline{\mathbf{\Lambda}})$ :

$$\begin{split} \nabla^2 f(\mathbf{\Lambda}) [\widetilde{\mathbf{\Lambda}}, \overline{\mathbf{\Lambda}}] &= -\frac{d}{dt} \left( \sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda} + t\bar{\mathbf{\Lambda}})^{-1}}^2} \right)^{-1} \cdot \sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda} + t\bar{\mathbf{\Lambda}})^{-1}}^2 \boldsymbol{\phi}^{\top} \mathbf{A} (\mathbf{\Lambda} + t\bar{\mathbf{\Lambda}})^{-1} \widetilde{\mathbf{\Lambda}} \mathbf{A} (\mathbf{\Lambda} + t\bar{\mathbf{\Lambda}})^{-1} \boldsymbol{\phi} \\ &= -\eta \left( \sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2} \right)^{-2} \left( \sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2 \boldsymbol{\phi}^{\top} \mathbf{A} (\mathbf{\Lambda})^{-1} \bar{\mathbf{\Lambda}} \mathbf{A} (\mathbf{\Lambda})^{-1} \boldsymbol{\phi} \right) \left( \sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2} \right)^{-1} \sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2} \left( \boldsymbol{\phi}^{\top} \mathbf{A} (\mathbf{\Lambda})^{-1} \bar{\mathbf{\Lambda}} \mathbf{A} (\mathbf{\Lambda})^{-1} \boldsymbol{\phi} \right) \left( \boldsymbol{\phi}^{\top} \mathbf{A} (\mathbf{\Lambda})^{-1} \widetilde{\mathbf{\Lambda}} \mathbf{A} (\mathbf{\Lambda})^{-1} \boldsymbol{\phi} \right) \\ &+ \left( \sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2} \right)^{-1} \sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2} \boldsymbol{\phi}^{\top} \mathbf{A} (\mathbf{\Lambda})^{-1} \bar{\mathbf{\Lambda}} \mathbf{A} (\mathbf{\Lambda})^{-1} \tilde{\mathbf{\Lambda}} \mathbf{A} (\mathbf{\Lambda})^{-1} \boldsymbol{\phi} \right) \end{aligned}$$

$$+ \left( \sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\boldsymbol{\Lambda})^{-1}}^2} \right)^{-1} \sum_{\boldsymbol{\phi} \in \Phi} e^{\eta \|\boldsymbol{\phi}\|_{\mathbf{A}(\boldsymbol{\Lambda})^{-1}}^2} \boldsymbol{\phi}^{\top} \mathbf{A}(\boldsymbol{\Lambda})^{-1} \widetilde{\boldsymbol{\Lambda}} \mathbf{A}(\boldsymbol{\Lambda})^{-1} \bar{\boldsymbol{\Lambda}} \mathbf{A}(\boldsymbol{\Lambda})^{-1} \boldsymbol{\phi}.$$

We can bound this as

$$\sup_{\boldsymbol{\Lambda},\widetilde{\boldsymbol{\Lambda}},\bar{\boldsymbol{\Lambda}}\succeq 0, \|\boldsymbol{\Lambda}\|_{\mathrm{op}}\leq 1, \|\widetilde{\boldsymbol{\Lambda}}\|_{\mathrm{op}}\leq 1, \|\bar{\boldsymbol{\Lambda}}\|_{\mathrm{op}}\leq 1} |\nabla^2 f(\boldsymbol{\Lambda})[\widetilde{\boldsymbol{\Lambda}},\bar{\boldsymbol{\Lambda}}]| \leq 2\eta \|\boldsymbol{\Lambda}_0^{-1}\|_{\mathrm{op}}^4 + 2\|\boldsymbol{\Lambda}_0^{-1}\|_{\mathrm{op}}^3.$$

Convexity of  $f(\mathbf{\Lambda})$  follows since it is the composition of a convex function with a strictly increasing convex function, so it is itself convex.

Lemma D.6. For  $\Lambda$  invertible,  $\frac{d}{dt}(\Lambda + te_i e_i^{\top})^{-1} = -\Lambda^{-1} e_i e_i^{\top} \Lambda^{-1}$ .

*Proof.* We can compute the gradient as

$$\frac{d}{dt}(\mathbf{\Lambda} + t\mathbf{e}_i\mathbf{e}_j^{\top})^{-1} = \lim_{t \to 0} \frac{(\mathbf{\Lambda} + t\mathbf{e}_i\mathbf{e}_j^{\top})^{-1} - \mathbf{\Lambda}^{-1}}{t}.$$

By the Sherman-Morrison formula,

$$(\mathbf{\Lambda} + t\mathbf{e}_i\mathbf{e}_j^{ op})^{-1} = \mathbf{\Lambda}^{-1} - rac{t\mathbf{\Lambda}^{-1}\mathbf{e}_i\mathbf{e}_j^{ op}\mathbf{\Lambda}^{-1}}{1 + t\mathbf{e}_j^{ op}\mathbf{\Lambda}^{-1}\mathbf{e}_i}$$

so as  $t \to 0$ ,

$$(\mathbf{\Lambda} + t\mathbf{e}_i\mathbf{e}_i^{\mathsf{T}})^{-1} 
ightarrow \mathbf{\Lambda}^{-1} - t\mathbf{\Lambda}^{-1}\mathbf{e}_i\mathbf{e}_i^{\mathsf{T}}\mathbf{\Lambda}^{-1}$$

Thus,

$$\lim_{t\to 0}\frac{(\boldsymbol{\Lambda}+t\boldsymbol{e}_i\boldsymbol{e}_j^\top)^{-1}-\boldsymbol{\Lambda}^{-1}}{t}=\lim_{t\to 0}\frac{\boldsymbol{\Lambda}^{-1}-t\boldsymbol{\Lambda}^{-1}\boldsymbol{e}_i\boldsymbol{e}_j^\top\boldsymbol{\Lambda}^{-1}-\boldsymbol{\Lambda}^{-1}}{t}=-\boldsymbol{\Lambda}^{-1}\boldsymbol{e}_i\boldsymbol{e}_j^\top\boldsymbol{\Lambda}^{-1}.$$

#### D.3Obtaining Well-Conditioned Covariates

# Algorithm 7 Collect Well-Conditioned Covariates (CONDITIONEDCOV)

- 1: **input**: Scale N, minimum eigenvalue  $\underline{\lambda}$ , confidence  $\delta$
- 2: **for**  $j = 1, 2, 3, \dots$  **do**
- $T_{j} \leftarrow \operatorname{poly}(2^{j}, d, H, \log 1/\delta)$   $\epsilon_{j} \leftarrow 2^{-j}, \, \gamma_{j}^{2} \leftarrow \frac{2^{-j}}{\max\{12544d \log \frac{2^{N}(2+32T_{j})}{\delta}, \underline{\lambda}\}}, \, \delta_{j} \leftarrow \delta/(4j^{2})$
- Run Algorithm 5 of Wagenmaker et al. (2022) with parameters  $(\epsilon_j, \gamma_j^2, \delta_j)$ , obtain covariates  $\dot{\Sigma}$  and store policies run as  $\dot{\Pi}$
- if  $\lambda_{\min}(\widetilde{\Lambda}) \geq \max\{12544d \log \frac{2N(2+32T_j)}{\delta}, \underline{\lambda}\}$  then
- 8: Rerun every policy  $\pi \in \widetilde{\Pi} [N/|\widetilde{\Pi}|]$  times, collect covariates  $\bar{\Sigma}$
- 9: return  $\Sigma + \bar{\Sigma}$

**Lemma D.7.** Consider running policies  $(\pi_{\tau})_{\tau=1}^{T}$ , for  $\pi_{\tau}$   $\mathcal{F}_{\tau-1}$ -measurable, and collecting covariance  $\Sigma_{T} = \sum_{\tau=1}^{T} \phi_{\tau} \phi_{\tau}^{\top}$ . Then as long as

$$\lambda_{\min}(\mathbf{\Sigma}_T) \ge 12544d \log \frac{2 + 32T}{\delta}.$$

with probability at least  $1-\delta$ , if we rerun each  $(\pi_{\tau})_{\tau=1}^T$ , we will collect covariates  $\widetilde{\Sigma}_T$  such that

$$\lambda_{\min}(\widetilde{oldsymbol{\Sigma}}_T) \geq rac{1}{2}\lambda_{\min}(oldsymbol{\Sigma}_T).$$

*Proof.* Let  $\mathcal{N}$  be an  $\frac{1}{8T}$ -net of  $\mathcal{S}^{d-1}$ . Let  $\Sigma \succeq 0$  be any matrix with  $\|\Sigma\|_{\text{op}} \leq T$  and let  $\boldsymbol{v}$  be the minimum eigenvalue of  $\Sigma$ . Let  $\widetilde{\boldsymbol{v}} \in \mathcal{N}$  be the element of  $\mathcal{N}$  closest to  $\boldsymbol{v}$  in the  $\ell_2$  norm. Then:

By the construction of  $\mathcal{N}$  and since  $\|\mathbf{\Sigma}\|_{\text{op}} \leq T$ , we can bound  $2\|\mathbf{\Sigma}\|_{\text{op}}\|\tilde{\mathbf{v}} - \mathbf{v}\|_2 \leq 1/4$ , so

$$\widetilde{\boldsymbol{v}}^{\top} \boldsymbol{\Sigma} \widetilde{\boldsymbol{v}} - 2 \| \boldsymbol{\Sigma} \|_{\text{op}} \| \widetilde{\boldsymbol{v}} - \boldsymbol{v} \|_2 \ge \widetilde{\boldsymbol{v}}^{\top} \boldsymbol{\Sigma} \widetilde{\boldsymbol{v}} - 1/4$$

which implies

$$\lambda_{\min}(\mathbf{\Sigma}) + 1/4 \ge \widetilde{\mathbf{v}}^{\top} \mathbf{\Sigma} \widetilde{\mathbf{v}} \ge \min_{\widetilde{\mathbf{v}} \in \mathcal{N}} \widetilde{\mathbf{v}}^{\top} \mathbf{\Sigma} \widetilde{\mathbf{v}}. \tag{D.2}$$

By Lemma A.1, we can bound  $|\mathcal{N}| \leq (1+16T)^d$ .

Note that  $\operatorname{Var}[\boldsymbol{v}^{\top}\boldsymbol{\phi}_{\tau}|\mathcal{F}_{\tau-1}] \leq \mathbb{E}_{\pi_{\tau}}[(\boldsymbol{v}^{\top}\boldsymbol{\phi}_{\tau})^2]$  so  $\sum_{\tau=1}^{T}\operatorname{Var}[\boldsymbol{v}^{\top}\boldsymbol{\phi}_{\tau}|\mathcal{F}_{\tau-1}] \leq \boldsymbol{v}^{\top}\mathbb{E}[\boldsymbol{\Sigma}_{T}|\pi_{1},\ldots,\pi_{T}]\boldsymbol{v}$  for  $\mathbb{E}[\boldsymbol{\Sigma}_{T}|\pi_{1},\ldots,\pi_{T}] = \sum_{\tau=1}^{T}\mathbb{E}_{\pi_{\tau}}[\boldsymbol{\phi}_{\tau}\boldsymbol{\phi}_{\tau}^{\top}]$ . By Freedman's Inequality (Lemma A.5), for all  $\boldsymbol{v} \in \mathcal{N}$  simultaneously, we will have, with probability at least  $1-\delta$ ,

$$\left| \boldsymbol{v}^{\top} \boldsymbol{\Sigma}_{T} \boldsymbol{v} - \boldsymbol{v}^{\top} \mathbb{E}[\boldsymbol{\Sigma}_{T} | \pi_{1}, \dots, \pi_{T}] \boldsymbol{v} \right| \leq 2 \sqrt{\boldsymbol{v}^{\top} \mathbb{E}[\boldsymbol{\Sigma}_{T} | \pi_{1}, \dots, \pi_{T}] \boldsymbol{v} \log \frac{2|\mathcal{N}|}{\delta}} + \log \frac{2|\mathcal{N}|}{\delta}$$
(D.3)

$$\left| \boldsymbol{v}^{\top} \widetilde{\boldsymbol{\Sigma}}_{T} \boldsymbol{v} - \boldsymbol{v}^{\top} \mathbb{E}[\boldsymbol{\Sigma}_{T} | \pi_{1}, \dots, \pi_{T}] \boldsymbol{v} \right| \leq 2 \sqrt{\boldsymbol{v}^{\top} \mathbb{E}[\boldsymbol{\Sigma}_{T} | \pi_{1}, \dots, \pi_{T}] \boldsymbol{v} \log \frac{2|\mathcal{N}|}{\delta}} + \log \frac{2|\mathcal{N}|}{\delta}.$$
 (D.4)

Rearranging (D.3), some algebra shows that

$$\mathbf{v}^{\top} \mathbb{E}[\mathbf{\Sigma}_{T} | \pi_{1}, \dots, \pi_{T}] \mathbf{v} \leq \mathbf{v}^{\top} \mathbf{\Sigma}_{T} \mathbf{v} + 3 \log \frac{2|\mathcal{N}|}{\delta} + 2 \sqrt{\mathbf{v}^{\top} \mathbf{\Sigma}_{T} \mathbf{v} \log \frac{2|\mathcal{N}|}{\delta}} + 2 \log^{2} \frac{2|\mathcal{N}|}{\delta}$$

$$\leq \mathbf{v}^{\top} \mathbf{\Sigma}_{T} \mathbf{v} + 6 \log \frac{2|\mathcal{N}|}{\delta} + 2 \sqrt{\mathbf{v}^{\top} \mathbf{\Sigma}_{T} \mathbf{v} \log \frac{2|\mathcal{N}|}{\delta}}$$

$$\leq 3 \mathbf{v}^{\top} \mathbf{\Sigma}_{T} \mathbf{v} + 8 \log \frac{2|\mathcal{N}|}{\delta}$$

where the last inequality uses  $\sqrt{ab} \leq \max\{a,b\}$ . Thus, if (D.3) and (D.4) hold, we have

$$v^{\top} \widetilde{\boldsymbol{\Sigma}}_{T} \boldsymbol{v} \geq \boldsymbol{v}^{\top} \boldsymbol{\Sigma}_{T} \boldsymbol{v} - 4 \sqrt{\boldsymbol{v}^{\top} \mathbb{E}[\boldsymbol{\Sigma}_{T} | \pi_{1}, \dots, \pi_{T}] \boldsymbol{v} \log \frac{2|\mathcal{N}|}{\delta}} - 2 \log \frac{2|\mathcal{N}|}{\delta}$$
$$\geq \boldsymbol{v}^{\top} \boldsymbol{\Sigma}_{T} \boldsymbol{v} - 4 \sqrt{3 \boldsymbol{v}^{\top} \boldsymbol{\Sigma}_{T} \boldsymbol{v} \log \frac{2|\mathcal{N}|}{\delta}} - 14 \log \frac{2|\mathcal{N}|}{\delta}$$

Therefore, as long as

$$\boldsymbol{v}^{\top} \boldsymbol{\Sigma}_T \boldsymbol{v} \geq 12544 \log \frac{2|\mathcal{N}|}{\delta},$$

we can lower bound

$$\boldsymbol{v}^{\top} \boldsymbol{\Sigma}_{T} \boldsymbol{v} - 4 \sqrt{3 \boldsymbol{v}^{\top} \boldsymbol{\Sigma}_{T} \boldsymbol{v} \log \frac{2|\mathcal{N}|}{\delta}} - 14 \log \frac{2|\mathcal{N}|}{\delta} \ge \frac{3}{4} \boldsymbol{v}^{\top} \boldsymbol{\Sigma}_{T} \boldsymbol{v} \ge \frac{3}{4} \lambda_{\min}(\boldsymbol{\Sigma}_{T})$$

so, for all  $v \in \mathcal{N}$ ,

$$oldsymbol{v}^{ op} \widetilde{oldsymbol{\Sigma}}_T oldsymbol{v} \geq rac{3}{4} \lambda_{\min}(oldsymbol{\Sigma}_T).$$

By assumption,  $\lambda_{\min}(\mathbf{\Sigma}_T) \geq 12544d \log \frac{2+32T}{\delta}$ , which implies, since  $|\mathcal{N}| \leq (1+16T)^d$ , that for all  $\mathbf{v} \in \mathcal{S}^{d-1}$ ,  $\mathbf{v}^{\top} \mathbf{\Sigma}_T \mathbf{v} \geq 12544 \log \frac{2|\mathcal{N}|}{\delta}$ , so the above condition will be met.

Since  $\|\widetilde{\Sigma}_T\|_{\text{op}} \leq T$ , we can apply (D.2) to then get that

$$\lambda_{\min}(\widetilde{\boldsymbol{\Sigma}}_T) \geq \frac{3}{4}\lambda_{\min}(\boldsymbol{\Sigma}_T) - 1/4 \geq \frac{1}{2}\lambda_{\min}(\boldsymbol{\Sigma}_T) + \frac{1}{4}(\lambda_{\min}(\boldsymbol{\Sigma}_T) - 1).$$

Since we have already establishes that  $\lambda_{\min}(\Sigma_T) \geq 12544d \log \frac{2+32T}{\delta}$ , we have  $\lambda_{\min}(\Sigma_T) - 1 \geq 0$ , so we can lower bound

$$\lambda_{\min}(\widetilde{oldsymbol{\Sigma}}_T) \geq rac{1}{2}\lambda_{\min}(oldsymbol{\Sigma}_T).$$

**Lemma D.8.** With probability at least  $1 - \delta$ , Algorithm 7 will terminate after at most

$$N + \operatorname{poly} \log \left( \frac{1}{\sup_{\pi} \lambda_{\min}(\boldsymbol{\Sigma}_{\pi})}, d, H, \underline{\lambda}, \log \frac{N}{\delta} \right) \cdot \left( \frac{d \max\{d \log \frac{N}{\delta}, \underline{\lambda}\}}{\sup_{\pi} \lambda_{\min}(\boldsymbol{\Sigma}_{\pi})^{2}} + \frac{d^{4}H^{3} \log^{7/2} \frac{1}{\delta}}{\sup_{\pi} \lambda_{\min}(\boldsymbol{\Sigma}_{\pi})} \right)$$

episodes, and will return covariates  $\Sigma$  such that

$$\lambda_{\min}(\mathbf{\Sigma}) \ge N \cdot \min\left\{\frac{\sup_{\pi} \lambda_{\min}(\mathbf{\Sigma}_{\pi})^{2}}{d}, \frac{\sup_{\pi} \lambda_{\min}(\mathbf{\Sigma}_{\pi})}{d^{3}H^{3}\log^{7/2}1/\delta}\right\} \cdot \operatorname{poly} \log\left(\frac{1}{\sup_{\pi} \lambda_{\min}(\mathbf{\Sigma}_{\pi})}, d, H, \underline{\lambda}, \log\frac{N}{\delta}\right)^{-1} + \max\{d\log 1/\delta, \underline{\lambda}\}$$

and

$$\|\mathbf{\Sigma}\|_{\text{op}} \leq N + \text{poly}\log\left(\frac{1}{\sup_{\pi} \lambda_{\min}(\mathbf{\Sigma}_{\pi})}, d, H, \underline{\lambda}, \log\frac{N}{\delta}\right) \cdot \left(\frac{d \max\{d \log\frac{N}{\delta}, \underline{\lambda}\}}{\sup_{\pi} \lambda_{\min}(\mathbf{\Sigma}_{\pi})^2} + \frac{d^4H^3 \log^{7/2}\frac{1}{\delta}}{\sup_{\pi} \lambda_{\min}(\mathbf{\Sigma}_{\pi})}\right).$$

*Proof.* By Theorem 4 of Wagenmaker et al. (2022), as long as Algorithm 5 of Wagenmaker et al. (2022) is run with parameters  $\epsilon$  and  $\gamma^2$ , it will terminate after at most

$$c_1 \cdot \frac{1}{\epsilon} \max \left\{ \frac{dm}{\gamma^2} \log \frac{dm}{\epsilon \gamma^2}, d^4 H^3 m^{7/2} \log^{3/2} (d/\gamma^2) \log^{7/2} \frac{c_2 m d H \log(d/\gamma^2)}{\delta} \right\}$$

episodes for  $m = \lceil \log(2/\epsilon) \rceil$  (to get the slightly more precise bound on the number of episodes collected than that given in Theorem 4 of Wagenmaker et al. (2022), we use the precise definition of

 $K_i$  given at the start of Appendix B). Furthermore, if  $\epsilon \leq \sup_{\pi} \lambda_{\min}(\Sigma_{\pi})$ , with probability at least  $1 - \delta$  it will collect covariates  $\widetilde{\Sigma}$  satisfying  $\lambda_{\min}(\widetilde{\Sigma}) \geq \epsilon/\gamma^2$ .

It follows that, by our choice of  $\epsilon_j = 2^{-j}$ ,  $\gamma_j^2 = \frac{2^{-j}}{\max\{12544d\log\frac{2^N(2+32T_j)}{\delta},\underline{\lambda}\}}$ , and  $\delta_j = \delta/(4j^2)$ , for every j we will collect at most

$$c_1 \cdot 2^{j} \max \left\{ 2^{j} dj^2 \max\{d \log \frac{2N(2+32T_j)}{\delta}, \underline{\lambda}\} \log(dja_j), d^4H^3j^5 \log^{3/2}(da_j) \log^{7/2} \frac{c_2 j^4 dH \log(da_j)}{\delta} \right\}$$

episodes, where we denote  $a_j := \max\{12544d\log\frac{2N(2+32T_j)}{\delta},\underline{\lambda}\}$ . Note that  $T_j$  is an upper bound on this complexity. Furthermore, once j is large enough that  $2^{-j} \leq \sup_{\pi} \lambda_{\min}(\Sigma_{\pi})$ , Theorem 4 of Wagenmaker et al. (2022) implies that the condition  $\lambda_{\min}(\widetilde{\Sigma}) \geq \epsilon_j/\gamma_j^2$  will be met. By our choice of  $\gamma_j^2$  and  $\epsilon_j$ , it follows that the if condition on Line 6 will be met once  $2^{-j} \leq \sup_{\pi} \lambda_{\min}(\Sigma_{\pi})$ . Since  $2^{-j}$  decreases by a factor of 2 each time, it follows that the if statement on Line 6 will have terminated once  $2^{-j} \geq \sup_{\pi} \lambda_{\min}(\Sigma_{\pi})/2$ . This implies that the total number of episodes collected before the if statement on Line 6 is met is bounded as

$$\operatorname{poly} \log \left( \frac{1}{\sup_{\pi} \lambda_{\min}(\mathbf{\Sigma}_{\pi})}, d, H, \underline{\lambda}, \log \frac{N}{\delta} \right) \cdot \left( \frac{d \max\{d \log \frac{N}{\delta}, \underline{\lambda}\}}{\sup_{\pi} \lambda_{\min}(\mathbf{\Sigma}_{\pi})^{2}} + \frac{d^{4}H^{3} \log^{7/2} \frac{1}{\delta}}{\sup_{\pi} \lambda_{\min}(\mathbf{\Sigma}_{\pi})} \right)$$
(D.5)

By Lemma D.7, since  $\lambda_{\min}(\widetilde{\Sigma}) \geq \max\{12544d \log \frac{2N(2+32T_j)}{\delta}, \underline{\lambda}\}$  and  $T_j$  is an upper bound on the number of episodes run at epoch j, every time we run all policies  $\pi \in \widetilde{\Pi}$ , with probability at least  $1 - \delta/(2N)$ , we will collect covariates  $\Sigma$  such that

$$\lambda_{\min}(\mathbf{\Sigma}) \geq \lambda_{\min}(\widetilde{\mathbf{\Sigma}})/2 \geq \frac{1}{2} \max\{12544d \log \frac{2N(2+32T_j)}{\delta}, \underline{\lambda}\}.$$

Thus, if we rerun every policy  $\lceil N/|\widetilde{\Pi}| \rceil$  times to create covariates  $\bar{\Sigma}$ , with probability at least  $1 - \delta/2$ , we have

$$\lambda_{\min}(\bar{\Sigma}) \ge \frac{N}{2|\widetilde{\Pi}|} \max\{12544d \log \frac{2N(2+32T_j)}{\delta}, \underline{\lambda}\}.$$

Note that this procedure will complete after at most  $N + |\widetilde{\Pi}|$  episodes. Furthermore,  $|\widetilde{\Pi}| \leq (D.5)$ , so we can lower bound

$$\lambda_{\min}(\bar{\boldsymbol{\Sigma}}) \geq N \cdot \min\left\{\frac{\sup_{\pi} \lambda_{\min}(\boldsymbol{\Sigma}_{\pi})^{2}}{d}, \frac{\sup_{\pi} \lambda_{\min}(\boldsymbol{\Sigma}_{\pi})}{d^{3}H^{3}\log^{7/2}1/\delta}\right\} \cdot \operatorname{poly}\log\left(\frac{1}{\sup_{\pi} \lambda_{\min}(\boldsymbol{\Sigma}_{\pi})}, d, H, \underline{\lambda}, \log\frac{N}{\delta}\right)^{-1}.$$

The final lower bound on the returned covariates follows since we return  $\bar{\Sigma} + \tilde{\Sigma}$ , and we know that  $\lambda_{\min}(\tilde{\Sigma}) \geq \max\{12544d \log \frac{2N(2+32T_j)}{\delta}, \underline{\lambda}\}$ . The upper bound on  $\|\bar{\Sigma} + \tilde{\Sigma}\|_{\text{op}}$  follows since every feature vector encountered has norm of at most 1.

The failure probability of each call to Algorithm 5 of Wagenmaker et al. (2022) is  $\delta/(4j^2)$ , so the total failure probability of Algorithm 7 is

$$\sum_{j=1}^{\infty} \frac{\delta}{4j^2} = \frac{\pi^2}{24} \delta \le \delta/2.$$

# D.4 Online XY-Optimal Design

**Theorem 9** (Full version of Theorem 5). Consider running OPTCOV with some  $\epsilon > 0$  and functions

$$f_i(\mathbf{\Lambda}) \leftarrow \widetilde{\mathsf{XY}}_{\mathrm{opt}}(\mathbf{\Lambda})$$

for  $\Lambda_0 \leftarrow (T_i K_i)^{-1} \Sigma_i =: \Lambda_i$  and

$$\eta_{i} = \frac{2}{\gamma_{\Phi}} \cdot (1 + \|\mathbf{\Lambda}_{i}\|_{\text{op}}) \cdot \log |\Phi| 
L_{i} = \|\mathbf{\Lambda}_{i}^{-1}\|_{\text{op}}^{2}, \quad \beta_{i} = 2\|\mathbf{\Lambda}_{i}^{-1}\|_{\text{op}}^{3} (1 + \eta_{i}\|\mathbf{\Lambda}_{i}^{-1}\|_{\text{op}}), \quad M_{i} = \|\mathbf{\Lambda}_{i}^{-1}\|_{\text{op}}^{2}$$

where  $\Sigma_i$  is the matrix returned by running ConditionedCov with  $N \leftarrow T_i K_i$ ,  $\delta \leftarrow \delta/(2i^2)$ , and some  $\underline{\lambda} \geq 0$ . Then with probability  $1 - 2\delta$ , this procedure will collect at most

$$20 \cdot \frac{\inf_{\mathbf{\Lambda} \in \mathbf{\Omega}} \max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{\mathbf{A}(\mathbf{\Lambda})^{-1}}^2}{\epsilon_{\exp}} + \operatorname{poly}\left(d, H, \log 1/\delta, \frac{1}{\lambda_{\min}^{\star}}, \frac{1}{\gamma_{\Phi}}, \underline{\lambda}, \log |\Phi|, \log \frac{1}{\epsilon_{\exp}}\right)$$

episodes, where

$$\mathbf{A}(\mathbf{\Lambda}) = \mathbf{\Lambda} + \min\left\{\frac{(\lambda_{\min}^{\star})^2}{d}, \frac{\lambda_{\min}^{\star}}{d^3 H^3 \log^{7/2} 1/\delta}\right\} \cdot \operatorname{poly} \log\left(\frac{1}{\lambda_{\min}^{\star}}, d, H, \underline{\lambda}, \log \frac{1}{\delta}\right)^{-1} \cdot I,$$

and will produce covariates  $\widehat{\Sigma} + \Sigma_i$  such that

$$\max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{(\widehat{\boldsymbol{\Sigma}} + \boldsymbol{\Sigma}_i)^{-1}}^2 \le \epsilon_{\exp}$$

and

$$\lambda_{\min}(\widehat{\Sigma} + \Sigma_i) \ge \max\{d \log 1/\delta, \underline{\lambda}\}.$$

*Proof.* Note that the total failure probability of our calls to CONDITIONEDCOV is at most

$$\sum_{i=1}^{\infty} \frac{\delta}{2i^2} = \frac{\pi^2}{12} \delta \le \delta.$$

For the remainder of the proof, we will then assume that we are on the success event of CONDITIONEDCOV, as defined in Lemma D.8.

By Lemma D.5,  $f_i(\Lambda)$  satisfies Definition 5.1 with constants

$$L_i = \|\mathbf{\Lambda}_i^{-1}\|_{\text{op}}^2, \quad \beta_i = 2\|\mathbf{\Lambda}_i^{-1}\|_{\text{op}}^3 (1 + \eta_i \|\mathbf{\Lambda}_i^{-1}\|_{\text{op}}), \quad M_i = \|\mathbf{\Lambda}_i^{-1}\|_{\text{op}}^2$$

for  $\Lambda_i \leftarrow (T_i K_i)^{-1} \Sigma_i$ .

By Lemma D.8, on the success event of Lemma D.8 we have that

$$\lambda_{\min}(\mathbf{\Lambda}_i) \ge \min\left\{\frac{(\lambda_{\min}^{\star})^2}{d}, \frac{\lambda_{\min}^{\star}}{d^3 H^3 \log^{7/2} 1/\delta}\right\} \cdot \operatorname{poly} \log\left(\frac{1}{\lambda_{\min}^{\star}}, d, H, \underline{\lambda}, i, \log \frac{1}{\delta}\right)^{-1}$$

(note that the poly  $\log(i)^{-1}$  dependence arises because we take  $N \leftarrow T_i K_i = 2^{4i}$ ). Thus, we can bound, for all i (using the upper bound on  $\|\mathbf{\Sigma}_i\|_{\text{op}}$  given in Lemma D.8 to upper bound  $\eta_i$ ),

$$L_i = M_i \le \max \left\{ \frac{d^2}{(\lambda_{\min}^{\star})^4}, \frac{d^6 H^6 \log^7 1/\delta}{(\lambda_{\min}^{\star})^2} \right\} \cdot \operatorname{poly} \log \left( \frac{1}{\lambda_{\min}^{\star}}, d, H, \underline{\lambda}, i, \log \frac{1}{\delta} \right),$$

$$\beta_i \leq \text{poly}\left(d, H, \log 1/\delta, \frac{1}{\lambda_{\min}^{\star}}, \frac{1}{\gamma_{\Phi}}, \underline{\lambda}, i, \log |\Phi|\right).$$

Assume that the termination condition of OPTCoV for  $\hat{i}$  satisfying

$$\hat{i} \le \log \left( \operatorname{poly} \left( \frac{1}{\epsilon_{\exp}}, d, H, \log 1/\delta, \frac{1}{\lambda_{\min}^{\star}}, \frac{1}{\gamma_{\Phi}}, \underline{\lambda}, \log |\Phi| \right) \right). \tag{D.6}$$

We assume this holds and justify it at the conclusion of the proof. For notational convenience, define

$$\iota := \operatorname{poly} \left( \log \frac{1}{\epsilon_{\exp}}, d, H, \log 1/\delta, \frac{1}{\lambda_{\min}^{\star}}, \frac{1}{\gamma_{\Phi}}, \underline{\lambda}, \log |\Phi| \right).$$

Given this upper bound on  $\hat{i}$ , set

$$L = M := \max \left\{ \frac{d^2}{(\lambda_{\min}^{\star})^4}, \frac{d^6 H^6 \log^7 1/\delta}{(\lambda_{\min}^{\star})^2} \right\} \cdot \text{poly} \log \iota, \qquad \beta := \iota.$$

With this choice of  $L, M, \beta$ , we have  $L_i \leq L, M_i \leq M, \beta_i \leq \beta$  for all  $i \leq \hat{i}$ .

Now take  $f(\mathbf{\Lambda}) \leftarrow \mathsf{XY}_{\mathrm{opt}}(\mathbf{\Lambda}; \eta, \mathbf{\Lambda}_0)$  with

$$\mathbf{\Lambda}_0 \leftarrow \min \left\{ \frac{(\lambda_{\min}^{\star})^2}{d}, \frac{\lambda_{\min}^{\star}}{d^3 H^3 \log^{7/2} 1/\delta} \right\} \cdot \frac{1}{\text{poly} \log \iota} \cdot I$$
 (D.7)

and

$$\eta = \frac{2\log|\Phi|}{\gamma_{\Phi}} \cdot \left(1 + \min\left\{\frac{(\lambda_{\min}^{\star})^2}{d}, \frac{\lambda_{\min}^{\star}}{d^3 H^3 \log^{7/2} 1/\delta}\right\} \cdot \frac{1}{\operatorname{poly} \log \iota}\right).$$

Note that in this case, we have  $\|\mathbf{\Lambda}_0\|_{\mathrm{op}} \leq \lambda_{\min}(\mathbf{\Lambda}_i)$  for all i, so  $\mathbf{\Lambda}_0 \leq \mathbf{\Lambda}_i$  and  $\eta \leq \eta_i$ . By the construction of  $\widetilde{\mathsf{XY}}_{\mathrm{opt}}$  and Lemma D.2, it follows that  $f(\mathbf{\Lambda}) \geq f_i(\mathbf{\Lambda})$  for all  $\mathbf{\Lambda} \succeq 0$ , so this is a valid choice of f, as required by Theorem 8. Furthermore, we can set R = 2, since  $\|\mathbf{\Lambda}_{\pi}\|_{\mathrm{F}} \leq 1$  for all  $\pi$ .

To apply Theorem 8, it remains only to find a suitable value of  $f_{\min}$ . By Lemma D.1 and Lemma D.3, we can lower bound  $f_i$  by  $\frac{\gamma_{\Phi}}{1+||\mathbf{A}_i||_{\mathrm{op}}}$ . By Lemma D.8, we can lower bound

$$\frac{\gamma_{\Phi}}{1 + \|\mathbf{\Lambda}_i\|_{\mathrm{op}}} \ge \frac{\gamma_{\Phi}}{2 + \mathrm{poly}\log\iota \cdot \left(\frac{d\max\{d\log\frac{1}{\delta}, \underline{\lambda}\}}{(\lambda_{\min}^{\star})^2} + \frac{d^4H^3\log^{7/2}\frac{1}{\delta}}{\lambda_{\min}^{\star}}\right)}.$$

We then take this as our choice of  $f_{\min}$ .

We can now apply Theorem 8, using the complexity for OPTCOV instantiated with FORCE given in Corollary 6, and get that with probability at least  $1 - \delta$ , OPTCOV will terminate in

$$N \le 5N^{\star} \left(\epsilon_{\rm exp}/2; f\right) + \iota$$

episodes, and will return (time-normalized) covariates  $\widehat{\mathbf{\Lambda}}$  such that

$$f_{\widehat{i}}(\widehat{\mathbf{\Lambda}}) \leq N\epsilon_{\text{exp}}.$$

By Lemma D.4, our choice of  $\eta$  and  $\Lambda_0$ , we can upper bound

$$N^{\star}\left(\epsilon_{\mathrm{exp}}/2;f\right) \leq 2N^{\star}\left(\epsilon_{\mathrm{exp}}/2;\mathsf{XY}_{\mathrm{opt}}\right) = \frac{4\inf_{\boldsymbol{\Lambda} \in \boldsymbol{\Omega}} \max_{\boldsymbol{\phi} \in \boldsymbol{\Phi}} \|\boldsymbol{\phi}\|_{\mathbf{A}(\boldsymbol{\Lambda})^{-1}}^2}{\epsilon_{\mathrm{exp}}}$$

where here  $\mathbf{A}(\mathbf{\Lambda}) = \mathbf{\Lambda} + \mathbf{\Lambda}_0$  for  $\mathbf{\Lambda}_0$  as in (D.7). Furthermore, by Lemma D.1 we have

$$\max_{\boldsymbol{\phi} \in \Phi} \|\boldsymbol{\phi}\|_{(\widehat{\boldsymbol{\Sigma}} + \boldsymbol{\Lambda}_0)^{-1}}^2 \leq f_{\widehat{i}}(\widehat{\boldsymbol{\Lambda}}).$$

The final upper bound on the number of episodes collected and the lower bound on the minimum eigenvalue of the covariates follows from Lemma D.8.

It remains to justify our bound on  $\hat{i}$ , (D.6). Note that by definition of OPTCOV, if we run for a total of  $\bar{N}$  episodes, we can bound  $\hat{i} \leq \frac{1}{4} \log_2(\bar{N})$ . However, we see that the bound on  $\hat{i}$  given in (D.6) upper bounds  $\frac{1}{4} \log_2(\bar{N})$  for  $\bar{N}$  the upper bound on the number of samples collected by OPTCOV stated above. Thus, our bound on  $\hat{i}$  is valid.

# E Suboptimality of Optimistic Algorithms

## E.1 Linear Bandit Construction

In the linear bandit setting, at each time step t, the learner chooses some  $z_t \in \mathcal{Z}$ , and observes  $y_t$ . We will consider the case when the noise is Bernoulli so that  $y_t \sim \text{Bernoulli}(\langle \theta_{\star}, z_t \rangle + 1/2)$ , and will set

$$\boldsymbol{\theta}_{\star} = \boldsymbol{e}_1, \quad \mathcal{Z} = \{\xi \boldsymbol{e}_1, \boldsymbol{e}_2, \dots, \boldsymbol{e}_d, \boldsymbol{x}_2, \dots, \boldsymbol{x}_d\}, \quad \boldsymbol{x}_i = (\xi - \Delta)\boldsymbol{e}_1 + \gamma \boldsymbol{e}_i$$

for some  $\xi, \Delta, \alpha$  to be chosen. In this setting, the optimal arm is  $z^* = \xi e_1$ , and  $\Delta(e_i) = \xi, i \geq 2$ ,  $\Delta(x_i) = \Delta$ .

We will assume:

$$\frac{1}{52d} \ge \xi \ge \max\{\gamma/\sqrt{d}, \sqrt{\Delta}\}, \quad \max\left\{\zeta := \frac{2C_1}{(d/\Delta^2)^{1-\alpha}} + \frac{2C_2\Delta^2}{d}, \Delta\right\} \le \gamma^2.$$
 (E.1)

We provide explicit values for  $\xi$ ,  $\Delta$ , and  $\gamma$  that satisfy this in Lemma E.3.

**Definition E.1** ( $\delta$ -correct). We say a stopping rule  $\tau$  is  $\delta$ -correct if  $\mathbb{P}[\hat{z}_{\tau} = z^{\star}] \geq 1 - \delta$ , where  $\hat{z}_{\tau}$  is the arm recommended at time  $\tau$ .

**Lemma E.1.** Consider running some low-regret algorithm satisfying Definition 4.1 on the linear bandit instance described above and let  $\tau$  be some stopping time. Then if  $\tau$  is  $\delta$ -correct, we must have that

$$\mathbb{E}[\tau] \ge \frac{d-1}{48\Delta^2} \cdot \log \frac{1}{2.4\delta}.$$

*Proof.* This proof follows closely the proof of Theorem 1 of Fiez et al. (2019) and relies on the Transportation Lemma of Kaufmann et al. (2016).

Bounding the number of pulls to  $\{e_2, \ldots, e_d\}$ . By assumption, we collect data with a low-regret algorithm satisfying Definition 4.1. Every time we pull  $e_i, i \geq 2$ , we incur a loss of 1/2. Thus, we can lower bound

$$\mathbb{E}[V_0^{\star} - V_0^{\pi_k}] \geq \frac{1}{2} \sum_{i=2}^d \mathbb{E}[\mathbb{P}_{\pi_k}[\boldsymbol{z}_k = \boldsymbol{x}_i]]$$

so, letting  $T(x_i)$  denote the total number of pulls to  $x_i$ , we have

$$C_1 K^{\alpha} + C_2 \ge \sum_{k=1}^{K} \mathbb{E}[V_0^{\star} - V_0^{\pi_k}] \ge \frac{1}{2} \sum_{k=1}^{K} \sum_{i=2}^{d} \mathbb{E}[\mathbb{P}_{\pi_k}[\boldsymbol{z}_k = \boldsymbol{x}_i]] = \frac{1}{2} \sum_{i=2}^{d} \mathbb{E}[T(\boldsymbol{x}_i)]. \tag{E.2}$$

Applying the Transportation Lemma. Let  $\Theta_{\text{alt}}$  denote the set of  $\boldsymbol{\theta}$  vectors such that  $\xi \boldsymbol{e}_1$  is not the optimal arm, that is,  $\max_{\boldsymbol{z} \in \mathcal{X}} \langle \boldsymbol{\theta}, \boldsymbol{z} \rangle > \langle \boldsymbol{\theta}, \xi \boldsymbol{e}_1 \rangle$ . Let  $\nu_{\boldsymbol{\theta}, \boldsymbol{z}} = \text{Bernoulli}(\langle \boldsymbol{\theta}, \boldsymbol{z} \rangle + 1/2)$ . Then by the Transportation Lemma of Kaufmann et al. (2016), for any  $\boldsymbol{\theta} \in \Theta_{\text{alt}}$ , assuming our stopping rule is  $\delta$ -correct, we have

$$\sum_{\boldsymbol{z} \in \mathcal{Z}} \mathbb{E}[T(\boldsymbol{z})] \mathrm{KL}(\nu_{\boldsymbol{\theta}_{\star}, \boldsymbol{z}} || \nu_{\boldsymbol{\theta}, \boldsymbol{z}}) \ge \log \frac{1}{2.4\delta}.$$

Combining this with our constraint (E.2), it follows that  $\sum_{z\in\mathcal{Z}} \mathbb{E}[T(z)] \geq \sum_{z\in\mathcal{Z}} t_z$  for any  $(t_z)_{z\in\mathcal{Z}}$  that is a feasible solution to

$$\min \sum_{z \in \mathcal{Z}} t_z \quad \text{s.t.} \quad \min_{\theta \in \Theta_{\text{alt}}} \sum_{z \in \mathcal{Z}} t_z \text{KL}(\nu_{\theta_{\star},z} || \nu_{\theta,z}) \ge \log \frac{1}{2.4\delta}, C_1(\sum_{z \in \mathcal{Z}} t_z)^{\alpha} + C_2 \ge \frac{1}{2} \sum_{i=2}^{d} t_{x_i}. \quad (E.3)$$

We can rearrange the second constraint to

$$\frac{2C_1}{(\sum_{z\in\mathcal{Z}} t_z)^{1-\alpha}} + \frac{2C_2}{\sum_{z\in\mathcal{Z}} t_z} \ge \frac{\sum_{i=2}^d t_{x_i}}{\sum_{z\in\mathcal{Z}} t_z}.$$

Assume that the optimal value of (E.3) satisfies  $\sum_{z \in \mathcal{Z}} t_z \geq \frac{d}{\Delta^2}$ , then this constraint can be weakened to

$$\zeta := \frac{2C_1}{(d/\Delta^2)^{1-\alpha}} + \frac{2C_2\Delta^2}{d} \ge \frac{\sum_{i=2}^d t_{\boldsymbol{x}_i}}{\sum_{\boldsymbol{z} \in \mathcal{Z}} t_{\boldsymbol{z}}}$$

It follows then that if the optimal value to

$$\min \sum_{z \in \mathcal{Z}} t_z \quad \text{s.t.} \quad \min_{\theta \in \Theta_{\text{alt}}} \sum_{z \in \mathcal{Z}} t_z \text{KL}(\nu_{\theta_*, z} || \nu_{\theta, z}) \ge \log \frac{1}{2.4\delta}, \zeta \ge \frac{\sum_{i=2}^d t_{x_i}}{\sum_{z \in \mathcal{Z}} t_z}$$
(E.4)

is at least  $d/\Delta^2$ , then the optimal value to (E.3) is also at least  $d/\Delta^2$ , so our assumption that  $\sum_{z\in\mathcal{Z}} t_z \geq \frac{d}{\Delta^2}$  will be justified.

For  $z \neq z^*$ , let  $\theta_z(\epsilon, t)$  denote the instance

$$oldsymbol{ heta_{\star}} - rac{(oldsymbol{y_z}^{ op}oldsymbol{ heta_{\star}} + \epsilon)\widetilde{\mathbf{A}}(t)^{-1}oldsymbol{y_z}}{oldsymbol{y_z}^{ op}\widetilde{\mathbf{A}}(t)^{-1}oldsymbol{y_z}}$$

for  $\boldsymbol{y_z} = \boldsymbol{z^*} - \boldsymbol{z}$ ,  $\widetilde{\mathbf{A}}(t) = \sum_{\boldsymbol{z} \in \mathcal{Z}} \frac{t_{\boldsymbol{z}}}{\sum_{\boldsymbol{z}' \in \mathcal{Z}} t_{\boldsymbol{z}'}} \boldsymbol{z} \boldsymbol{z}^{\top} + \operatorname{diag}([\xi^2, \gamma^2/d, \dots, \gamma^2/d])$ , and  $\epsilon \leq \min\{\Delta, \xi\}$ . Note that  $\boldsymbol{y_z}^{\top} \boldsymbol{\theta_z}(\epsilon, t) = -\epsilon < 0$  which implies that  $\boldsymbol{\theta_z}(\epsilon, t) \in \Theta_{\text{alt}}$ . Furthermore, we can bound:

Claim E.2. For all  $z, v \in \mathcal{Z}$ ,

$$\mathrm{KL}(\nu_{\boldsymbol{\theta_{\star}},\boldsymbol{v}}||\nu_{\boldsymbol{\theta_{z}}(\epsilon,t),\boldsymbol{v}}) \leq 16(\boldsymbol{y_{z}^{\top}}\boldsymbol{\theta_{\star}} + \epsilon)^{2} \frac{\boldsymbol{y_{z}^{\top}}\widetilde{\mathbf{A}}(t)^{-1}\boldsymbol{v}\boldsymbol{v}^{\top}\widetilde{\mathbf{A}}(t)^{-1}\boldsymbol{y_{z}}}{(\boldsymbol{y_{z}^{\top}}\widetilde{\mathbf{A}}(t)^{-1}\boldsymbol{y_{z}})^{2}}.$$

This implies that, for any t,

$$\sum_{\boldsymbol{v}\in\mathcal{Z}} t_{\boldsymbol{v}} \mathrm{KL}(\nu_{\boldsymbol{\theta}_{\star},\boldsymbol{v}}||\nu_{\boldsymbol{\theta}_{\boldsymbol{z}}(\epsilon,t),\boldsymbol{v}}) \leq 16 \sum_{\boldsymbol{v}\in\mathcal{Z}} t_{\boldsymbol{v}} (\boldsymbol{y}_{\boldsymbol{z}}^{\top}\boldsymbol{\theta}_{\star} + \epsilon)^{2} \frac{\boldsymbol{y}_{\boldsymbol{z}}^{\top} \tilde{\mathbf{A}}(t)^{-1} \boldsymbol{v} \boldsymbol{v}^{\top} \tilde{\mathbf{A}}(t)^{-1} \boldsymbol{y}_{\boldsymbol{z}}}{(\boldsymbol{y}_{\boldsymbol{z}}^{\top} \tilde{\mathbf{A}}(t)^{-1} \boldsymbol{y}_{\boldsymbol{z}})^{2}}$$

$$= 16 \sum_{\boldsymbol{v} \in \mathcal{Z}} t_{\boldsymbol{v}} \cdot (\boldsymbol{y}_{\boldsymbol{z}}^{\top} \boldsymbol{\theta}_{\star} + \epsilon)^{2} \frac{\boldsymbol{y}_{\boldsymbol{z}}^{\top} \widetilde{\mathbf{A}}(t)^{-1} (\sum_{\boldsymbol{v} \in \mathcal{Z}} \frac{t_{\boldsymbol{v}}}{\sum_{\boldsymbol{v}' \in \mathcal{Z}} t_{\boldsymbol{v}'}} \boldsymbol{v} \boldsymbol{v}^{\top}) \widetilde{\mathbf{A}}(t)^{-1} \boldsymbol{y}_{\boldsymbol{z}}}{(\boldsymbol{y}_{\boldsymbol{z}}^{\top} \widetilde{\mathbf{A}}(t)^{-1} \boldsymbol{y}_{\boldsymbol{z}})^{2}}$$

$$\leq 16 \sum_{\boldsymbol{v} \in \mathcal{Z}} t_{\boldsymbol{v}} \cdot (\boldsymbol{y}_{\boldsymbol{z}}^{\top} \boldsymbol{\theta}_{\star} + \epsilon)^{2} \frac{\boldsymbol{y}_{\boldsymbol{z}}^{\top} \widetilde{\mathbf{A}}(t)^{-1} \widetilde{\mathbf{A}}(t) \widetilde{\mathbf{A}}(t)^{-1} \boldsymbol{y}_{\boldsymbol{z}}}{(\boldsymbol{y}_{\boldsymbol{z}}^{\top} \widetilde{\mathbf{A}}(t)^{-1} \boldsymbol{y}_{\boldsymbol{z}})^{2}}$$

$$= \sum_{\boldsymbol{v} \in \mathcal{Z}} t_{\boldsymbol{v}} \cdot \frac{16(\boldsymbol{y}_{\boldsymbol{z}}^{\top} \boldsymbol{\theta}_{\star} + \epsilon)^{2}}{\|\boldsymbol{y}_{\boldsymbol{z}}\|_{\widetilde{\mathbf{A}}(t)^{-1}}^{2}}$$

Thus:

$$(\mathbf{E}.4) \ge \min \sum_{\boldsymbol{v} \in \mathcal{Z}} t_{\boldsymbol{v}} \quad \text{s.t.} \quad \min_{\boldsymbol{z} \neq \boldsymbol{z}^{\star}} \sum_{\boldsymbol{v} \in \mathcal{Z}} t_{\boldsymbol{v}} \text{KL}(\nu_{\boldsymbol{\theta}_{\star}, \boldsymbol{v}} || \nu_{\boldsymbol{\theta}_{\boldsymbol{z}}(\epsilon, t), \boldsymbol{v}}) \ge \log \frac{1}{2.4\delta}, \zeta \ge \frac{\sum_{i=2}^{d} t_{\boldsymbol{x}_{i}}}{\sum_{\boldsymbol{v} \in \mathcal{Z}} t_{\boldsymbol{v}}}$$

$$\ge \min \sum_{\boldsymbol{v} \in \mathcal{Z}} t_{\boldsymbol{v}} \quad \text{s.t.} \quad \sum_{\boldsymbol{v} \in \mathcal{Z}} t_{\boldsymbol{v}} \ge \max_{\boldsymbol{z} \neq \boldsymbol{z}^{\star}} \frac{\|\boldsymbol{y}_{\boldsymbol{z}}\|_{\widetilde{\mathbf{A}}(t)^{-1}}^{2}}{16(\boldsymbol{y}_{\boldsymbol{z}}^{\top} \boldsymbol{\theta}_{\star} + \epsilon)^{2}} \cdot \log \frac{1}{2.4\delta}, \zeta \ge \frac{\sum_{i=2}^{d} t_{\boldsymbol{x}_{i}}}{\sum_{\boldsymbol{v} \in \mathcal{Z}} t_{\boldsymbol{v}}}$$

$$= \inf_{\lambda \in \widetilde{\Delta}} \max_{\boldsymbol{z} \neq \boldsymbol{z}^{\star}} \frac{\|\boldsymbol{y}_{\boldsymbol{z}}\|_{\widetilde{\mathbf{A}}(\lambda)^{-1}}^{2}}{16(\boldsymbol{y}_{\boldsymbol{z}}^{\top} \boldsymbol{\theta}_{\star} + \epsilon)^{2}} \cdot \log \frac{1}{2.4\delta}$$

where  $\widetilde{\mathbf{A}}(\lambda) = \sum_{\boldsymbol{z} \in \mathcal{Z}} \lambda_{\boldsymbol{z}} \boldsymbol{z} \boldsymbol{z}^{\top}$  and  $\widetilde{\triangle} = \{\lambda \in \triangle_{\mathcal{Z}} : \zeta \geq \sum_{i=2}^{d} \lambda_{\boldsymbol{x}_i} \}$ . We can further lower bound this by

$$\geq \inf_{\lambda \in \widetilde{\Delta}} \max_{i \geq 2} \frac{\|\boldsymbol{z}^{\star} - \boldsymbol{x}_i\|_{\widetilde{\mathbf{A}}(\lambda)^{-1}}^2}{16((\boldsymbol{z}^{\star} - \boldsymbol{x}_i)^{\top} \boldsymbol{\theta}_{\star} + \epsilon)^2} \cdot \log \frac{1}{2.4\delta}$$

$$= \inf_{\lambda \in \widetilde{\Delta}} \max_{i \geq 2} \frac{\|\Delta \boldsymbol{e}_1 - \gamma \boldsymbol{e}_i\|_{\widetilde{\mathbf{A}}(\lambda)^{-1}}^2}{16(\Delta + \epsilon)^2} \cdot \log \frac{1}{2.4\delta}.$$

By Lemma E.4, we have

$$\begin{split} &\inf_{\lambda \in \widetilde{\Delta}} \max_{i \geq 2} \|\Delta \boldsymbol{e}_{1} - \gamma \boldsymbol{e}_{i}\|_{\widetilde{\mathbf{A}}(\lambda)^{-1}}^{2} \\ &\geq \inf_{\lambda \in \Delta_{d}} \max_{i \geq 2} (\Delta \boldsymbol{e}_{1} - \gamma \boldsymbol{e}_{i})^{\top} \left( 2\xi^{2} \boldsymbol{e}_{1} \boldsymbol{e}_{1}^{\top} + 2 \max\{\zeta, \gamma^{2}\} \lambda_{i} \boldsymbol{e}_{i} \boldsymbol{e}_{i}^{\top} + \operatorname{diag}([\xi^{2}, \gamma^{2}/d, \dots, \gamma^{2}/d]) \right)^{-1} (\Delta \boldsymbol{e}_{1} - \gamma \boldsymbol{e}_{i}) \\ &\geq \inf_{\lambda \in \Delta_{d}} \max_{i \geq 2} (\Delta \boldsymbol{e}_{1} - \gamma \boldsymbol{e}_{i})^{\top} \left( 3\xi^{2} \boldsymbol{e}_{1} \boldsymbol{e}_{1}^{\top} + (2 \max\{\zeta, \gamma^{2}\} \lambda_{i} + \gamma^{2}/d) \boldsymbol{e}_{i} \boldsymbol{e}_{i}^{\top} \right)^{-1} (\Delta \boldsymbol{e}_{1} - \gamma \boldsymbol{e}_{i}) \\ &= \frac{\Delta^{2}}{3\xi^{2}} + \inf_{\lambda \in \Delta_{d}} \max_{i \geq 2} \frac{1}{2\lambda_{i} + 1/d} \end{split}$$

where in the final equality we have used that  $\zeta \leq \gamma^2$ . However, this is clearly minimized by choosing  $\lambda_i = 1/(d-1)$ , which gives a lower bound of

$$\frac{1}{2/(d-1)+1/d} \ge \frac{d-1}{3}.$$

Putting all of this together, we have shown that any feasibly solution  $(t_z)_{z\in\mathcal{Z}}$  to (E.3) must satisfy

$$\sum_{z \in \mathcal{Z}} t_z \ge \frac{d-1}{48(\Delta + \epsilon)^2} \cdot \log \frac{1}{2.4\delta}.$$

Using that any feasible solution to (E.3) lower bounds  $\sum_{z\in\mathcal{Z}} \mathbb{E}[T(z)]$  and taking  $\epsilon\to 0$  gives the result.

**Lemma E.3.** Take some  $\Delta > 0$  satisfying:

$$\Delta \le \min \left\{ \frac{1}{2704d^2}, \sqrt{\frac{1}{10816C_2}}, \left(\frac{1}{10816d^{\alpha}C_1}\right)^{\frac{1}{2(1-\alpha)}} \right\}$$

and set

$$\xi = \frac{1}{52d}, \qquad \gamma = \max\left\{\frac{2\mathcal{C}_1}{(d/\Delta^2)^{1-\alpha}} + \frac{2\mathcal{C}_2\Delta^2}{d}, d\Delta\right\}.$$

Then this choice of  $\xi, \gamma, \delta$  satisfies (E.1) and, furthermore,  $\|\mathbf{z}\|_2 \leq 1$  for all  $\mathbf{z} \in \mathcal{Z}$ .

*Proof.* To satisfy (E.1), we must have  $\frac{1}{52\sqrt{d}} \geq \gamma$ . Thus, if

$$\Delta \le \frac{1}{2704d^2}, \quad \Delta \le \sqrt{\frac{1}{10816C_2}}, \quad \Delta \le \left(\frac{1}{10816d^{\alpha}C_1}\right)^{\frac{1}{2(1-\alpha)}},$$

some computation shows that choosing  $\gamma$  as prescribed will meet the constraint  $\frac{1}{52\sqrt{d}} \geq \gamma$  and will also satisfy  $\gamma \geq \sqrt{d\Delta}$ . The norm bound follows by our choice of  $\xi$  and since  $\xi \geq \gamma/\sqrt{d} \geq \sqrt{\Delta}$ .

# E.1.1 Additional Proofs

Proof of Claim E.2. We first show that  $|\langle \boldsymbol{\theta}_{\boldsymbol{z}}(\epsilon,t), \boldsymbol{v} \rangle| \leq 13d\xi$  for all  $\boldsymbol{z}, \boldsymbol{v} \in \mathcal{Z}$ . Note that for all  $\boldsymbol{v} \in \mathcal{Z}, |\langle \boldsymbol{v}, \boldsymbol{\theta}_{\star} \rangle| \leq \xi$ .

Case 1:  $z = z^*$ . In this case,  $\langle y_z, \theta_* \rangle = 0$  so the result follows from our condition on  $\epsilon$ .

Case 2:  $z = e_i, i \geq 2$ . Let  $\widetilde{\triangle} = \{\lambda \in \triangle_{\mathcal{Z}} : \zeta \geq \sum_{i=2}^d \lambda_{x_i} \}$ . In this case,  $\langle z, \theta_{\star} \rangle = 0$  and  $\langle y_z, \theta_{\star} \rangle = \xi$ . Furthermore, by Lemma E.4,

$$\begin{aligned} \boldsymbol{y}_{\boldsymbol{z}}^{\top} \widetilde{\mathbf{A}}(t)^{-1} \boldsymbol{y}_{\boldsymbol{z}} &\geq \inf_{\lambda \in \widetilde{\Delta}} \boldsymbol{y}_{\boldsymbol{z}}^{\top} \left( 2 \sum_{\boldsymbol{z}' \in \mathcal{Z}} \lambda_{\boldsymbol{z}'} \mathrm{diag}([(\boldsymbol{z}')^2]) + \mathrm{diag}([\xi^2, \gamma^2/d, \dots, \gamma^2/d]) \right)^{-1} \boldsymbol{y}_{\boldsymbol{z}} \\ &\geq \boldsymbol{y}_{\boldsymbol{z}}^{\top} \left( 2 \xi^2 \boldsymbol{e}_1 \boldsymbol{e}_1^{\top} + 2 \max\{\zeta, \gamma^2\} \boldsymbol{e}_i \boldsymbol{e}_i^{\top} + \mathrm{diag}([\xi^2, \gamma^2/d, \dots, \gamma^2/d]) \right)^{-1} \boldsymbol{y}_{\boldsymbol{z}} \\ &\geq \xi^2 \frac{3}{\xi^2} + \frac{1}{2 \max\{\zeta, \gamma^2\} + \gamma^2/d} \\ &= 3 + \frac{1}{2 \max\{\zeta, \gamma^2\} + \gamma^2/d} \\ &\geq \frac{1}{3\gamma^2} \end{aligned}$$

where the last inequality follows from our assumption that  $\zeta \leq \gamma^2$ . In the other direction, we can bound

$$\boldsymbol{v}^{\top}\widetilde{\mathbf{A}}(t)^{-1}\boldsymbol{y_z} \leq \boldsymbol{v}^{\top}\mathrm{diag}([\xi^2,\gamma^2/d,\ldots,\gamma^2/d])^{-1}\boldsymbol{y_z} \leq \frac{1}{\xi^2} + \frac{d}{\gamma^2} \leq \frac{2d}{\gamma^2}$$

where the last inequality follows by our assumption that  $\xi \geq \gamma/\sqrt{d}$ . Putting this together, we have

$$|\langle \boldsymbol{\theta}_{\boldsymbol{z}}(\epsilon, t), \boldsymbol{v} \rangle| \le \xi + (\xi + \epsilon) \frac{2d/\gamma^2}{1/3\gamma^2} \le 13d\xi.$$

Case 3:  $z = x_i$ . In this case  $\langle y_z, \theta_{\star} \rangle = \Delta$ . We can apply a calculation analogous to above to lower bound  $y_z^{\top} \widetilde{\mathbf{A}}(t)^{-1} y_z$ , but in this case obtain

$$\boldsymbol{y}_{\boldsymbol{z}}^{\top}\widetilde{\mathbf{A}}(t)^{-1}\boldsymbol{y}_{\boldsymbol{z}} \geq \Delta^{2}\frac{3}{\xi^{2}} + \frac{\gamma^{2}}{2\max\{\zeta,\gamma^{2}\} + \gamma^{2}/d} \geq \frac{1}{3}.$$

Similarly, we can upper bound

$$\boldsymbol{v}^{\top}\widetilde{\mathbf{A}}(t)^{-1}\boldsymbol{y_z} \leq \boldsymbol{v}^{\top}\mathrm{diag}([\xi^2,\gamma^2/d,\ldots,\gamma^2/d])^{-1}\boldsymbol{y_z} \leq \frac{\Delta}{\xi^2} + \frac{d\gamma}{\gamma^2} \leq \frac{2d}{\gamma}.$$

This gives a final upper bound of

$$|\langle \boldsymbol{\theta}_{z}(\epsilon, t), \boldsymbol{v} \rangle| \leq \xi + (\Delta + \epsilon) \frac{6d}{\gamma} \leq \xi + \frac{12d\Delta}{\sqrt{\Delta}} \leq 13d\xi.$$

Combining these three cases gives that  $|\langle \boldsymbol{\theta}_{\boldsymbol{z}}(\epsilon,t), \boldsymbol{v} \rangle| \leq 13d\xi$  for all  $\boldsymbol{z}, \boldsymbol{v} \in \mathcal{Z}$ . By our assumption that  $\xi \leq \frac{1}{52d}$ , it follows that  $|\langle \boldsymbol{\theta}_{\boldsymbol{z}}(\epsilon,t), \boldsymbol{v} \rangle| \leq 1/4$  for all  $\boldsymbol{z}, \boldsymbol{v} \in \mathcal{Z}$ .

By Lemma D.2 of Wagenmaker et al. (2022), as long as  $\langle \boldsymbol{\theta}_{z}(\epsilon,t), \boldsymbol{v} \rangle + 1/2 \in (0,1)$  and  $\langle \boldsymbol{\theta}_{\star}, \boldsymbol{v} \rangle + 1/2 \in (0,1)$ , which will be the case by the definition of  $\boldsymbol{\theta}_{\star}$  and since  $|\langle \boldsymbol{\theta}_{z}(\epsilon,t), \boldsymbol{v} \rangle| \leq 1/4$  as noted above, we have

$$\mathrm{KL}(\nu_{\boldsymbol{\theta_{\star}},\boldsymbol{v}}||\nu_{\boldsymbol{\theta_{z}}(\epsilon,t),\boldsymbol{v}}) \leq \frac{\langle \boldsymbol{\theta_{z}}(\epsilon,t) - \boldsymbol{\theta_{\star}}, \boldsymbol{v} \rangle^{2}}{(\langle \boldsymbol{\theta_{z}}(\epsilon,t), \boldsymbol{v} \rangle + 1/2)(1/2 - \langle \boldsymbol{\theta_{z}}(\epsilon,t), \boldsymbol{v} \rangle)}.$$

Using what we have just shown, we can upper bound this as

$$\frac{\langle \boldsymbol{\theta}_{\boldsymbol{z}}(\epsilon,t) - \boldsymbol{\theta}_{\star}, \boldsymbol{v} \rangle^{2}}{(\langle \boldsymbol{\theta}_{\boldsymbol{z}}(\epsilon,t), \boldsymbol{v} \rangle + 1/2)(1/2 - \langle \boldsymbol{\theta}_{\boldsymbol{z}}(\epsilon,t), \boldsymbol{v} \rangle)} \leq \frac{\langle \boldsymbol{\theta}_{\boldsymbol{z}}(\epsilon,t) - \boldsymbol{\theta}_{\star}, \boldsymbol{v} \rangle^{2}}{(-1/4 + 1/2)(1/2 - 1/4)}$$

$$= 16\langle \boldsymbol{\theta}_{\boldsymbol{z}}(\epsilon,t) - \boldsymbol{\theta}_{\star}, \boldsymbol{v} \rangle^{2}.$$

By our choice of  $\theta_z(\epsilon, t)$ , this is equal to:

$$16(\boldsymbol{y}_{\boldsymbol{z}}^{\top}\boldsymbol{\theta}_{\star} + \epsilon)^{2} \frac{\boldsymbol{y}_{\boldsymbol{z}}^{\top}\widetilde{\mathbf{A}}(t)^{-1}\boldsymbol{v}\boldsymbol{v}^{\top}\widetilde{\mathbf{A}}(t)^{-1}\boldsymbol{y}_{\boldsymbol{z}}}{(\boldsymbol{y}_{\boldsymbol{z}}^{\top}\widetilde{\mathbf{A}}(t)^{-1}\boldsymbol{y}_{\boldsymbol{z}})^{2}}$$

which completes the proof.

### Lemma E.4.

$$\sum_{\boldsymbol{z} \in \mathcal{Z}} \lambda_{\boldsymbol{z}} \boldsymbol{z} \boldsymbol{z}^\top \preceq 2 \sum_{\boldsymbol{z} \in \mathcal{Z}} \lambda_{\boldsymbol{z}} \mathrm{diag}(\boldsymbol{z}^2).$$

*Proof.* This follows since every  $z \in \mathcal{Z}$  has at most two non-zero entries, and since  $(ax + by)(ax + by)^{\top} \leq 2a^2xx^{\top} + 2b^2yy^{\top}$ .

# E.2 Mapping to Linear MDPs

We can map this linear bandit (with parameters chose as in Lemma E.3) to a linear MDP with state space  $S = \{s_0, s_1, \bar{s}_2, \dots, \bar{s}_{d+1}\}$ , action space  $A = Z \cup \{e_{d+1}/2\}$ , parameters

$$m{ heta}_1 = m{0}, \quad m{ heta}_2 = m{e}_1 \ m{\mu}_1(s_1) = [2m{ heta}_\star, 1], \quad m{\mu}_1(ar{s}_i) = rac{1}{d}[-2m{ heta}_\star, 1],$$

and feature vectors

$$\phi(s_0, e_{d+1}) = e_{d+1}/2, \quad \phi(s_0, z) = [z/2, 1/2], \quad \forall z \in \mathcal{Z}$$
  
 $\phi(s_1, z) = e_1, \quad \phi(\bar{s}_i, z) = e_i, i \ge 2, \quad \forall z \in \mathcal{A}.$ 

Note that, if we take action z in state  $s_0$ , our expected episode reward is

$$P_1(s_1|s_0, z) \cdot 1 + \sum_{i=2}^{d+1} P_1(\bar{s}_i|s_0, z) \cdot 0 = \langle \theta_{\star}, z \rangle + 1/2$$

since we always acquire a reward of 1 in any state  $s_1$ , and a reward of 0 in any state  $\bar{s}_i$ , and the reward distribution is Bernoulli.

**Lemma E.5.** The MDP constructed above is a valid linear MDP as defined in Definition 3.1.

*Proof.* For  $z \in \mathcal{Z}$  we have.

$$P_1(s_1|s_0, \mathbf{z}) = \langle \boldsymbol{\phi}(s_0, \mathbf{z}), \boldsymbol{\mu}_1(s_1) \rangle = \langle \boldsymbol{\theta}_{\star}, \mathbf{z} \rangle + 1/2 \ge 0$$

$$P_1(\bar{s}_i|s_0, \mathbf{z}) = \langle \boldsymbol{\phi}(s_0, \mathbf{z}), \boldsymbol{\mu}_1(\bar{s}_i) \rangle = \frac{1}{d} (-\langle \boldsymbol{\theta}_{\star}, \mathbf{z} \rangle + 1/2) \ge 0$$

where the inequality follows since  $|\langle \boldsymbol{\theta}_{\star}, \boldsymbol{z} \rangle| \leq \mathcal{O}(1/d)$  for all  $\boldsymbol{z} \in \mathcal{Z}$ . In addition,

$$P_1(s_1|s_0, z) + \sum_{i=2}^{d+1} P_1(\bar{s}_i|s_0, z) = \langle \theta_{\star}, z \rangle + 1/2 + d \cdot \frac{1}{d} (-\langle \theta_{\star}, z \rangle + 1/2) = 1.$$

Thus,  $P_1(\cdot|s_0, \mathbf{z})$  is a valid probability distribution for  $\mathbf{z} \in \mathcal{Z}$ . A similar calculation shows the same for  $\mathbf{z} = \mathbf{e}_{d+1}/2$ .

It remains to check the normalization bounds. Clearly, by our construction of  $\mathcal{Z}$ ,  $\|\phi(s,a)\|_2 \leq 1$  for all s and a. It is also obvious that  $\|\theta_0\|_2 \leq \sqrt{d}$  and  $\|\theta_1\|_2 \leq \sqrt{d}$ . Finally,

$$\||\boldsymbol{\mu}_1(\mathcal{S})|\|_2 = \left\|\sum_{s \in \mathcal{S} \setminus s_0} |\boldsymbol{\mu}_1(s)|\right\|_2 = \|[2\boldsymbol{\theta}_{\star}, 1] + d \cdot \frac{1}{d}[2\boldsymbol{\theta}_{\star}, 1]\|_2 \le \sqrt{d}.$$

Thus, all normalization bounds are met, so this is a valid linear MDP.

Proof of Proposition 2. If we assume that the learner has prior access to the feature vectors, and also knows this is a linear MDP, then, even with no knowledge of the dynamics, we can guarantee an optimal policy is contained in the set of policies  $\pi^{z,z'}$  defined as:

$$\pi_1^{z,z'}(s_0) = z, \pi_2^{z,z'}(s_0) = z', \pi_h^{z,z'}(s_1) = \xi e_1, \pi_h^{z,z'}(\bar{s}_i) = \xi e_1$$

This holds because in states  $s_1$  and  $\bar{s}_i$ , the performance of each action is identical since the feature vectors are identical, so it doesn't matter which action we choose in these states. In this case, we can bound  $|\Pi| \leq |\mathcal{Z}|^2 \leq 4d^2$ .

Now, for  $z \in A$ ,  $z \neq e_{d+1}/2$ , we have

$$\begin{split} \boldsymbol{\phi}_{\pi^{\boldsymbol{z},\boldsymbol{z}'},1} &= [\boldsymbol{z}/2,1/2] \\ \boldsymbol{\phi}_{\pi^{\boldsymbol{z},\boldsymbol{z}'},2} &= (\langle \boldsymbol{\theta}_{\star},\boldsymbol{z} \rangle + 1/2)\boldsymbol{e}_1 + \frac{1}{d}(-\langle \boldsymbol{\theta}_{\star},\boldsymbol{z} \rangle + 1/2)\sum_{i \geq 2}\boldsymbol{e}_i \end{split}$$

and if  $z = e_{d+1}/2$ ,  $\phi_{\pi^{z,z'},1} = e_{d+1}/2$ ,  $\phi_{\pi^{z,z'},2} = e_{1}/2 + \frac{1}{2d} \sum_{i \geq 2} e_i$ . Let  $\pi_{\exp}$  be the policy that plays action  $e_2$  in state  $s_0$  at step h = 1. Then,

$$oldsymbol{\Lambda}_{\pi_{ ext{exp}},2} = rac{1}{2}oldsymbol{e}_1oldsymbol{e}_1 + rac{1}{2d}\sum_{i\geq 3}oldsymbol{e}_ioldsymbol{e}_i^{ op}.$$

Since  $\langle \boldsymbol{\theta}_{\star}, \boldsymbol{z} \rangle \leq \mathcal{O}(1/d)$  and  $[\boldsymbol{z}]_1 \leq \mathcal{O}(1/d)$  for all  $\boldsymbol{z}$  by construction, it follows that we can bound, for all  $\boldsymbol{z}, \boldsymbol{z}'$ ,

$$\|\boldsymbol{\phi}_{\pi^{\boldsymbol{z},\boldsymbol{z}'},2}\|_{\boldsymbol{\Lambda}_{\pi_{\mathrm{exp}},2}^{-1}}^{2} = \mathcal{O}\left(1 + \sum_{i \geq 2} \frac{1}{d^{2}} \cdot d\right) = \mathcal{O}(1)$$

SO

$$\inf_{\pi_{\text{exp}}} \max_{\pi \in \Pi} \frac{\|\boldsymbol{\phi}_{\pi,2}\|_{\boldsymbol{\Sigma}_{\pi \text{exp}},2}^{-1}}{\max\{V_0^{\star} - V_0^{\pi}, \Delta_{\min}^{\Pi}, \epsilon\}^2} \leq \mathcal{O}(1/\epsilon^2).$$

Now let  $\pi_{\text{exp}}$  be the policy that, at step h = 1, plays  $\xi e_1$  with probability 1/4,  $e_{d+1}$  with probability 1/4, and plays  $e_i$  with probability  $\frac{1}{4(d-1)}$  for  $i \geq \{2, \ldots, d\}$ . In this setting, we have

$$\boldsymbol{\Lambda}_{\pi_{\text{exp}},1} = \frac{1}{4} \xi^2 \boldsymbol{e}_1 \boldsymbol{e}_1^\top + \frac{1}{4} \boldsymbol{e}_{d+1} \boldsymbol{e}_{d+1} + \frac{1}{4(d-1)} \sum_{i \in \{2,...,d\}} \boldsymbol{e}_i \boldsymbol{e}_i^\top.$$

Note that  $V_0^{\star} - V_0^{\pi^{z,z'}} = \xi - \langle \boldsymbol{\theta}_{\star}, \boldsymbol{z} \rangle$ , so for  $\boldsymbol{z} = \boldsymbol{e}_2, \dots, \boldsymbol{e}_{d+1}$ , we have  $V_0^{\star} - V_0^{\pi^{z,z'}} = \xi = \mathcal{O}(1/d)$ , while for  $\boldsymbol{z} = \xi \boldsymbol{e}_1, \boldsymbol{x}_2, \dots, \boldsymbol{x}_d$ , we have  $V_0^{\star} - V_0^{\pi^{z,z'}} = \Delta$ .

It's easy to see that for  $z = e_2, \ldots, e_{d+1}$ , we have  $\|\phi_{\pi^{z,z'},1}\|_{\Lambda_{\text{mexp},1}^{-1}} \leq \mathcal{O}(d)$ , and for  $z = \xi e_1, x_2, \ldots, x_d$ ,  $\|\phi_{\pi^{z,z'},1}\|_{\Lambda_{\text{mexp},1}^{-1}} \leq \mathcal{O}(1+d\gamma^2) = \mathcal{O}(1)$ . Combining these bounds with the gap values, we conclude that

$$\inf_{\pi_{\text{exp}}} \max_{\pi \in \Pi} \frac{\|\phi_{\pi,1}\|_{\mathbf{\Sigma}_{\pi \text{exp},1}^{-1}}^2}{\max\{V_0^{\star} - V_0^{\pi}, \Delta_{\min}^{\Pi}, \epsilon\}^2} \leq \mathcal{O}(1/\epsilon^2 + \text{poly}(d)).$$

The result then follows by Theorem 7.

Lower bounding the performance of low-regret algorithms. Assume that we have access to the linear bandit instance constructed in Appendix E.1 with parameters chosen as in Lemma E.3. That is, at every timestep t we can choose an arm  $z_t \in \mathcal{Z}$  and obtain and observe reward  $y_t \sim \text{Bernoulli}(\langle \theta_{\star}, z_t \rangle + 1/2)$ . Using the mapping up, we can use this bandit to simulate a linear MDP as follows:

- 1. Start in state  $s_0$  and choose any action  $z_t \in \mathcal{A}$
- 2. Play action  $z_t$  in our linear bandit. If reward obtained is  $y_t = 1$ , then in MDP transition to any of the states  $s_1$ . If reward obtained is  $y_t = 0$  transition to any of the states  $\bar{s}_2, \ldots, \bar{s}_{d+1}$ , each with probability 1/d. If the chosen action was  $z_t = e_{d+1}/2$ , then play any action in the linear bandit and transition to state  $s_1$  with probability 1/2 and  $\bar{s}_2, \ldots, \bar{s}_{d+1}$  with probability 1/2d, regardless of  $y_t$
- 3. Take any action in the state in which you end up, and receive reward of 1 if you are in  $s_1$ , and reward of 0 if you are in  $\bar{s}_2, \ldots, \bar{s}_{d+1}$ .

Note that this MDP has precisely the transition and reward structure as the MDP constructed above.

**Lemma E.6.** Assume  $\pi$  is  $\epsilon < \Delta/2$ -optimal in the MDP constructed above. Then,  $\mathbf{z}^* = \arg\max_{\mathbf{z} \in \mathcal{A}} \pi_1(\mathbf{z}|s_0)$ .

*Proof.* Note that the value of  $\pi$  in the linear MDP is given by  $V_0^{\pi} = \sum_{z \in \mathcal{Z}} \pi_1(z|s_0)(\langle z, \theta_{\star} \rangle + 1/2) + \pi_1(e_{d+1}/2|s_0)/2$  and the optimal policy is  $\pi_1(z^{\star}|s_0) = 1$  and has value  $V_0^{\star} = \langle z^{\star}, \theta_{\star} \rangle + 1/2$ . It follows that if  $\pi$  is  $\epsilon$ -optimal, then

$$\sum_{\boldsymbol{z}\in\mathcal{Z}} \pi_1(\boldsymbol{z}|s_0)(\langle \boldsymbol{z},\boldsymbol{\theta}_{\star}\rangle + 1/2) + \pi_1(\boldsymbol{e}_{d+1}/2|s_0)/2 \ge \langle \boldsymbol{z}^{\star},\boldsymbol{\theta}_{\star}\rangle + 1/2 - \epsilon$$

$$\implies \pi_1(\boldsymbol{z}^{\star}|s_0)(\xi + 1/2) + \sum_{\boldsymbol{z}\in\mathcal{A},\boldsymbol{z}\neq\boldsymbol{z}^{\star}} \pi_1(\boldsymbol{z}|s_0)(\xi - \Delta + 1/2) \ge \xi + 1/2 - \epsilon$$

$$\implies -\Delta \sum_{\boldsymbol{z}\in\mathcal{A},\boldsymbol{z}\neq\boldsymbol{z}^{\star}} \pi_1(\boldsymbol{z}|s_0) \ge -\epsilon$$

$$\implies \epsilon \ge \Delta \sum_{\boldsymbol{z}\in\mathcal{A},\boldsymbol{z}\neq\boldsymbol{z}^{\star}} \pi_1(\boldsymbol{z}|s_0).$$

If  $\epsilon < \Delta/2$ , this implies that  $\sum_{\boldsymbol{z} \in \mathcal{A}, \boldsymbol{z} \neq \boldsymbol{z}^*} \pi_1(\boldsymbol{z}|s_0) < 1/2$ , so it must be the case that  $\pi_1(\boldsymbol{z}^*|s_0) > 1/2$ .

Proof of Proposition 3. Consider running the above procedure for some number of steps. By Lemma E.6, if we can identify an  $\epsilon < \Delta/2$ -optimal policy in this MDP, we can use it to determine  $z^*$ , the optimal arm in the linear bandit. As we have used no extra information other than samples from the linear bandit to construct this, it follows that to find an  $\epsilon < \Delta/2$ -optimal policy in the MDP, we must take at least the number of samples prescribed by Lemma E.1.