

---

# Sparse Optimistic Information Directed Sampling

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Many high-dimensional decision-making problems can be modeled as stochastic  
2 sparse linear bandits. Most existing algorithms are designed to achieve opti-  
3 mal worst-case regret in either the data-rich regime, where polynomial depen-  
4 dence on the ambient dimension is unavoidable, or the data-poor regime, where  
5 dimension-independence is possible at the cost of worse dependence on the num-  
6 ber of rounds. In contrast, the Bayesian approach of Information Directed Sam-  
7 pling (IDS) achieves the best of both worlds: a Bayesian regret bound that has  
8 the optimal rate in both regimes simultaneously. In this work, we explore the use  
9 of Sparse Optimistic Information Directed Sampling (SOIDS) to achieve the best  
10 of both worlds in the worst-case setting, without Bayesian assumptions. Through  
11 a novel analysis that enables the use of a time-dependent learning rate, we show  
12 that SOIDS can optimally balance information and regret. Our results extend the  
13 theoretical guarantees of IDS, providing the first algorithm that simultaneously  
14 achieves optimal worst-case regret in both the data-rich and data-poor regimes. In  
15 addition, we empirically demonstrate the good performance of SOIDS.

## 16 1 Introduction

17 In stochastic linear bandits, one assumes that the mean reward associated with each action is linear  
18 in an unknown  $d$ -dimensional parameter vector [Abe and Long, 1999, Auer, 2002, Dani et al., 2008,  
19 Abbasi-Yadkori et al., 2011]. Under standard conditions, it is known that the minimax regret in this  
20 setting is of the order  $\mathcal{O}(d\sqrt{T})$  [Dani et al., 2008, Rusmevichientong and Tsitsiklis, 2010]. Nu-  
21 merous follow-up works have investigated the possibility of reduced regret under various structural  
22 assumptions on the unknown parameter vector, the noise, or the shape of the decision set [Valko  
23 et al., 2014, Chu et al., 2011, Kirschner and Krause, 2018], [Lattimore and Szepesvári, 2020, Chap-  
24 ter 22]. One such assumption is that the unknown parameter vector is *sparse*, which means that it  
25 has only  $s \ll d$  non-zero components. This setting is called *sparse linear bandits* and  $s$  is referred to  
26 as the *sparsity level*. In this setting, previous work has established the existence of algorithms with  
27 regret scaling as  $\mathcal{O}(\sqrt{sdT})$  [Abbasi-Yadkori et al., 2012]. This result is complemented by a lower  
28 bound, which says that this rate cannot be improved as long as  $T \geq d^\alpha$  for some  $\alpha > 0$  [Lattimore  
29 and Szepesvári, 2020]. We refer to this scenario as the *data-rich regime*. Since this bound scales  
30 polynomially with the dimension  $d$ , many researchers have considered this to be a negative result,  
31 interpreting it as a sign that sparsity cannot be effectively exploited in linear bandit problems. This  
32 interpretation has been challenged by a more recent observation that, when the action set admits  
33 an *exploratory distribution*, simple “explore-then-commit” algorithms enjoy regret bounds of order  
34  $\mathcal{O}(\text{poly}(s)T^{\frac{2}{3}})$  [Hao et al., 2020, Jang et al., 2022]. These bounds scale only logarithmically with  
35 the dimension, and constitute a major improvement over the previously mentioned rate in the *data-*  
36 *poor regime*, where  $T \ll (\frac{d}{s})^3$ . Most known algorithms are specialized to either the data-poor or  
37 data-rich regime, and perform poorly in the other one. A notable exception is the *sparse Information*  
38 *Directed Sampling* algorithm introduced in Hao et al. [2021], which performs almost optimally in  
39 both regimes. However, Hao et al. [2021] only provide *Bayesian* regret bounds for sparse IDS.

40 In this work, we lift this assumption and develop an algorithm that can adapt to both regimes in a  
 41 “frequentist” sense. The algorithm is an adaptation of the Optimistic Information Directed Sampling  
 42 (OIDS) algorithm of [Neu, Papini, and Schwartz \[2024\]](#). Our contribution is as follows:

- 43 • We extend the analysis of the optimistic posterior to allow the use of time-dependent learn-  
 44 ing rates and history-dependent learning rates. This removes the need to know the horizon  
 45 in advance and allows us to update the learning rate based on data observed by the agent  
 46 instead of some loose theoretical constant, a necessity for efficient algorithms.
- 47 • We demonstrate that the SOIDS algorithm recovers almost optimal rates in both the data-  
 48 poor and data-rich regimes. This is the first algorithm to do so in a frequentist setting.

## 49 2 Preliminaries

50 **Sparse linear bandits.** We consider the following decision-making game, in which a learning  
 51 agent interacts with an environment over a sequence of  $T$  rounds. At the start of each round  $t$ , the  
 52 learner selects an action  $A_t \in \mathcal{A} \subset \mathbb{R}^d$  according to a randomized policy  $\pi_t \in \Delta(\mathcal{A})$ . In response,  
 53 the environment generates a stochastic reward  $Y_t = r(A_t) + \epsilon_t$ , where  $r : \mathcal{A} \rightarrow \mathbb{R}$  is a fixed reward  
 54 function and  $\epsilon_t$  is zero-mean, conditionally 1-sub-Gaussian noise. We assume that the action set  $\mathcal{A}$   
 55 is finite, and that the reward function can be written as

$$r(a) = \langle \theta_0, a \rangle,$$

56 where  $\theta_0 \in \mathbb{R}^d$  is an unknown parameter vector. We make the mild boundedness assumptions  
 57 that  $\max_{a \in \mathcal{A}} \|a\|_\infty \leq 1$  and  $\|\theta_0\|_1 \leq 1$ . We study the special case of this problem in which the  
 58 parameter vector  $\theta_0$  is  $s$ -sparse in the sense that at most  $s \ll d$  of its components are non-zero. In  
 59 other words, we assume that  $\theta_0$  belongs to the following *sparse parameter space*:

$$\Theta = \left\{ \theta \in \mathbb{R}^d : \sum_{j=1}^d \mathbb{I}_{\{\theta_j \neq 0\}} \leq s, \|\theta\|_1 \leq 1 \right\}.$$

60 We assume that the sparsity level  $s$  is known to the agent. The performance of the agent is evaluated  
 61 in terms of the *regret*, which is defined as

$$R_T = T \max_{a \in \mathcal{A}} \langle \theta_0, a \rangle - \mathbb{E} \left[ \sum_{t=1}^T r(A_t, \theta_0) \right], \quad (1)$$

62 where the expectation is taken with respect to both the random choices of the agent and the random  
 63 noise in the observed rewards. We note that the regret is implicitly a function of the true parameter  
 64  $\theta_0$ . Our focus is on proving regret bounds that hold for arbitrary choices of  $\theta_0 \in \Theta$ .

65 **The data-rich and data-poor regimes.** As mentioned in the introduction, it is known there exist  
 66 algorithms for sparse linear bandits with worst-case regret of the order  $\mathcal{O}(\sqrt{sdT})$  [[Abbasi-Yadkori](#)  
 67 [et al., 2012](#)]. This regret bound is only meaningful when the dimension  $d$  is smaller than the number  
 68 of rounds  $T$ , a situation referred to as the data-rich regime. Under the assumption that there exists  
 69 an exploratory policy, [Hao et al. \[2020\]](#) showed that there is a simple algorithm that satisfies a  
 70 problem-dependent regret bound, which can be meaningful in the so-called data-poor regime, where  
 71  $d$  is much larger than  $T$ . Formally, we say that there exists an exploratory policy if the action set  $\mathcal{A}$   
 72 is such that

$$C_{\min} := \max_{\mu \in \Delta(\mathcal{A})} \sigma_{\min} \left( \int_{\mathcal{A}} aa^\top d\mu(a) \right) > 0,$$

73 which is equivalent to the condition that  $\mathcal{A}$  spans  $\mathbb{R}^d$ . The exploratory policy, is the distribution  
 74 on  $\mathcal{A}$  that achieves the maximum (which is guaranteed to exist when  $\mathcal{A}$  is finite). The Explore  
 75 the Sparsity Then Commit (ESTC) algorithm was shown to satisfy a regret bound of the order  
 76  $\mathcal{O}(s^{2/3}T^{2/3}C_{\min}^{-2/3})$  [[Hao et al., 2020](#)]. The transition between the  $T^{2/3}$  rate in the data-poor  
 77 regime and the  $\sqrt{T}$  rate in the data-rich regime also appears in an existing lower bound of the  
 78 order  $\Omega(\min(s^{1/3}T^{2/3}C_{\min}^{-1/3}, \sqrt{dT}))$  [[Hao et al., 2020](#)].

79 **The sparse optimal action condition.** Part of our analysis requires that a certain technical condi-  
 80 tion is satisfied. This condition comes from prior work [[Hao et al., 2021](#)], and is used to bound the  
 81 regret in the data-poor regime (cf. Lemma 7).

82 **Definition 1.** For a given prior  $Q_1^+$ , an action set  $\mathcal{A}$  has sparse optimal actions if with probability 1  
83 over the random draw of  $\theta$  from  $Q_1^+$ , there exists  $a' \in \arg \max_{a \in \mathcal{A}} r(a, \theta)$  such that  $\|a'\|_0 \leq s$ .

84 We use a prior that only assigns positive probability to  $s$ -sparse vectors, which means the sparse  
85 optimal action property is satisfied whenever the action set is an  $\ell_p$  ball. Note that the hard  
86 instances in both the  $\sqrt{sdT}$  lower bound in Theorem 24.3 of [Lattimore and Szepesvári \[2020\]](#) and the  
87  $s^{2/3}T^{2/3}$  lower bound in Theorem 5 of [Jang et al. \[2022\]](#) satisfy the sparse optimal action property<sup>1</sup>.  
88 Therefore, this additional condition does not trivialize the problem.

89 **Notation.** We conclude this section by introducing some additional notation that will be used in the  
90 subsequent sections. For any candidate parameter vector (or model)  $\theta \in \mathbb{R}^d$ , we let  $r(a, \theta) = \langle \theta, a \rangle$   
91 denote the corresponding linear reward function. In addition, we define  $a^*(\theta) = \arg \max_{a \in \mathcal{A}} r(a, \theta)$   
92 (with ties broken arbitrarily) and  $r^*(\theta) = r(a^*(\theta), \theta)$  to be the optimal action and maximum reward  
93 for the model  $\theta$ . The gap of an action  $a$  for a model  $\theta$  is  $\Delta(a, \theta) = r^*(\theta) - r(a, \theta)$ . Similarly, the  
94 gap for a policy  $\pi \in \Delta(\mathcal{A})$  and a model distribution  $Q \in \Delta(\Theta)$  is  $\Delta(\pi, Q) = \int_{\mathcal{A} \times \Theta} \Delta(a, \theta) d\pi \otimes$   
95  $Q(a, \theta)$ , and we let  $\Delta_t = \Delta(\pi_t, \theta_0)$  denote the gap of the policy played by the agent in round  $t$   
96 under the true model  $\theta_0$ . Using this notation, the regret can be written as  $R_T = \mathbb{E}[\sum_{t=1}^T \Delta_t]$ . We  
97 define the unnormalized Gaussian likelihood function  $p(y|\theta, a) = \exp(-\frac{(y - \langle \theta, a \rangle)^2}{2})$ . Finally, we  
98 let  $\mathcal{F}_t = \sigma(A_1, Y_1, \dots, A_t, Y_t)$  denote the  $\sigma$ -algebra generated by the interaction between the agent  
99 and the environment up to the end of round  $t$ .

### 100 3 Sparse Optimistic Information Directed Sampling

101 We develop an extension of the Optimistic Information Directed Sampling (OIDS) algorithm pro-  
102 posed by [Neu, Papini, and Schwartz \[2024\]](#). The main difference between OIDS and IDS is that  
103 the Bayesian posterior is replaced by an appropriately adjusted *optimistic posterior*. For an arbitrary  
104 prior  $Q_1^+ \in \Delta(\Theta)$ , the optimistic posterior is defined by the following update rule:

$$\frac{dQ_{t+1}^+}{dQ_1^+}(\theta) \propto \prod_{s=1}^t (p(Y_s | \theta, A_s))^\eta \cdot \exp\left(\lambda_t \sum_{s=1}^t \Delta(A_s, \theta)\right). \quad (2)$$

105 Here,  $\eta$  is a positive constant that should be thought of as “large”, and  $(\lambda_t)_t$  is a decreasing se-  
106 quence of positive real numbers that decays to 0, and should be thought of as “small”. Note that  
107 when  $\eta = 1$  and  $\lambda_t = 0$ , the optimistic posterior coincides with the Bayesian posterior. When  
108  $\lambda_t > 0$ , the  $\Delta(A_s, \theta)$  term promotes “overestimation” of the true gaps, driving exploration towards  
109 parameters that promise rewards much higher than whatever would have been accrued by the agent.  
110 This construction is closely related to the optimistic posterior update described in [Zhang \[2022\]](#) and  
111 [Neu, Papini, and Schwartz \[2024\]](#). To describe our algorithm, we must first define the *surrogate*  
112 *information gain* and the *surrogate regret*. For any round  $t$  and any policy  $\pi \in \Delta(\mathcal{A})$ , the surrogate  
113 information gain is defined as

$$\overline{\text{IG}}_t(\pi) = \frac{1}{2} \sum_{a \in \mathcal{A}} \pi(a) \int_{\Theta} (\langle \theta - \bar{\theta}(Q_t^+), a \rangle)^2 dQ_t^+(\theta),$$

114 where for any  $Q \in \Delta(\Theta)$ ,  $\bar{\theta}(Q) = \mathbb{E}_{\theta \sim Q}[\theta]$  is the mean parameter under distribution  $Q$ . The  
115 surrogate regret is defined as

$$\widehat{\Delta}_t(\pi) = \sum_{a \in \mathcal{A}} \pi(a) \int_{\Theta} \Delta(a, \theta) dQ_t^+(\theta).$$

116 For any policy  $\pi$  and any  $\gamma \geq 2$ , we define the *surrogate generalized information ratio* as

$$\overline{\text{IR}}_t^{(\gamma)}(\pi) = \frac{(\widehat{\Delta}_t(\pi))^\gamma}{\overline{\text{IG}}_t(\pi)} = 2 \cdot \frac{(\sum_{a \in \mathcal{A}} \pi(a) \int_{\Theta} \langle \theta, a^*(\theta) - a \rangle dQ_t^+(\theta))^\gamma}{\sum_{a \in \mathcal{A}} \pi(a) \int_{\Theta} (\langle \theta - \bar{\theta}(Q_t^+), a \rangle)^2 dQ_t^+(\theta)}. \quad (3)$$

117 We can at last define our algorithm: Sparse Optimistic Information Directed Sampling (SOIDS). In  
118 each round  $t$ , the policy played by SOIDS is defined to be the distribution on  $\mathcal{A}$  that minimizes the  
119 2-information ratio:

$$\pi_t^{(\text{SOIDS})} = \arg \min_{\pi \in \Delta(\mathcal{A})} \overline{\text{IR}}_t^{(2)}(\pi). \quad (4)$$

120 The choice of  $\gamma = 2$  is motivated by the fact that the minimizer of the 2-information ratio is an  
121 approximate minimizer of surrogate generalized information ratio for all  $\gamma \geq 2$ .

<sup>1</sup>The optimal actions in the hard instance used to prove Theorem 5 in [Jang et al. \[2022\]](#) are  $2s$ -sparse, which still allows us to prove the same bound on the surrogate 3-information ratio, up to constant factors.

122 **Lemma 1.** For all  $\gamma \geq 2$ ,

$$\overline{IR}_t^{(\gamma)}(\pi_t^{(\text{SOIDS})}) \leq 2^{\lambda-2} \min_{\pi \in \Delta(\mathcal{A})} \overline{IR}_t^{(\gamma)}(\pi).$$

123 This fact was discovered for the Bayesian IDS policy by [Lattimore and György \[2021\]](#). We provide  
 124 a proof in Appendix F.2 for completeness. Finally, we remark that the "sparse" part of the name  
 125 SOIDS refers to the choice of the prior  $Q_1^+$ . We use the subset selection prior from Section 3 of  
 126 [Alquier and Lounici \[2011\]](#), which is described in Appendix B.2.

## 127 4 Main results

128 In this section, we state our main results. First, we relate the true regret of any policy sequence to  
 129 the surrogate regret of the same policy sequence. In combination with Lemma 1 and the fact that the  
 130 surrogate regret is controlled by both the 2 and 3-information ratios, this allows us to show that with  
 131 properly tuned parameters, SOIDS has optimal worst-case regret in both the data-poor and data-rich  
 132 regimes. Finally, we show that SOIDS can be tuned in a data-dependent manner, such that its regret  
 133 bound scales with the cumulative observed information ratio instead of the time horizon. The strong  
 134 empirical performance of SOIDS is demonstrated in Appendix J.

### 135 4.1 General bound for the optimistic posterior

136 We start with a generic worst-case regret bound relating the true regret of any algorithm to its sur-  
 137 rogate regret. Since the surrogate regret is defined with respect to the optimistic posterior, which  
 138 is known to the learner, it can be controlled with standard Bayesian techniques. This result is an  
 139 extension of the bounds stated in [Neu et al. \[2024\]](#), [Zhang \[2022\]](#). To our knowledge it is the first  
 140 result of its kind which is compatible with time-dependent or data-dependent learning rates. The  
 141 stated result is specialized to the setting of sparse linear bandits, but the techniques used to deal with  
 142 time-dependent and data-dependent learning rates are applicable beyond this setting.

143 **Theorem 1.** Assume that the optimistic posterior is computed with  $\eta = \frac{1}{4}$  and a sequence of de-  
 144 creasing learning rates  $\lambda_t$  satisfying  $\forall t \geq 1, \lambda_t \leq \frac{1}{2}$ . Set  $\lambda_0 = \frac{1}{2}$ . If the learning rates do not  
 145 depend on the history, then the regret of any sequence of policies  $\pi_t$  satisfies

$$R_T \leq \mathbb{E} \left[ \frac{5+2s \log \frac{edT}{s}}{\lambda_{T-1}} - \sum_{t=1}^T \frac{3}{32} \cdot \frac{\overline{IG}_t(\pi_t)}{\lambda_{t-1}} + 2 \sum_{t=1}^T \widehat{\Delta}_t(\pi_t) \right]. \quad (5)$$

146 Otherwise, if the learning rates depend on the history, let  $C_{1,T}$  be a deterministic upper bound on  
 147  $\frac{1}{\lambda_t} - \frac{1}{\lambda_{t-1}}$  valid for all  $t \leq T$ , and  $C_{2,T}$  be a deterministic upper bound on  $\frac{1}{\lambda_{T-1}}$ . The regret of any  
 148 sequence of policies  $\pi_t$  satisfies

$$R_T \leq \mathbb{E} \left[ \frac{2+s \log \frac{4e^3 d^2 T^3 C_{1,T}^2 C_{2,T}}{s^2}}{\lambda_{T-1}} - \sum_{t=1}^T \frac{3}{32} \cdot \frac{\overline{IG}_t(\pi_t)}{\lambda_{t-1}} + 2 \sum_{t=1}^T \widehat{\Delta}_t(\pi_t) \right] + 2. \quad (6)$$

### 149 4.2 Best of both worlds guarantees for Sparse Optimistic Information Directed Sampling

150 Next, we show that the SOIDS algorithm with properly tuned parameters achieves the optimal rate  
 151 in both the data-rich and data-poor regimes.

152 **Theorem 2.** Assume that our problem satisfies the spare optimal action condition described in

153 *definition 1*. Let  $\lambda_t^{(2)} = \sqrt{\frac{3C_{t+1}}{128d(t+1)}}$  and  $\lambda_t^{(3)} = \frac{1}{4 \cdot 6^{\frac{1}{3}}} \left( \frac{C_{t+1} \sqrt{C_{\min}}}{(t+1)\sqrt{s}} \right)^{\frac{2}{3}}$ , with  $C_t = 5 + 2s \log \frac{edT}{s}$ .

154 Now, set  $\lambda_t = \min(\frac{1}{2}, \max(\lambda_t^{(2)}, \lambda_t^{(3)}))$ , then the regret of SOIDS run with parameter  $\lambda_t$  is upper  
 155 bounded by

$$\begin{aligned} R_T &\leq \min \left( 27 \sqrt{(5 + 2s \log \frac{edT}{s}) dT}, 30 \left( 5 + 2s \log \frac{edT}{s} \right)^{\frac{1}{3}} \left( \frac{T\sqrt{s}}{\sqrt{C_{\min}}} \right)^{\frac{2}{3}} \right) + \mathcal{O}(\sqrt{s} \log \frac{d}{s}) \quad (7) \\ &= \min \left( \mathcal{O} \left( \sqrt{sdT \log \frac{edT}{s}} \right), \mathcal{O} \left( (sT)^{\frac{2}{3}} \left( \log \frac{edT}{s} \right)^{\frac{1}{3}} \right) \right), \end{aligned}$$

156 where  $\mathcal{O}(\sqrt{s} \log \frac{d}{s})$  represents an absolute constant independent of  $T$ .

157 We observe that our algorithm enjoys both the  $\tilde{O}(\sqrt{sdT})$  and the  $\tilde{O}(s^{\frac{2}{3}}T^{\frac{2}{3}})$  rates. Unlike the  
 158 Bayesian regret bound for the sparse IDS algorithm of [Hao et al. \[2021\]](#), our regret bound holds  
 159 in a “worst-case” sense for any value of  $\theta_0 \in \Theta$ . To our knowledge, this makes our method the first  
 160 algorithm to achieve optimal worst-case regret in both regimes simultaneously.

### 161 4.3 Instance dependent guarantees

162 The bounds presented in the previous sections are minimax in nature, meaning they hold uniformly  
 163 over all problem instances. We present a bound in which the scaling with respect to the horizon  
 164  $T$  is replaced with the cumulative surrogate-information ratio. Those quantities are always upper  
 165 bounded by Lemma 7 but could be much smaller in “easier” instances leading to better guarantees.

166 **Theorem 3.** *Assume that our problem satisfies the sparse optimal action condition described in def-*

167 *inition 1 and that  $s \leq \frac{d}{2}$ . Let  $\lambda_t^{(2)} = \sqrt{\frac{s}{2d + \sum_{s=1}^t \overline{IR}_s^{(2)}(\pi_s)}}$  and  $\lambda_t^{(3)} = \left( \frac{s}{\frac{3\sqrt{6}s}{\sqrt{C_{\min}}} + \sum_{s=1}^t \sqrt{\overline{IR}_s^{(3)}(\pi_s)}} \right)^{\frac{2}{3}}$ .*

168 *Then the regret of SOIDS run with parameter  $\lambda_t = \max(\lambda_t^{(3)}, \lambda_t^{(2)})$  satisfies the following regret*  
 169 *bound*

$$\begin{aligned}
 R_T &\leq \left( \frac{2}{s} + \frac{80}{3} + 5 \log \frac{edT}{s} \right) \min \left( \sqrt{s \left( 2d + \sum_{t=1}^{T-1} \overline{IR}_t^{(2)}(\pi_t) \right)}, s^{\frac{1}{3}} \left( \frac{3\sqrt{6}s}{\sqrt{C_{\min}}} + \sum_{t=1}^T \sqrt{\overline{IR}_t^{(3)}(\pi_t)} \right)^{\frac{2}{3}} \right) \\
 &= \mathcal{O} \left( \log \frac{edT}{s} \min \left( \sqrt{s \left( d + \sum_{t=1}^{T-1} \overline{IR}_t^{(2)}(\pi_t) \right)}, s^{\frac{1}{3}} \left( \frac{s}{\sqrt{C_{\min}}} + \sum_{t=1}^T \sqrt{\overline{IR}_t^{(3)}(\pi_t)} \right)^{\frac{2}{3}} \right) \right).
 \end{aligned} \tag{8}$$

170 History dependent learning rates can be used with our novel analysis, making this type of result  
 171 possible. A full proof of that result is provided in Appendix D. Note that this means that our algo-  
 172 rithm is fully adaptive to which of the two regimes is best. Because our analysis requires decreasing  
 173 learning rates, we are forced to leave the  $\log(t)$  terms out of the learning rates and our logarithmic  
 174 term has a worse power than in the bound of Theorem 2. An interesting open question is whether  
 175 it is possible to improve the dependency on this logarithmic term while still using data-dependent  
 176 learning rates.

## 177 5 Analysis

### 178 5.1 Proof of Theorem 1

179 A key observation is that the optimistic posterior can be interpreted as a learner playing an auxiliary  
 180 online learning game over distributions  $\Delta(\Theta)$ . The loss of that game is a weighted sum of neg-  
 181 ative log-likelihood and estimation error losses. We define  $L_t^{(1)}(\theta) = \sum_{s=1}^t \log \left( \frac{1}{p(Y_s|\theta, A_s)} \right) =$   
 182  $\sum_{s=1}^t \frac{1}{2} (\langle \theta, A_s \rangle - Y_s)^2$  to be the *cumulative negative log-likelihood loss* of  $\theta$  and  $L_t^{(2)}(\theta) =$   
 183  $\sum_{s=1}^t -\Delta(A_s, \theta)$  to be the *cumulative estimation error loss* of  $\theta$ . In addition, we define the regular-  
 184 izer  $\Phi : \Delta(\Theta) \rightarrow \mathbb{R}$  by the mapping  $P \mapsto \mathcal{D}_{\text{KL}}(P \| Q_1^+)$ , which is the KL-divergence with respect  
 185 to the prior  $Q_1^+$ . With those notations, the optimistic posterior can be understood as an instance  
 186 of the Follow the Regularized Leader (FTRL) algorithm introduced by [Hazan and Kale \[2010\]](#) and  
 187 [Abernethy et al. \[2008\]](#). In particular, the optimistic posterior can be written as

$$Q_{t+1}^+ = \arg \min_{P \in \Delta(\Theta)} \langle P, \eta L_t^{(1)} + \lambda_t L_t^{(2)} \rangle + \Phi(P).$$

188 This formulation enables the application of tools from convex analysis and online learning, such as  
 189 Fenchel duality, to derive regret bounds for this auxiliary online learning game and to understand  
 190 the interplay between the two losses under the learning rates  $\eta$  and  $\lambda_t$ . Here, we focus on the case in  
 191 which the learning rates  $\lambda_t$  are history-independent, and relegate the analysis of history-dependent  
 192 learning rates to Appendix C. The following lemma provides a bound on the regret under an arbitrary  
 193 comparator distribution  $P$ .

194 **Lemma 2.** *Let  $P \in \Delta(\Theta)$  be any comparator, then the following bound holds*

$$\sum_{t=1}^T \Delta(P, A_t) \leq \frac{\mathcal{D}_{\text{KL}}(P \| Q_1^+)}{\lambda_T} + \frac{\Phi^*(\eta(L_T^{(1)}(\theta_T) - L_T^{(1)}(\cdot)) - \lambda_T L_T^{(2)}(\cdot))}{\lambda_T} + \frac{\eta}{\lambda_T} (P \cdot L_T^{(1)} - L_T^{(1)}(\theta_T)).$$

195 Here  $\theta_t = \arg \min_{\theta \in \Theta} L_t^{(1)}(\theta)$  is the maximum likelihood estimator at time  $t$ , and  $\Phi^*(L) =$   
 196  $\log \int_{\Theta} \exp(L(\theta)) dQ_1^+(\theta)$  is the Fenchel dual of the regularizer  $\Phi$ . A proof is provided in ap-  
 197 pendix B.1.1. We aim to choose the comparator  $P$  and the prior  $Q_1^+$  such that  $P$  is concentrated  
 198 around  $\theta_0$  and  $\mathcal{D}_{\text{KL}}(P \| Q_1^+)$  is small. To achieve this, we exploit the sparsity of  $\theta_0$ . We choose  $Q_1^+$   
 199 to be a subset-selection prior and  $P$  to be the uniform distribution on a sparse neighborhood of  $\theta_0$ .

200 **Lemma 3.** *The subset-selection prior  $Q_1^+ \in \Delta(\Theta)$  verifies that for any  $\epsilon > 0$  and  $\theta \in \Theta$ , there is a*  
 201 *comparator  $P(\theta) \in \Delta(\Theta)$  satisfying both*

$$\forall \theta' \in \text{supp}(P(\theta)), \|\theta - \theta'\|_1 \leq \epsilon \quad \text{and} \quad \mathcal{D}_{\text{KL}}(P(\theta) \| Q_1^+) \leq s \log \frac{2ed}{\epsilon s}.$$

202 The proof of this lemma, as well as the exact choice of the prior  $Q_1^+$  and the comparator  $P(\theta_0)$ ,  
 203 are provided in Appendix B.2. In Appendix I (cf. Lemma 21), we establish that both  $L_T^{(2)}(\cdot)$  and  
 204  $\mathbb{E} \left[ L_T^{(1)}(\cdot) \right]$  are  $2T$ -Lipschitz with respect to the  $\ell_1$ -norm. Hence,

$$\mathbb{E} \left[ \frac{|P \cdot L_T^{(1)} - L_T^{(1)}(\theta_0)|}{\lambda_T} \right] \leq \frac{2T\epsilon}{\lambda_T}, \quad \text{and} \quad \sum_{t=1}^T |\Delta(\theta_0, A_t) - \Delta(P, A_t)| \leq 2T\epsilon.$$

205 In combination with Lemma 2, we obtain the following bound on the cumulative regret:

$$R_T \leq \mathbb{E} \left[ \frac{s \log \frac{2ed}{\epsilon s} + 2T(\lambda_T + \eta)\epsilon}{\lambda_T} + \frac{\Phi^*(-\eta(L_T^{(1)}(\cdot) - L_T^{(1)}(\theta_T)) - \lambda_T L_T^{(2)}(\cdot))}{\lambda_T} \right] \\ + \mathbb{E} \left[ \frac{\eta}{\lambda_T} (L_T^{(1)}(\theta_0) - L_T^{(1)}(\theta_T)) \right].$$

206 The first term balances model complexity and approximation via  $\epsilon$ . In the usual FTRL analysis,  
 207  $\lambda \rightarrow \frac{\phi^*(\lambda L)}{\lambda}$  is non decreasing for any  $L \in \mathbb{R}^\Theta$ , and the term involving  $\Phi^*$  can be telescoped.  
 208 Things are more complex here because only part of the loss is weighted by the time varying learning  
 209 rate  $\lambda_T$ . Through a careful analysis involving the maximum likelihood estimator, we can decompose  
 210 the  $\Phi^*$  term into a telescoping sum and a remainder term.

**Lemma 4.**

$$\frac{\Phi^*(\eta(L_T^{(1)}(\theta_T) - L_T^{(1)}(\cdot)) - \lambda_T L_T^{(2)}(\cdot))}{\lambda_T} \\ \leq \mathbb{E} \left[ \sum_{t=1}^T \frac{\Phi^*(\eta(L_t^{(1)}(\theta_0) - L_t^{(1)}(\cdot)) - \lambda_{t-1} L_t^{(2)}(\cdot))}{\lambda_{t-1}} - \frac{\Phi^*(\eta(L_{t-1}^{(1)}(\theta_0) - L_{t-1}^{(1)}(\cdot)) - \lambda_{t-1} L_{t-1}^{(2)}(\cdot))}{\lambda_{t-1}} \right] \quad (9) \\ + \frac{\eta(6 + s \log \frac{edT}{s})}{\lambda_T}. \quad (10)$$

211 A detailed proof of this result is provided in Appendix B.1.4. Finally, the remaining sum can be  
 212 handled by looking at the explicit formula for  $\Phi^*$ . The terms related to the likelihood and the gap  
 213 estimates can be separated using Hölder's inequality, as is done in Zhang [2022] and Neu, Papini,  
 214 and Schwartz [2024]. More explicitly, by choosing  $\eta = \frac{1}{4}$ , we obtain the following result.

**Lemma 5.**

$$\mathbb{E} \left[ \sum_{t=1}^T \frac{\Phi^*(\eta(L_t^{(1)}(\theta_0) - L_t^{(1)}(\cdot)) - \lambda_{t-1} L_t^{(2)}(\cdot))}{\lambda_{t-1}} - \frac{\Phi^*(\eta(L_{t-1}^{(1)}(\theta_0) - L_{t-1}^{(1)}(\cdot)) - \lambda_{t-1} L_{t-1}^{(2)}(\cdot))}{\lambda_{t-1}} \right] \\ \leq \mathbb{E} \left[ -\sum_{t=1}^T \frac{3\overline{G}_t(\pi_t)}{32\lambda_{t-1}} + 2 \sum_{t=1}^T \widehat{\Delta}(\pi_t) \right]. \quad (11)$$

215 A full proof of this result is provided in Appendix B.1.4. Combining Lemmas 2, 3, 4 and 5, and  
 216 setting  $\epsilon = \frac{2}{T}$ , we obtain the desired regret bound stated in Theorem 1.

## 217 5.2 Proof of Theorem 2

218 We show how Theorem 1 can be combined with bounds on the surrogate regret to control the true  
 219 regret. The first important fact is that the surrogate regret of any policy can always be controlled in  
 220 terms of the 2 or the 3-surrogate information ratio of that policy.

221 **Lemma 6.** Let  $\lambda > 0$ , then we have that for any policy  $\pi \in \Delta(\mathcal{A})$

$$\widehat{\Delta}_t(\pi) \leq \frac{\overline{G}_t(\pi)}{\lambda} + \min \left( \frac{1}{4} \lambda \overline{IR}_t^{(2)}(\pi), c_3^* \sqrt{\lambda \overline{IR}_t^{(3)}(\pi)} \right),$$

222 where  $c_3^* < 2$  is an absolute constant defined in Lemma 27.

223 This result is a consequence of a simple generalization of the AM-GM inequality, and is proved in  
 224 Appendix F.1. Combining this lemma with  $\lambda = \frac{6^d}{3} \lambda_{t-1}$  and Theorem 1, we can further upper bound  
 225 the regret of a sequence of policies  $\pi_t$  as

$$\begin{aligned} R_T &\leq \mathbb{E} \left[ \frac{5+2s \log \frac{edT}{s}}{\lambda_{T-1}} - \sum_{t=1}^T \frac{3\overline{G}_t(\pi_t)}{32\lambda_{t-1}} + 2 \sum_{t=1}^T \widehat{\Delta}_t(\pi_t) \right] \\ &\leq \mathbb{E} \left[ \frac{C_T}{\lambda_{T-1}} + \sum_{t=1}^T \min \left( \frac{32}{3} \lambda_{t-1} \overline{IR}_t^{(2)}(\pi_t), \frac{16}{3} c_3^* \sqrt{3\lambda_{t-1} \overline{IR}_t^{(3)}(\pi_t)} \right) \right], \end{aligned} \quad (12)$$

226 where  $C_T = 5 + 2s \log \frac{edT}{s}$ . Usually, bounds on the 2-information ratio can be converted to  $\mathcal{O}(\sqrt{T})$   
 227 bounds and bounds on the 3-information ratio can be converted to  $\mathcal{O}(T^{\frac{2}{3}})$  bounds. Hence we will  
 228 use the 2-information ratio to control the regret in the data-rich regime and the 3-information ratio  
 229 to control the regret in the data-poor regime. Due to Lemma 1, the SOIDS policy minimizes the  
 230 2-information ratio and approximately minimizes the 3-information ratio. As a result, if there exists  
 231 a "forerunner" algorithm with bounded 2-information ratio or 3-information ratio, SOIDS inherits  
 232 these bounds automatically. In particular, we can use a different forerunner for each regime and  
 233 SOIDS will match the regret guarantees of the best forerunner in each regime. The forerunner  
 234 we consider for the 2-information ratio is the *Feel-Good Thompson Sampling* (FGTS) algorithm  
 235 of Zhang [2022]. For the 3-information ratio, we consider a mixture of the FGTS policy and an  
 236 exploratory policy. The following lemma provides bounds on the surrogate information ratios of the  
 237 SOIDS algorithm.

238 **Lemma 7.** The 2- and 3-surrogate-information ratio of the SOIDS algorithm satisfy for any  $t \geq 0$

$$\overline{IR}_t^{(2)}(\pi_t^{(\text{SOIDS})}) \leq \overline{IR}_t^{(2)}(\pi_t^{(\text{FGTS})}) \leq 2d, \quad (13)$$

239 and

$$\overline{IR}_t^{(3)}(\pi_t^{(\text{SOIDS})}) \leq 2\overline{IR}_t^{(3)}(\pi_t^{(\text{mix})}) \leq \frac{54s}{C_{\min}}. \quad (14)$$

240 The explicit definition of both forerunner algorithms as well as the proof of this lemma are deferred  
 241 to appendix F.3. Note that this bound on the 3-information ratio is the only part of our analysis in  
 242 which the sparse optimal action condition (cf. Definition 1) is required. Finally, we must pick the  
 243 learning rate  $\lambda_t$ . The following lemma describes the appropriate learning rate for the data-poor and  
 244 the data-rich regimes separately.

245 **Lemma 8.** The learning rates  $\lambda_t^{(2)} = \sqrt{\frac{3C_{t+1}}{128d(t+1)}}$  and  $\lambda_t^{(3)} = \frac{1}{4 \cdot 6^{\frac{1}{3}}} \left( \frac{C_{t+1} \sqrt{C_{\min}}}{(t+1)\sqrt{s}} \right)^{\frac{2}{3}}$  guarantee

$$\frac{C_T}{\lambda_{T-1}^{(2)}} + \frac{32}{3} \sum_{t=1}^T \lambda_{t-1}^{(2)} \overline{IR}_t^{(2)}(\pi_t) \leq 16 \sqrt{\frac{2}{3}} C_T d T,$$

246 and

$$\frac{C_T}{\lambda_{T-1}^{(3)}} + \frac{16}{3} c_3^* \sum_{t=1}^T \sqrt{3\lambda_{t-1}^{(3)} \overline{IR}_t^{(3)}(\pi_t)} \leq 12 \cdot 6^{\frac{1}{3}} (C_T)^{\frac{1}{3}} \left( \frac{T\sqrt{s}}{\sqrt{C_{\min}}} \right)^{\frac{2}{3}}.$$

247 The proof is deferred to appendix G.2. It remains to analyze what happens when the learning rate  
 248  $\lambda_t = \min(\frac{1}{2}, \max(\lambda_t^{(2)}, \lambda_t^{(3)}))$  is chosen. We defer this to appendix G.4.

## 249 6 Conclusion

250 There remain several interesting questions that our work leaves open for future research. In our  
 251 experiments, we have made use of an approximate implementation of OIDS adapted from Hao et al.  
 252 [2021]. The initial success we have seen in our experiments suggests that this approximation might  
 253 be viable in more challenging settings, and worthy of an attempt at a solid theoretical analysis. More  
 254 broadly, the results indicate a potential advantage of IDS-style methods over DEC-inspired methods  
 255 [Foster et al., 2022b, Kirschner et al., 2023]. Indeed, we are not aware of any general methods for  
 256 approximating the optimization problems that the E2D algorithm of Foster et al. [2022b] requires to  
 257 solve, in contrast to our results that indicate that IDS-inspired algorithms may very well be amenable  
 258 to practical implementation. Whether the concrete approximation we used in our experiments is the  
 259 best possible one or not remains to be seen.

## 260 References

- 261 Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic  
262 bandits. *Advances in neural information processing systems*, 24, 2011.
- 263 Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Online-to-confidence-set conversions and  
264 application to sparse stochastic bandits. volume 22 of *JMLR Proceedings*, pages 1–9, 2012. URL  
265 <http://proceedings.mlr.press/v22/abbasi-yadkori12.html>.
- 266 Naoki Abe and Philip M Long. Associative reinforcement learning using linear probabilistic con-  
267 cepts. In *ICML*, pages 3–11. Citeseer, 1999.
- 268 Jacob D. Abernethy, Elad Hazan, and Alexander Rakhlin. Competing in the dark: An efficient  
269 algorithm for bandit linear optimization. pages 263–274, 2008. URL <http://colt2008.cs.helsinki.fi/papers/127-Abernethy.pdf>.
- 271 Pierre Alquier and Karim Lounici. Pac-bayesian theorems for sparse regression estimation with  
272 exponential weights. *Electronic Journal of Statistics*, 5:127–145, 2011.
- 273 Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine*  
274 *Learning Research*, 3(Nov):397–422, 2002.
- 275 Hamsa Bastani and Mohsen Bayati. Online decision making with high-dimensional covariates.  
276 *Operations Research*, 68(1):276–294, 2020.
- 277 Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration Inequalities - A*  
278 *Nonasymptotic Theory of Independence*. 2013. ISBN 978-0-19-953525-5. doi: 10.1093/  
279 ACPROF:OSO/9780199535255.001.0001. URL [https://doi.org/10.1093/acprof:oso/](https://doi.org/10.1093/acprof:oso/9780199535255.001.0001)  
280 [9780199535255.001.0001](https://doi.org/10.1093/acprof:oso/9780199535255.001.0001).
- 281 Sébastien Bubeck and Mark Sellke. First-Order Bayesian Regret Analysis of Thompson Sampling,  
282 2022. URL <http://arxiv.org/abs/1902.00681>.
- 283 Sunrit Chakraborty, Saptarshi Roy, and Ambuj Tewari. Thompson sampling for high-dimensional  
284 sparse linear contextual bandits. In *International Conference on Machine Learning*, pages  
285 3979–4008. PMLR, 2023.
- 286 Wei Chu, Lihong Li, Lev Reyzin, and Robert E. Schapire. Contextual bandits with linear payoff  
287 functions. volume 15 of *JMLR Proceedings*, pages 208–214, 2011. URL <http://proceedings.mlr.press/v15/chu11a/chu11a.pdf>.
- 289 Eugenio Clerico, Hamish Flynn, Wojciech Kotowski, and Gergely Neu. Confidence sequences for  
290 generalized linear models via regret analysis, 2025. URL [https://arxiv.org/abs/2504.](https://arxiv.org/abs/2504.16555)  
291 [16555](https://arxiv.org/abs/2504.16555).
- 292 Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit  
293 feedback. In *COLT*, volume 2, page 3, 2008.
- 294 Dylan J. Foster and Alexander Rakhlin. Beyond UCB: Optimal and Efficient Contextual Bandits  
295 with Regression Oracles, 2020. URL <http://arxiv.org/abs/2002.04926>.
- 296 Dylan J. Foster, Noah Golowich, Jian Qian, Alexander Rakhlin, and Ayush Sekhari. A Note on  
297 Model-Free Reinforcement Learning with the Decision-Estimation Coefficient, 2022a. URL  
298 <http://arxiv.org/abs/2211.14250>.
- 299 Dylan J. Foster, Sham M. Kakade, Jian Qian, and Alexander Rakhlin. The Statistical Complexity of  
300 Interactive Decision Making, 2022b. URL <http://arxiv.org/abs/2112.13487>.
- 301 Dylan J. Foster, Alexander Rakhlin, Ayush Sekhari, and Karthik Sridharan. On the Complexity of  
302 Adversarial Decision Making, 2022c. URL <http://arxiv.org/abs/2206.13063>.
- 303 Sébastien Gerchinovitz. Sparsity regret bounds for individual sequences in online linear regression.  
304 *The Journal of Machine Learning Research*, 14(1):729–769, 2013.

- 305 Botao Hao and Tor Lattimore. Regret bounds for information-directed reinforcement  
306 learning. 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/  
307 b733cdd80ed2ae7e3156d8c33108c5d5-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/b733cdd80ed2ae7e3156d8c33108c5d5-Abstract-Conference.html).
- 308 Botao Hao, Tor Lattimore, and Mengdi Wang. High-dimensional sparse linear band-  
309 dits. 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/  
310 7a006957be65e608e863301eb98e1808-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/7a006957be65e608e863301eb98e1808-Abstract.html).
- 311 Botao Hao, Tor Lattimore, and Wei Deng. Information directed sampling for sparse linear bandits.  
312 pages 16738–16750, 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/  
313 8ba6c657b03fc7c8dd4dff8e45defcd2-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/8ba6c657b03fc7c8dd4dff8e45defcd2-Abstract.html).
- 314 Botao Hao, Tor Lattimore, and Chao Qin. Contextual Information-Directed Sampling, 2022. URL  
315 <http://arxiv.org/abs/2205.10895>.
- 316 Elad Hazan and Satyen Kale. Extracting certainty from uncertainty: Regret bounded by variation  
317 in costs. *Machine Learning*, 80(2-3):165–188, 2010. doi: 10.1007/S10994-010-5175-X. URL  
318 <https://doi.org/10.1007/s10994-010-5175-x>.
- 319 Kyoungseok Jang, Chicheng Zhang, and Kwang-Sung Jun. Popart: Efficient sparse regression and  
320 experimental design for optimal sparse linear bandits. *Advances in Neural Information Processing  
321 Systems*, 35:2102–2114, 2022.
- 322 Gi-Soo Kim and Myunghee Cho Paik. Doubly-robust lasso bandit. *Advances in Neural Information  
323 Processing Systems*, 32, 2019.
- 324 Johannes Kirschner and Andreas Krause. Information Directed Sampling and Bandits with Het-  
325 eroscedastic Noise, 2018. URL <http://arxiv.org/abs/1801.09667>.
- 326 Johannes Kirschner, Tor Lattimore, and Andreas Krause. Information directed sampling for linear  
327 partial monitoring. volume 125 of *Proceedings of Machine Learning Research*, pages 2328–2369,  
328 2020. URL <http://proceedings.mlr.press/v125/kirschner20a.html>.
- 329 Johannes Kirschner, Tor Lattimore, Claire Vernade, and Csaba Szepesvári. Asymptotically optimal  
330 information-directed sampling. volume 134 of *Proceedings of Machine Learning Research*, pages  
331 2777–2821, 2021. URL <http://proceedings.mlr.press/v134/kirschner21a.html>.
- 332 Johannes Kirschner, Seyed Alireza Bakhtiari, Kushagra Chandak, Volodymyr  
333 Tkachuk, and Csaba Szepesvári. Regret minimization via saddle point optimiza-  
334 tion. 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/  
335 6eaf8c729af4fbeb18006dc2e6a41d9b-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/6eaf8c729af4fbeb18006dc2e6a41d9b-Abstract-Conference.html).
- 336 Tor Lattimore and András György. Mirror Descent and the Information Ratio. volume 134 of *Pro-  
337 ceedings of Machine Learning Research*, pages 2965–2992, 2021. URL [http://proceedings.  
338 mlr.press/v134/lattimore21b.html](http://proceedings.mlr.press/v134/lattimore21b.html).
- 339 Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- 340 Gergely Neu, Julia Olkhovskaya, Matteo Papini, and Ludovic Schwartz. Lifting the Information  
341 Ratio: An Information-Theoretic Analysis of Thompson Sampling for Contextual Bandits, 2022.  
342 URL <http://arxiv.org/abs/2205.13924>.
- 343 Gergely Neu, Matteo Papini, and Ludovic Schwartz. Optimistic information directed sampling.  
344 volume 247 of *Proceedings of Machine Learning Research*, pages 3970–4006, 2024. URL  
345 <https://proceedings.mlr.press/v247/neu24a.html>.
- 346 Min-hwan Oh, Garud Iyengar, and Assaf Zeevi. Sparsity-agnostic lasso bandit. In *International  
347 Conference on Machine Learning*, pages 8271–8280. PMLR, 2021.
- 348 Francesco Orabona. A modern introduction to online learning. *CoRR*, abs/1912.13213, 2019. URL  
349 <http://arxiv.org/abs/1912.13213>.
- 350 Paat Rusmevichientong and John N Tsitsiklis. Linearly parameterized bandits. *Mathematics of  
351 Operations Research*, 35(2):395–411, 2010.

- 352 Daniel Russo and Benjamin Van Roy. An information-theoretic analysis of thompson sampling.  
353 *Journal of Machine Learning Research*, 17:68:1–68:30, 2016. URL [https://jmlr.org/  
354 papers/v17/14-087.html](https://jmlr.org/papers/v17/14-087.html).
- 355 Daniel Russo and Benjamin Van Roy. Learning to Optimize via Information-Directed Sampling,  
356 2017. URL <http://arxiv.org/abs/1403.5556>.
- 357 Michal Valko, Rémi Munos, Branislav Kveton, and Tomáš Kocák. Spectral bandits for smooth graph  
358 functions. volume 32 of *JMLR Workshop and Conference Proceedings*, pages 46–54, 2014. URL  
359 <http://proceedings.mlr.press/v32/valko14.html>.
- 360 M.J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series  
361 in Statistical and Probabilistic Mathematics. 2019. ISBN 978-1-108-49802-9. URL [https:  
362 //books.google.es/books?id=8C8nuQEACAAJ](https://books.google.es/books?id=8C8nuQEACAAJ).
- 363 Xue Wang, Mingcheng Wei, and Tao Yao. Minimax concave penalized multi-armed bandit model  
364 with high-dimensional covariates. In *International Conference on Machine Learning*, pages  
365 5200–5208. PMLR, 2018.
- 366 Tong Zhang. Feel-good thompson sampling for contextual bandits and reinforcement learning. *SIAM  
367 Journal on Mathematics of Data Science*, 4(2):834–857, 2022. doi: 10.1137/21M140924X. URL  
368 <https://doi.org/10.1137/21m140924x>.

369 **A Related work**

370 The first algorithms and regret bounds for sparse linear bandits were designed for the data-rich  
 371 regime. [Abbasi-Yadkori et al. \[2012\]](#) developed an online-to-confidence-set conversion for linear  
 372 models, which converts any algorithm for online linear regression into a linear bandit algorithm  
 373 whose regret depends on the regret of the online regression algorithm. When the SeqSEW algorithm  
 374 [[Gerchinovitz, 2013](#)] is used in this conversion, the result is a sparse linear bandit algorithm with  
 375 a regret bound of the order  $\mathcal{O}(\sqrt{sdT})$  (ignoring logarithmic factors). [Lattimore and Szepesvári](#)  
 376 [\[2020\]](#) established a matching lower bound for the data-rich regime, showing that this rate cannot  
 377 be improved.

378 More recently, several works have studied the data-poor regime, in which the dimension  $d$  is much  
 379 larger than the number of rounds  $T$ . [Hao et al. \[2020\]](#) showed that an explore-then-commit algorithm  
 380 satisfies a regret bound of the order  $\mathcal{O}((sT)^{2/3}C_{\min}^{-2/3})$ , and established a lower bound of order  
 381  $\Omega(\min(s^{1/3}T^{2/3}C_{\min}^{-1/3}, \sqrt{dT}))$ . Subsequently, [Jang et al. \[2022\]](#) proposed the PopArt estimator for  
 382 sparse linear regression, and showed that an explore-then-commit algorithm that uses this estimator  
 383 achieves a regret bound of the order  $\mathcal{O}(s^{2/3}T^{2/3}H_{\star}^{2/3})$ , where  $H_{\star}$  is another problem-dependent  
 384 quantity that satisfies  $H_{\star}^2 \leq C_{\min}^{-1}$ . In addition, [Jang et al. \[2022\]](#) established a lower bound of order  
 385  $\Omega(s^{2/3}T^{2/3}C_{\min}^{-1/3})$ , showing that the optimal rate for the data-poor regime is  $s^{2/3}T^{2/3}$ . [Hao et al.](#)  
 386 [\[2021\]](#) showed that sparse IDS has a Bayesian best of both worlds/regimes regret bound.

387 A number of works have considered the setting of sparse contextual linear bandits, in which the  
 388 action set  $\mathcal{A}$  changes in each round  $t$ . In the case where the actions sets are chosen by an adaptive  
 389 adversary, the upper and lower bounds of the order  $\sqrt{sdT}$  by [Abbasi-Yadkori et al. \[2012\]](#) and [Lat-](#)  
 390 [timore and Szepesvári \[2020\]](#) respectively still hold. Under the assumption that the action sets are  
 391 generated randomly, and such that either a uniform or greedy policy is (with high probability) ex-  
 392 ploratory, several methods have been shown to achieve nearly dimension-free regret bounds [Bastani](#)  
 393 [and Bayati \[2020\]](#), [Wang et al. \[2018\]](#), [Kim and Paik \[2019\]](#), [Oh et al. \[2021\]](#), [Chakraborty et al.](#)  
 394 [\[2023\]](#).

395 The concept of balancing instantaneous regret and information gain through the information ratio  
 396 was first introduced by [Russo and Roy \[2016\]](#) in the context of analyzing Thompson Sampling.  
 397 Building upon this, the Information-Directed Sampling (IDS) algorithm was proposed by [Russo and](#)  
 398 [Van Roy \[2017\]](#) to directly minimize the information ratio, thereby optimizing the trade-off between  
 399 regret and information gain. These foundational ideas have since been extended and applied to  
 400 a variety of settings including bandits [[Bubeck and Sellke, 2022](#)], contextual bandits [[Neu et al.,](#)  
 401 [2022](#), [Hao et al., 2022](#)], reinforcement learning [[Hao and Lattimore, 2022](#)], and sparse linear bandits  
 402 [[Hao et al., 2021](#)]. However, these works are primarily situated in the Bayesian framework and focus  
 403 on Bayesian regret bounds that hold only in expectation with respect to the prior distribution.

404 A key challenge in extending these methods to the frequentist setting lies in estimating the instanta-  
 405 neous regret and define a meaningful notion of information gain. Both of those things are naturally  
 406 possible in Bayesian analysis but difficult when the true model is unknown. Moreover, Bayesian  
 407 posteriors may inadequately represent model uncertainty from a frequentist perspective. We high-  
 408 light three strands of research that have attempted to address this challenge:

409 Confidence-set based information ratio approaches: Works such as [Kirschner and Krause \[2018\]](#),  
 410 [Kirschner et al. \[2020\]](#), and [Kirschner et al. \[2021\]](#) extend the notion of the information ratio to  
 411 frequentist settings by constructing high-probability confidence sets for the instantaneous regret and  
 412 information gain. These results are mostly limited to setting with some linear structure.

413 Distributionally robust and worst-case information-regret trade-offs: The Decision-to-Estimation-  
 414 Coefficient(DEC) line of work of [[Foster et al., 2022b](#), [Foster and Rakhlin, 2020](#), [Foster et al.,](#)  
 415 [2022c,a](#), [Kirschner et al., 2023](#)] explores the frequentist setting by analyzing worst-case trade-offs  
 416 between regret and information gain. One limitation is that the DEC is an inherently worst-case  
 417 measure of complexity. Moreover, algorithms based on the DEC usually require solving complex  
 418 min-max optimization problems at each time step, making their practical implementation challeng-  
 419 ing and unclear.

420 Optimistic posterior approaches for frequentist guarantees: The approach most closely related to  
 421 our work modifies the Bayesian posterior to provide frequentist guarantees. Introduced by [Zhang](#)  
 422 [\[2022\]](#), the optimistic posterior is a modification of the Bayesian posterior which enables frequentist  
 423 regret bounds for a variant of Thompson Sampling. Subsequently, [Neu et al. \[2024\]](#) studied the

424 optimistic posterior framework in greater depth, defining a frequentist analog of the information  
 425 ratio to extend IDS to frequentist settings. A notable limitation of these works is their restriction to  
 426 constant learning rates in the optimistic posterior, which limits adaptivity, an issue that we address  
 427 in this paper.

## 428 **B Analysis of the Optimistic posterior**

429 This section provides further details about the prior underlying the optimistic posterior and guaran-  
 430 tees on the posterior updates.

### 431 **B.1 Follow the regularized leader analysis**

432 The main step in our analysis of the optimistic posterior is to leverage the follow the regularized  
 433 leader formulation of our optimistic posterior update

$$Q_{t+1}^+ = \arg \min_{P \in \Delta(\Theta)} \langle P, \eta L_t^{(1)} + \lambda_t L_t^{(2)} \rangle + \Phi(P).$$

#### 434 **B.1.1 Proof of lemma 2**

435 As is usual in the analysis of the follow the regularized leader algorithm, we introduce the Fenchel  
 436 conjugate of the regularization function  $\Phi = \mathcal{D}_{\text{KL}}(\cdot \| Q_1^+)$  as the function  $\Phi^* : \mathbb{R}^\Theta \rightarrow \mathbb{R}$  taking  
 437 values  $\Phi^*(L) = \sup_{P \in \Delta(\Theta)} \{\langle P, L \rangle - \Phi(P)\}$ . The Fenchel–Young inequality guarantees that for  
 438 any  $P \in \Delta(\Theta), L \in \mathbb{R}^\Theta$ , we have

$$\langle P, L \rangle \leq \Phi(P) + \Phi^*(L)$$

439 We now introduce the maximum likelihood estimator  $\theta_t = \arg \min_{\theta \in \Theta} L_t^{(1)}(\theta)$  and let  $L =$   
 440  $-\eta(L_T^{(1)}(\cdot) - L_T^{(1)}(\theta_T)) - \lambda_T L_T^{(2)}(\cdot)$ . Since  $\lambda_T$  is never used by the algorithm, we can further  
 441 assume that  $\lambda_T = \lambda_{T-1}$ . The role of the maximum likelihood estimator is to make sure that the  
 442 term  $L_t^{(1)}(\theta) - L_t^{(1)}(\theta_t)$  is always non-negative. Applying Fenchel–Young to  $L$  gives us the follow-  
 443 ing bound:

$$\eta \left( L_T^{(1)}(\theta_T) - \langle P, L_T^{(1)} \rangle \right) - \lambda_T \langle P, L_T^{(2)} \rangle \leq \Phi(P) + \Phi^* \left( -\eta(L_T^{(1)}(\cdot) - L_T^{(1)}(\theta_T)) - \lambda_T L_T^{(2)}(\cdot) \right)$$

444 Noticing that  $\langle P, L_T^{(1)} \rangle = -\sum_{t=1}^T \Delta(P, A_t)$  and rearranging the terms concludes the proof.

#### 445 **B.1.2 Proof of Lemma 4**

446 We start by rewriting the potential function in the form of the following telescopic sum:

$$\begin{aligned} & \frac{\Phi^*(-\eta(L_T^{(1)}(\cdot) - L_T^{(1)}(\theta_T)) - \lambda_T L_T^{(2)}(\cdot))}{\lambda_T} \\ &= \sum_{t=1}^T \frac{\Phi^*(-\eta(L_t^{(1)}(\cdot) - L_t^{(1)}(\theta_t)) - \lambda_t L_t^{(2)}(\cdot))}{\lambda_t} - \frac{\Phi^*(-\eta(L_{t-1}^{(1)}(\cdot) - L_{t-1}^{(1)}(\theta_{t-1})) - \lambda_{t-1} L_{t-1}^{(2)}(\cdot))}{\lambda_{t-1}}. \end{aligned}$$

447 In the usual follow-the-regularized-leader analysis, we use the fact that  $\lambda \rightarrow \frac{\phi^*(\lambda L)}{\lambda}$  is non-  
 448 decreasing for any  $L \in \mathbb{R}^\Theta$ . Here however, only some of the linear loss is scaled by  $\lambda_t$  and the  
 449 usual FTRL analysis fails. Crucially, because we introduced the maximum likelihood estimator  $\theta_t$ ,  
 450 we have that  $L_t^{(1)}(\cdot) - L_t^{(1)}(\theta_t) \geq 0$  and we can instead use the following lemma that guarantees  
 451 that a scaled and shifted dual is monotonous.

452 **Lemma 9.** *Let  $\Phi \geq 0, \Phi^*$  be a convex function and its dual as defined previously,  $L_1, L_2 \in \mathbb{R}^\Theta$   
 453 with  $L_1 \geq 0$ , then  $\lambda \in \mathbb{R}^{+*} \rightarrow \frac{\Phi^*(-L_1 + \lambda L_2)}{\lambda}$  is a non-decreasing function.*

454 *Proof.* By definition, we have

$$\begin{aligned} \frac{\Phi^*(-L_1 + \lambda L_2)}{\lambda} &= \frac{\sup_{P \in \Delta(\Theta)} \langle P, -L_1 + \lambda L_2 \rangle - \Phi(P)}{\lambda} \\ &= \sup_{P \in \Delta(\Theta)} \langle P, L_2 \rangle - \frac{\langle P, L_1 \rangle + \Phi(P)}{\lambda}. \end{aligned}$$

455 For any  $P \in \Delta(\Theta)$ , we have that  $\Phi(P) + \langle P, L_1 \rangle \geq 0$  and the term inside the supremum is  
 456 non-decreasing with respect to lambda. Since the supremum of non-decreasing functions is also  
 457 non-decreasing, this concludes the proof.  $\square$

458 Applying the previous lemma, we upper bound the previous sum by replacing each  $\lambda_t$  factor by  
 459  $\lambda_{t-1}$  (using the convention  $\lambda_0 = 1/2$ ), and then we replace the maximum likelihood estimator  $\theta_t$   
 460 by  $\theta_0$  inside  $\Phi^*$  to obtain

$$\begin{aligned} & \sum_{t=1}^T \frac{\Phi^*(-\eta(L_t^{(1)}(\cdot) - L_t^{(1)}(\theta_t)) - \lambda_t L_t^{(2)}(\cdot))}{\lambda_t} - \frac{\Phi^*(-\eta(L_{t-1}^{(1)}(\cdot) - L_{t-1}^{(1)}(\theta_{t-1})) - \lambda_{t-1} L_{t-1}^{(2)}(\cdot))}{\lambda_{t-1}} \\ & \leq \sum_{t=1}^T \frac{\Phi^*(-\eta(L_t^{(1)}(\cdot) - L_t^{(1)}(\theta_t)) - \lambda_t L_t^{(2)}(\cdot))}{\lambda_{t-1}} - \frac{\Phi^*(-\eta(L_{t-1}^{(1)}(\cdot) - L_{t-1}^{(1)}(\theta_{t-1})) - \lambda_{t-1} L_{t-1}^{(2)}(\cdot))}{\lambda_{t-1}} \\ & = \sum_{t=1}^T \frac{\Phi^*(-\eta(L_t^{(1)}(\cdot) - L_t^{(1)}(\theta_0)) - \lambda_t L_t^{(2)}(\cdot))}{\lambda_{t-1}} - \frac{\Phi^*(-\eta(L_{t-1}^{(1)}(\cdot) - L_{t-1}^{(1)}(\theta_0)) - \lambda_{t-1} L_{t-1}^{(2)}(\cdot))}{\lambda_{t-1}} \\ & + \frac{\eta}{\lambda_{t-1}} (L_t^{(1)}(\theta_t) - L_t^{(1)}(\theta_0) + L_{t-1}^{(1)}(\theta_0) - L_{t-1}^{(1)}(\theta_{t-1})). \end{aligned}$$

461 It remains to bound the difference of the negative log likelihood of the true parameter and the max-  
 462 imum likelihood estimator. This is done via the following result (whose proof we relegate to ap-  
 463 pendix E.1.1).

464 **Lemma 10.** *For any  $t \geq 1$ , we have*

$$0 \leq \mathbb{E} \left[ L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_t) \right] \leq \inf_{\rho} \left\{ 2\rho t + s \log \frac{ed(1+2/\rho)}{s} \right\} \leq 6 + s \log \frac{edt}{s} \quad (15)$$

465 Using this lemma, we can further bound the previously considered expression as the following  
 466 telescopic sum:

$$\begin{aligned} & \mathbb{E} \left[ \sum_{t=1}^T \frac{\eta}{\lambda_{t-1}} (L_t^{(1)}(\theta_t) - L_t^{(1)}(\theta_0) + L_{t-1}^{(1)}(\theta_0) - L_{t-1}^{(1)}(\theta_{t-1})) + \frac{\eta}{\lambda_T} (L_T^{(1)}(\theta_0) - L_T^{(1)}(\theta_T)) \right] \\ & = \mathbb{E} \left[ \sum_{t=1}^T \frac{\eta}{\lambda_{t-1}} (L_t^{(1)}(\theta_t) - L_t^{(1)}(\theta_0)) - \sum_{t=1}^T \frac{\eta}{\lambda_t} (L_t^{(1)}(\theta_t) - L_t^{(1)}(\theta_0)) \right] \\ & \leq \eta \cdot \sum_{t=1}^T \mathbb{E} \left[ (L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_t)) \right] \left( \frac{1}{\lambda_t} - \frac{1}{\lambda_{t-1}} \right) \\ & \leq \frac{\eta(6 + s \log \frac{edT}{s})}{\lambda_T}. \end{aligned}$$

467 Here, the first inequality comes from the non-negativity of  $L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_t)$  by definition of  $\theta_t$   
 468 and the second one is from Lemma 10 just above and a telescoping argument. Finally we obtain the  
 469 claim of Lemma 4.

### 470 B.1.3 Controlling the losses separately

471 The focus of this section is to understand how to control  $\Phi^*(-L)$  where  $L$  is either the negative-  
 472 likelihood loss or the estimation-error loss. We start by analyzing the negative-likelihood loss. As  
 473 was done in Neu, Papini, and Schwartz [2024], we will relate the negative-likelihood loss to the  
 474 surrogate information gain.

475 For this analysis, we define the *true information gain* as

$$\text{IG}_t(\pi) = \frac{1}{2} \sum_{a \in \mathcal{A}} \pi(a) \int_{\Theta} \langle \theta - \theta_0, a \rangle^2 dQ_t^+(\theta), \quad (16)$$

476 and note that, by linearity reward function, the surrogate information gain is always smaller than the  
 477 true information gain. This is stated formally below.

478 **Proposition 1.** For any policy  $\pi \in \Delta(\mathcal{A})$  and any  $t \geq 1$  we have that

$$\overline{IG}_t(\pi) \leq IG_t(\pi) \quad (17)$$

479 The proof is provided in Appendix I.1. This result can then be used to relate the surrogate and the  
480 true information gain to the negative-likelihood loss. This result and its proof are identical to the  
481 proof of Lemma 17 in Neu, Papini, and Schwartz [2024].

482 **Lemma 11.** Assume that the noise  $\epsilon_t$  is conditionally 1-sub-Gaussian, then for any  $t \geq 1, \eta, \alpha \geq 0$   
483 such that  $\gamma = \frac{\eta\alpha}{2} (1 - \eta\alpha) > 0$ , the following inequality holds

$$\mathbb{E} \left[ \log \int_{\Theta} \left( \frac{p(Y_t|\theta, A_t)}{p(Y_t|\theta_0, A_t)} \right)^{\eta\alpha} dQ_t^+(\theta) \right] \leq -2\gamma(1 - 2\gamma) \mathbb{E} [IG_t(\pi_t)] \quad (18)$$

$$\leq -2\gamma(1 - 2\gamma) \mathbb{E} [\overline{IG}_t(\pi_t)]. \quad (19)$$

484 In particular, the constant  $2\gamma(1 - 2\gamma)$  can be maximized to the value  $\frac{3}{16}$  by the choice  $\eta\alpha = \frac{1}{2}$ .

485 *Proof.* By the tower rule of expectation and Jensen's inequality applied to the logarithm, we have

$$\begin{aligned} \mathbb{E} \left[ -\log \int_{\Theta} \left( \frac{p(Y_t|\theta, A_t)}{p(Y_t|\theta_0, A_t)} \right)^{\eta\alpha} \right] &= \mathbb{E} \left[ \mathbb{E} \left[ -\log \int_{\Theta} \left( \frac{p(Y_t|\theta, A_t)}{p(Y_t|\theta_0, A_t)} \right)^{\eta\alpha} \middle| \mathcal{F}_t, A_t \right] \right] \\ &\leq \mathbb{E} \left[ -\log \mathbb{E} \left[ \int_{\Theta} \left( \frac{p(Y_t|\theta, A_t)}{p(Y_t|\theta_0, A_t)} \right)^{\eta\alpha} \middle| \mathcal{F}_t, A_t \right] \right] \\ &= \mathbb{E} \left[ -\log \int_{\Theta} \mathbb{E} \left[ \exp \left( -\eta\alpha \left( \frac{(Y_t - \langle \theta, A_t \rangle)^2}{2} - \frac{(Y_t - \langle \theta_0, A_t \rangle)^2}{2} \right) \right) \middle| \mathcal{F}_t, A_t \right] \right]. \end{aligned}$$

486 Now, we fix some  $\theta \in \Theta$  and to simplify the notation, we let  $r_0 = \langle \theta_0, A_t \rangle$  and  $r = \langle \theta, A_t \rangle$ . Using  
487 some elementary manipulations and the conditional sub-gaussianity of  $\epsilon_t$  and  $Y_t = r_0 + \epsilon_t$  which im-  
488 plies that for any  $(\mathcal{F}_t, A_t)$ -measurable  $\zeta_t$ ,  $\mathbb{E} [\exp(Y_t \zeta_t) | \mathcal{F}_t, A_t] = \exp(r_0 \zeta_t) \mathbb{E} [\exp(\epsilon_t \zeta_t) | \mathcal{F}_t, A_t] \leq$   
489  $\exp(r_0 \zeta_t) \exp\left(\frac{\zeta_t^2}{2}\right)$ , we have

$$\begin{aligned} &\mathbb{E} \left[ \exp \left( -\eta\alpha \left( \frac{(Y_t - r)^2}{2} - \frac{(Y_t - r_0)^2}{2} \right) \right) \middle| \mathcal{F}_t, A_t \right] \\ &= \mathbb{E} \left[ \exp \left( -\frac{\eta\alpha}{2} (2Y_t - r - r_0)(r_0 - r) \right) \middle| \mathcal{F}_t, A_t \right] \\ &= \exp \left( \eta\alpha \frac{r_0^2 - r^2}{2} \right) \mathbb{E} [\exp(\eta\alpha Y_t (r - r_0)) | \mathcal{F}_t, A_t] \\ &\leq \exp \left( \eta\alpha \frac{r_0^2 - r^2}{2} \right) \cdot \exp(\eta\alpha r_0 (r - r_0)) \exp \left( \frac{\eta^2 \alpha^2}{2} (r - r_0)^2 \right) \\ &= \exp \left( -(r - r_0)^2 \cdot \frac{\eta\alpha}{2} (1 - \eta\alpha) \right). \end{aligned}$$

490 Further, defining  $\gamma = \frac{\eta\alpha}{2} (1 - \eta\alpha)$ , we have

$$\begin{aligned} &\mathbb{E} \left[ \exp \left( -\eta\alpha \left( \frac{(Y_t - r)^2}{2} - \frac{(Y_t - r_0)^2}{2} \right) \right) \middle| \mathcal{F}_t, A_t \right] \\ &\leq \exp(-\gamma(r - r_0)^2) \\ &\leq 1 - \gamma(r - r_0)^2 + \frac{\gamma^2}{2}(r - r_0)^4 \\ &\leq 1 - \gamma(r - r_0)^2 + 2\gamma^2(r - r_0)^2 \\ &\leq 1 - \gamma(1 - 2\gamma)(r - r_0)^2. \end{aligned}$$

491 Here, we used the elementary inequality  $\exp(x) \leq 1 + x + \frac{x^2}{2}$  for  $x \leq 0$  and then used  $|r - r_0| \leq 2$ .  
492 Finally, using that  $\log x \leq x - 1$  for any  $x > 0$ , and taking the integral over  $\Theta$ , we get that

$$\begin{aligned} \mathbb{E} \left[ -\log \int_{\Theta} \left( \frac{p(Y_t|\theta, A_t)}{p(Y_t|\theta_0, A_t)} \right)^{\eta\alpha} \right] &\leq -\gamma(1 - 2\gamma) \mathbb{E} \left[ \sum_{a \in \mathcal{A}} \pi_t(A) \int_{\Theta} \langle \theta - \theta_0, a \rangle^2 \right] dQ_t^+(\theta) \\ &= -2\gamma(1 - 2\gamma) \mathbb{E} [IG_t(\pi_t)]. \end{aligned}$$

493 Rearranging and combining the result with Proposition 1 yields the claim of the lemma.  $\square$

494 We now turn our focus to the estimation error loss and relate it to the surrogate regret through the  
 495 following lemma, whose proof is a straightforward application of Lemma 23.

496 **Lemma 12.** *For any  $t \geq 1, \beta > 1$ , if  $\beta\lambda_{t-1} \leq 1$ , we have*

$$\mathbb{E} \left[ \frac{1}{\beta\lambda_{t-1}} \log \int_{\Theta} \exp(\beta\lambda_{t-1}\Delta(a_t, \theta)) dQ_t^+(\theta) \right] \leq \mathbb{E} \left[ 2\widehat{\Delta}_t(\pi_t) \right]. \quad (20)$$

#### 497 B.1.4 Separation of the two losses: proof of Lemma 5

498 We now make use of the fact that the Fenchel dual of  $\Phi$  can be explicitly written as  $\Phi^*(L) =$   
 499  $\log \int_{\Theta} \exp(L(\theta)) dQ_1(\theta)$ . As a result, we have

$$\begin{aligned} & \mathbb{E} \left[ \sum_{t=1}^T \frac{\Phi^*(-\eta(L_t^{(1)}(\cdot) - L_t^{(1)}(\theta_0)) - \lambda_{t-1}L_t^{(2)}(\cdot))}{\lambda_{t-1}} - \frac{\Phi^*(-\eta(L_{t-1}^{(1)}(\cdot) - L_{t-1}^{(1)}(\theta_0)) - \lambda_{t-1}L_{t-1}^{(2)}(\cdot))}{\lambda_{t-1}} \right] \\ &= \mathbb{E} \left[ \sum_{t=1}^T \frac{1}{\lambda_{t-1}} \log \frac{\int_{\Theta} \left( \frac{p(Y_t|\theta, A_t)}{p(Y_t|\theta_0, A_t)} \right)^\eta \exp(\lambda_{t-1}\Delta(A_t, \theta)) \exp(-\eta L_{t-1}^{(1)}(\theta) - \lambda_{t-1}L_{t-1}^{(2)}(\theta)) dQ_1(\theta)}{\int_{\Theta} \exp(-\eta L_{t-1}^{(1)}(\theta) - \lambda_{t-1}L_{t-1}^{(2)}(\theta)) dQ_1(\theta)} \right] \\ &= \mathbb{E} \left[ \sum_{t=1}^T \frac{1}{\lambda_{t-1}} \log \int_{\Theta} \left( \frac{p(Y_t|\theta, A_t)}{p(Y_t|\theta_0, A_t)} \right)^\eta \exp(\lambda_{t-1}\Delta(A_t, \theta)) dQ_t^+(\theta) \right] \\ &\leq \mathbb{E} \left[ \sum_{t=1}^T \frac{1}{\alpha\lambda_{t-1}} \log \int_{\Theta} \left( \frac{p(Y_t|\theta, A_t)}{p(Y_t|\theta_0, A_t)} \right)^{\eta\alpha} + \frac{1}{\beta\lambda_{t-1}} \log \int_{\Theta} \exp(\beta\lambda_{t-1}\Delta(A_t, \theta)) dQ_t^+(\theta) \right], \end{aligned}$$

500 where the last equality is by definition of the optimistic posterior and the last inequality follows from  
 501 using Hölder's inequality with the two real numbers  $\alpha, \beta > 1$  that satisfy  $\frac{1}{\alpha} + \frac{1}{\beta} = 1$ . Combining  
 502 Lemma 11 and Lemma 12 with the choice  $\alpha = \beta = 2$ , the fact that  $\eta = \frac{1}{4}$  and the last inequality  
 503 yields the claim of the Lemma.  $\square$

#### 504 B.2 Choice of the prior and comparator distribution: proof of Lemma 3

505 In order to construct the prior  $Q_1$  and the comparator  $P$  for the regret analysis, we need to take into  
 506 account two criteria: that  $\mathcal{D}_{\text{KL}}(P||Q_1)$  be controlled and that  $|\langle P, L \rangle - L(\theta_0)|$  be small. Note that  
 507 the comparator should be a function of the unknown parameter  $\theta_0$ , and thus we denote it by  $P(\theta_0)$ .  
 508 As for the prior, it should take into account the sparsity level of the unknown  $\theta_0$ , but should have no  
 509 access to its support.

510 For the prior, we first design a distribution  $\Pi$  over the set of all subsets of  $[d] = \{1, \dots, d\}$ , which  
 511 have cardinality at most  $s$ . We choose the distribution such that: a) the probability assigned to each  
 512 subset depends only on its cardinality; b) the probability assigned to the set of all subsets of size  $k$   
 513 is proportional to  $2^{-k}$ , where  $1 \leq k \leq s$ . In other words, we prefer smaller subsets and have no  
 514 preference over which indices in  $[d]$  are included. The distribution that satisfies these requirements is

$$\Pi(S) = \frac{2^{-|S|}}{\binom{d}{|S|} \sum_{k=1}^s 2^{-k}}. \quad (21)$$

515 For  $S = \emptyset$ , we set  $\Pi(S) = 0$ . Doing so only complicates matters if the support of  $\theta_0$  is empty (i.e.,  
 516  $\theta_0 = 0$ ). However, in this case, the reward function is 0 everywhere, which means any algorithm  
 517 would have 0 regret. We therefore continue under the assumption that  $\theta_0 \neq 0$ . The most impor-  
 518 tant property of this distribution, which we will use later, is that for any subset  $S$  of cardinality  $s$ ,  
 519  $\log(1/\Pi(S)) \leq s \log(2ed/s)$ . For each subset  $S$ , we define  $Q_S$  to be the uniform distribution on  
 520  $\Theta_S$ . The prior is defined to be

$$Q_1 = \sum_{S \subset [d]: |S| \leq s} \Pi(S) Q_S.$$

521 As for the comparator distribution  $P(\theta_0)$ , we would ideally like to take a Dirac measure on  $\theta_0$ , but  
 522 this would make the KL divergence appearing in the bound blow up. Thus, we pick a comparator  
 523  $P$  which dilutes its mass around  $\theta_0$ . For any  $\theta \in \Theta$ , with support  $\bar{S}$ , and any  $\epsilon \in (0, 1)$ , we define  
 524 the set  $(1 - \epsilon)\theta + \epsilon\Theta_{\bar{S}} = \{(1 - \epsilon)\theta + \epsilon\theta' : \theta' \in \Theta_{\bar{S}}\} \subset \Theta_{\bar{S}}$ . We will choose  $P$  to be the uniform  
 525 distribution on  $(1 - \epsilon)\theta_0 + \epsilon\Theta_{S_0}$ . We now bound  $\Phi(P) = \mathcal{D}_{\text{KL}}(P||Q_1)$  for this choice of  $P$  in the  
 526 following lemma, from which the claim of Lemma 3 then directly follows.

527 **Lemma 13.** For any  $\bar{\theta} \in \Theta$ , let  $\bar{S}$  denote its support, and let  $|\bar{S}| = s$ . If, for  $\epsilon \in (0, 1)$ ,  $P =$   
528  $\mathcal{U}((1 - \epsilon)\bar{\theta} + \epsilon\Theta_{\bar{S}})$  and  $Q_1 = \sum_{S \subset [d]; |S|=s} \Pi(S)Q_S$ , then  $\mathcal{D}_{\text{KL}}(P||Q_1) \leq s \log \frac{2ed}{\epsilon s}$ .

529 *Proof.* We notice that  $(1 - \epsilon)\bar{\theta} + \epsilon\Theta_{\bar{S}}$  is an  $s$ -dimensional L1 ball of radius  $\epsilon$ , which is contained in  
530  $\Theta_{\bar{S}}$ . Therefore, on the support of  $P$ ,  $\frac{dP}{dQ_{\bar{S}}}$  is equal to the ratio of the volumes of a unit L1 ball and  
531 an L1 ball of radius  $\epsilon$ , which is  $(1/\epsilon)^s$ . Thus,

$$\mathcal{D}_{\text{KL}}(P||Q_1) = \int \log \frac{dP}{\sum_S \Pi(S)dQ_S} dP \leq \int \log \frac{dP}{\Pi(\bar{S})dQ_{\bar{S}}} dP \leq s \log \frac{1}{\epsilon} + \log \frac{1}{\Pi(\bar{S})}.$$

532 Using the definition of  $\Pi$  and the bound  $\binom{d}{s} \leq \left(\frac{ed}{s}\right)^s$  on the binomial coefficient, we have

$$\log \frac{1}{\Pi(\bar{S})} = \log \binom{d}{s} + s \log(2) + \log \sum_{k=1}^s 2^{-k} \leq s \log \frac{2ed}{s}.$$

533 Combining everything, we obtain

$$\mathcal{D}_{\text{KL}}(P||Q_1) \leq s \log \frac{1}{\epsilon} + s \log \frac{2ed}{s} = s \log \frac{2ed}{\epsilon s}, \quad (22)$$

534 as advertised.  $\square$

## 535 C Proof of the history-dependent part of Theorem 1

536 We now focus on the case in which  $\lambda_t$  is allowed to depend on the history. Following the original  
537 analysis, we arrive again at equation 2

$$\Delta(P, a_t) \leq \frac{\mathcal{D}_{\text{KL}}(P||Q_1)}{\lambda_T} + \frac{\Phi^*(-\eta L_T^{(1)}(\cdot) + \eta L_T^{(1)}(\theta_T) + \lambda_T L_T^{(2)}(\cdot))}{\lambda_T} + \frac{\eta}{\lambda_T} (P \cdot L_T^{(1)} - L_T^{(1)}(\theta_T)),$$

538 where  $P \in \Delta(\Theta)$  can be any comparator distribution. Lemma 3 is still valid and we can chose the  
539 same prior as before. We can still choose a comparator distribution supported on an  $\epsilon$ -ball around  $\theta_0$ .

540 However, because  $\lambda_t$  depends on the history, we can no longer upper bound  $\mathbb{E} \left[ \frac{|P \cdot L_T^{(1)} - L_T^{(1)}(\theta_0)|}{\lambda_{T-1}} \right]$

541 by  $\mathbb{E} \left[ \frac{2T\epsilon}{\lambda_T} \right]$ . Using Lemma 21, we still have that  $L_T^{(2)}(\cdot)$  is  $2T$ -Lipschitz and  $\mathbb{E} \left[ L_T^{(1)}(\cdot) \right]$  is  $2T$ -  
542 Lipschitz. Hence,

$$\mathbb{E} \left[ \frac{|P \cdot L_T^{(1)} - L_T^{(1)}(\theta_0)|}{\lambda_{T-1}} \right] \leq 2T\epsilon C_{2,T}, \quad \text{and} \quad \sum_{t=1}^T |\Delta(\theta_0, a_t) - \Delta(P, a_t)| \leq 2T\epsilon,$$

543 where we used  $C_{2,T}$ , a deterministic upper bound on  $\frac{1}{\lambda_{T-1}}$ . Exactly the same telescoping of  $\Phi^*$  can  
544 be done, however because the learning rate is history-dependent, the difference between the negative  
545 log likelihood of  $\theta_0$  and  $\theta_t$  must be treated with more care. We have the following lemma

546 **Lemma 14.** Let  $C_{1,T}$  be a deterministic upper bound on  $\left(\frac{1}{\lambda_{t+1}} - \frac{1}{\lambda_t}\right)$  that holds for all  $t < T$ , then

$$\begin{aligned} & \mathbb{E} \left[ \sum_{t=1}^T \frac{\eta}{\lambda_{t-1}} (L_t^{(1)}(\theta_t) - L_t^{(1)}(\theta_0) + L_{t-1}^{(1)}(\theta_0) - L_{t-1}^{(1)}(\theta_{t-1})) + \frac{\eta}{\lambda_T} (L_T^{(1)}(\theta_0) - L_T^{(1)}(\theta_T)) \right] \\ & \leq \mathbb{E} \left[ \frac{\eta(12 + 3s \log \frac{2e^2 d T^2 C_{1,T}^2}{s})}{2\lambda_{T-1}} \right]. \end{aligned} \quad (23)$$

547 A complete proof of that result can be found in appendix E.2.1.

548 Finally, as was the case in the history independent version the telescoping sum can be handled by  
549 looking at the explicit formula for  $\Phi^*$  and Lemma 5 still holds. Applying Lemma 5 and setting  
550  $\epsilon = \frac{1}{TC_{2,T}}$  yields the claim of the theorem.

551 **D Proof of Theorem 3**

552 We turn our attention to data-dependent bounds (that will scale with the cumulative information  
 553 ratio rather than the time horizon). Combining the second part of Theorem 1 with Lemma 6 and the  
 554 choice  $\lambda = \frac{64}{3}\lambda_{t-1}$ , we have that for any non-increasing sequence of learning rates  $\lambda_t$  satisfying  
 555  $\lambda_0 \leq \frac{1}{2}$ , the following holds

$$R_T \leq \mathbb{E} \left[ \frac{C_T}{\lambda_{T-1}} + \min \left( \sum_{t=1}^T \frac{32}{3} \lambda_{t-1} \overline{\text{IR}}_t^{(2)}(\pi_t), \frac{16}{3} c_3^* \sqrt{3\lambda_{t-1} \overline{\text{IR}}_t^{(3)}(\pi_t)} \right) \right], \quad (24)$$

556 where  $C_T = 2 + s \log \frac{4e^3 d^2 T^3 C_{1,T}^2 C_{2,T}}{s^2}$  and  $C_{1,T}$ , respectively  $C_{2,T}$  are deterministic upper bounds  
 557 on  $\frac{1}{\lambda_t} - \frac{1}{\lambda_{t-1}}$ , respectively  $\frac{1}{\lambda_{T-1}}$ .

558 We let  $\lambda_t^{(2)} = \sqrt{\frac{s}{2d + \sum_{s=1}^t \overline{\text{IR}}_s^{(2)}(\pi_s)}}$  and  $\lambda_t^{(3)} = \left( \frac{s}{\frac{3\sqrt{6}s}{\sqrt{C_{\min}}} + \sum_{s=1}^t \sqrt{\overline{\text{IR}}_s^{(3)}(\pi_s)}} \right)^{\frac{3}{2}}$ , and verify that  $\lambda_t =$   
 559  $\max(\lambda_t^{(2)}, \lambda_t^{(3)})$  is decreasing and always smaller than  $\frac{1}{2}$ . We also verify that  $C_{1,T} = C_{2,T} = \sqrt{\frac{dT}{s}}$   
 560 are valid upper bounds. As a result, we have the following upper bound

$$C_T = 2 + s \log \frac{4e^3 d^2 T^3 C_{1,T}^2 C_{2,T}}{s^2} \leq 2 + s \log 4e^3 T^{4.5} \left( \frac{d}{s} \right)^{3.5} \leq 2 + 5s \log \left( \frac{edT}{s} \right). \quad (25)$$

561 We now focus on bounding the sum containing the information ratios. Applying Lemma 7, we  
 562 obtain that for all  $t \geq 1$ ,  $\overline{\text{IR}}_t^{(2)}(\pi_t) \leq 2d$  and for any  $T \geq 1$

$$\begin{aligned} \sum_{t=1}^T \lambda_{t-1}^{(2)} \overline{\text{IR}}_t^{(2)}(\pi) &= \sqrt{s} \sum_{t=1}^T \frac{\overline{\text{IR}}_t^{(2)}(\pi_t)}{\sqrt{2d + \sum_{s=1}^{t-1} \overline{\text{IR}}_s^{(2)}(\pi_s)}} \\ &\leq \sqrt{s} \sum_{t=1}^T \frac{\overline{\text{IR}}_t^{(2)}(\pi_t)}{\sqrt{\sum_{s=1}^t \overline{\text{IR}}_s^{(2)}(\pi_s)}} \\ &\leq 2 \sqrt{s \sum_{t=1}^T \overline{\text{IR}}_t^{(2)}(\pi_t)} \\ &\leq 2 \sqrt{s \left( 2d + \sum_{t=1}^{T-1} \overline{\text{IR}}_t^{(2)}(\pi_t) \right)}, \end{aligned}$$

563 where we applied Lemma 19 with the function  $f(x) = \frac{1}{\sqrt{x}}$  and  $a_i = \overline{\text{IR}}_i^{(2)}(\pi_i)$  to get the second  
 564 inequality. This can be seen as a generalization of the usual  $\sum_{t=1}^T \frac{1}{\sqrt{t}} \leq 2\sqrt{T}$  inequality. We  
 565 now define  $R_T^{(2)} = \sqrt{s \left( 2d + \sum_{t=1}^{T-1} \overline{\text{IR}}_t^{(2)}(\pi_t) \right)}$ , the constant-free regret rate associated to the 2-  
 566 surrogate-information ratio.

567 We now turn our attention to the 3-information ratio. Applying Lemma 7 we obtain that for all  
 568  $t \geq 1$ ,  $\overline{\text{IR}}_t^{(3)}(\pi_t) \leq 54 \frac{s}{C_{\min}} \leq 54 \frac{s^2}{C_{\min}}$  and for any  $T \geq 1$

$$\begin{aligned} \sum_{t=1}^T \sqrt{\lambda_{t-1}^{(3)} \overline{\text{IR}}_t^{(3)}(\pi_t)} &= s^{\frac{1}{3}} \sum_{t=1}^T \frac{\sqrt{\overline{\text{IR}}_t^{(3)}(\pi_t)}}{\left( \frac{3\sqrt{6}s}{\sqrt{C_{\min}}} + \sum_{s=1}^{t-1} \sqrt{\overline{\text{IR}}_s^{(3)}(\pi_s)} \right)^{\frac{1}{3}}} \\ &\leq s^{\frac{1}{3}} \sum_{t=1}^T \frac{\sqrt{\overline{\text{IR}}_t^{(3)}(\pi_t)}}{\left( \sum_{s=1}^t \sqrt{\overline{\text{IR}}_s^{(3)}(\pi_s)} \right)^{\frac{1}{3}}} \\ &\leq \frac{3}{2} s^{\frac{1}{3}} \left( \sum_{t=1}^T \sqrt{\overline{\text{IR}}_t^{(3)}(\pi_t)} \right)^{\frac{2}{3}} \\ &\leq \frac{3}{2} s^{\frac{1}{3}} \left( \frac{3\sqrt{6}s}{\sqrt{C_{\min}}} + \sum_{t=1}^{T-1} \sqrt{\overline{\text{IR}}_t^{(3)}(\pi_t)} \right), \end{aligned}$$

569 where we applied Lemma 19 with the function  $f(x) = \frac{1}{x^{\frac{2}{3}}}$  and  $a_i = \sqrt{\overline{\text{IR}}_i^{(3)}(\pi_i)}$  to get the  
 570 second inequality. This can be seen as a generalization of the usual  $\sum_{t=1}^T \frac{1}{t^{\frac{2}{3}}} \leq \frac{3}{2} T^{\frac{2}{3}}$ . We  
 571 now define  $R_T^{(3)} = s^{\frac{1}{3}} \left( \frac{3\sqrt{6}s}{\sqrt{C_{\min}}} + \sum_{t=1}^{T-1} \sqrt{\overline{\text{IR}}_t^{(3)}(\pi_t)} \right)^{\frac{2}{3}}$ , the constant-free regret rate associated  
 572 to the 3-surrogate-information ratio. We now consider the last time that the learning rates  $\lambda_t^{(3)}$   
 573 and  $\lambda_t^{(2)}$  have been used. More specifically, we denote  $T_2 = \max\{t \leq T, \lambda_{t-1}^{(2)} \geq \lambda_{t-1}^{(3)}\}$ , and  
 574  $T_3 = \max\{t \leq T, \lambda_{t-1}^{(3)} \geq \lambda_{t-1}^{(2)}\}$ . Coming back to the bound of Equation 24 and using the defini-  
 575 tion  $\lambda_t = \max(\lambda_t^{(2)}, \lambda_t^{(3)})$ , the following bound holds

$$\begin{aligned} R_T &\leq \mathbb{E} \left[ \frac{C_T}{\lambda_{T-1}} + \sum_{t=1}^T \min \left( \frac{32}{3} \lambda_{t-1} \overline{\text{IR}}_t^{(2)}(\pi_t), \frac{16}{3} c_3^* \sqrt{3 \lambda_{t-1} \overline{\text{IR}}_t^{(3)}(\pi_t)} \right) \right] \\ &\leq \mathbb{E} \left[ C_T \min \left( \frac{1}{\lambda_{T-1}^{(2)}}, \frac{1}{\lambda_{T-1}^{(3)}} \right) + \sum_{t=1}^T \min \left( \frac{32}{3} \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)}) \overline{\text{IR}}_t^{(2)}(\pi_t), \frac{16}{3} c_3^* \sqrt{3 \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)}) \overline{\text{IR}}_t^{(3)}(\pi_t)} \right) \right]. \end{aligned}$$

576 We can now separate the sum obtained at the last line based on which learning rate was used at time  
 577  $t$ .

$$\begin{aligned} &\sum_{t=1}^T \min \left( \frac{32}{3} \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)}) \overline{\text{IR}}_t^{(2)}(\pi_t), \frac{16}{3} c_3^* \sqrt{3 \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)}) \overline{\text{IR}}_t^{(3)}(\pi_t)} \right) \\ &\leq \sum_{\lambda_{t-1}^{(2)} \geq \lambda_{t-1}^{(3)}} \frac{32}{3} \lambda_{t-1}^{(2)} \overline{\text{IR}}_t^{(2)}(\pi_t) + \sum_{\lambda_{t-1}^{(3)} \geq \lambda_{t-1}^{(2)}} \frac{16}{3} c_3^* \sqrt{3 \lambda_{t-1}^{(3)} \overline{\text{IR}}_t^{(3)}(\pi_t)} \\ &\leq \sum_{t=1}^{T_2} \frac{32}{3} \lambda_{t-1}^{(2)} \overline{\text{IR}}_t^{(2)}(\pi_t) + \sum_{t=1}^{T_3} \frac{16}{3} c_3^* \sqrt{3 \lambda_{t-1}^{(3)} \overline{\text{IR}}_t^{(3)}(\pi_t)}. \end{aligned}$$

578 We further bound  $\sum_{t=1}^{T_2} \frac{32}{3} \lambda_{t-1}^{(2)} \overline{\text{IR}}_t^{(2)}(\pi_t) \leq \frac{64}{3} R_{T_2}^{(2)}$  and  $\sum_{t=1}^{T_3} \frac{16}{3} c_3^* \sqrt{3 \lambda_{t-1}^{(3)} \overline{\text{IR}}_t^{(3)}(\pi_t)} \leq \frac{16}{3} R_{T_3}^{(3)}$   
 579 (Using the explicit value  $c_3^* = \frac{2}{3^{\frac{3}{2}}}$ ).

580 The crucial observation is that which of  $\lambda_T^{(3)}$  or  $\lambda_T^{(2)}$  is bigger will determine whether  $R_T^{(2)}$  or  
 581  $R_T^{(3)}$  is the term of leading order (up to some constants). More specifically, Let  $T$  be such that

582  $\lambda_{T-1}^{(2)} \geq \lambda_{T-1}^{(3)}$  which means that  $\sqrt{\frac{s}{2d + \sum_{t=1}^{T-1} \overline{\mathbf{R}}_t^{(2)}(\pi_t)}} \geq \left( \frac{s}{\frac{3\sqrt{6}s}{\sqrt{C_{\min}}} + \sum_{t=1}^{T-1} \sqrt{\overline{\mathbf{R}}_t^{(3)}(\pi_t)}} \right)^{\frac{2}{3}}$ . Rearrang-  
583 ing, this implies that  $\sqrt{s} \left( 2d + \sum_{s=1}^{T-1} \overline{\mathbf{R}}_t^{(2)}(\pi_t) \right) \leq s^{\frac{2}{3}} \left( \frac{3\sqrt{6}s}{\sqrt{C_{\min}}} + \sum_{t=1}^{T-1} \sqrt{\overline{\mathbf{R}}_t^{(3)}(\pi_t)} \right)^{\frac{3}{2}}$ , which  
584 means that  $R_T^{(2)} \leq R_T^{(3)}$ . Following the exact same steps, we also have that  $\lambda_{T-1}^{(3)} \geq \lambda_{T-1}^{(2)}$  implies  
585 that  $R_T^{(3)} \leq R_T^{(2)}$ . We apply this to the time  $T_2$  in which  $\lambda_{T_2-1}^{(2)} \geq \lambda_{T_2-1}^{(3)}$  by definition. we have that  
586  $R_{T_2}^{(2)} \leq R_{T_2}^{(3)}$  and putting this together with the previous bound, we have

$$\begin{aligned} R_T &\leq \mathbb{E} \left[ \frac{C_T}{\lambda_{T-1}^{(3)}} + \frac{64}{3} R_{T_2}^{(2)} + \frac{16}{3} R_{T_3}^{(3)} \right] \\ &\leq \mathbb{E} \left[ \frac{C_T}{s} R_T^{(3)} + \frac{64}{3} R_{T_2}^{(2)} + \frac{16}{3} R_{T_3}^{(3)} \right] \\ &\leq \mathbb{E} \left[ \frac{C_T}{s} R_T^{(3)} + \frac{64}{3} R_{T_2}^{(3)} + \frac{16}{3} R_{T_3}^{(3)} \right] \\ &\leq \mathbb{E} \left[ \frac{C_T}{s} R_T^{(3)} + \frac{64}{3} R_T^{(3)} + \frac{16}{3} R_T^{(3)} \right] \\ &\leq \mathbb{E} \left[ \left( \frac{C_T}{s} + \frac{80}{3} \right) R_T^{(3)} \right], \end{aligned}$$

587 where we use the fact that  $T \rightarrow R_T^{(2)}$  and  $T \rightarrow R_T^{(3)}$  are non-decreasing and  $T_2 \leq T, T_3 \leq T$   
588 Similarly by definition of  $T_3$ , we have that  $\lambda_{T_3-1}^{(3)} \geq \lambda_{T_3-1}^{(2)}$  and we can conclude that  $R_{T_3}^{(3)} \leq R_{T_3}^{(2)}$ .  
589 Putting this together, with the previous bound, we have

$$\begin{aligned} R_T &\leq \mathbb{E} \left[ \frac{C_T}{\lambda_{T-1}^{(3)}} + \frac{64}{3} R_{T_2}^{(2)} + \frac{16}{3} R_{T_3}^{(3)} \right] \\ &\leq \mathbb{E} \left[ \frac{C_T}{s} R_T^{(2)} + \frac{64}{3} R_{T_2}^{(2)} + \frac{16}{3} R_{T_3}^{(3)} \right] \\ &\leq \mathbb{E} \left[ \frac{C_T}{s} R_T^{(2)} + \frac{64}{3} R_{T_2}^{(2)} + \frac{16}{3} R_{T_3}^{(2)} \right] \\ &\leq \mathbb{E} \left[ \frac{C_T}{s} R_T^{(2)} + \frac{64}{3} R_T^{(2)} + \frac{16}{3} R_T^{(2)} \right] \\ &\leq \mathbb{E} \left[ \left( \frac{C_T}{s} + \frac{80}{3} \right) R_T^{(2)} \right], \end{aligned}$$

590 where we use the fact that  $T \rightarrow R_T^{(2)}$  and  $T \rightarrow R_T^{(3)}$  are non-decreasing and  $T_2 \leq T, T_3 \leq T$ .  
591 Putting both of those bounds together with Equation 25 yields the claim of the Theorem.

## 592 E Maximum likelihood estimation

593 The focus of this section is to bound the difference between the log-likelihoods associated with the  
594 true parameter and the maximum likelihood estimator (MLE). We start by establishing an upper  
595 bound that holds in expectation which suffices to handle history-independent learning rates. Then,  
596 we move on to high-probability bounds that will allow us to deal with data-dependent learning rates.

### 597 E.1 Bound in expectation

598 We start with the case in which the maximum likelihood estimator is computed on a finite subset of  
599 the parameter space  $\Theta$ .

600 **Lemma 15.** *Let  $t \geq 1$ , and  $\Theta'$  be a finite subset of  $\Theta$ , we define the MLE over  $\Theta'$  as*

$$\theta_{MLE,t}(\Theta') = \arg \min_{\theta \in \Theta'} L_t^{(1)}(\theta).$$

601 Then,

$$\mathbb{E} \left[ L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_{MLE,t}(\Theta')) \right] \leq \log |\Theta'| \quad (26)$$

602 *Proof.* By the concavity of the logarithm and Jensen's inequality, we have

$$\begin{aligned} \mathbb{E} \left[ L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_{MLE,t}(\Theta')) \right] &\leq \log \mathbb{E} \left[ \prod_{s=1}^t \frac{p(Y_s | \theta_{MLE,t}(\Theta'), A_s)}{p(Y_s | \theta_0, A_s)} \right] \\ &= \log \mathbb{E} \left[ \max_{\theta \in \Theta'} \prod_{s=1}^t \frac{p(Y_s | \theta, A_s)}{p(Y_s | \theta_0, A_s)} \right] \leq \log \mathbb{E} \left[ \sum_{\theta \in \Theta'} \prod_{s=1}^t \frac{p(Y_s | \theta, A_s)}{p(Y_s | \theta_0, A_s)} \right] \\ &= \log \sum_{\theta \in \Theta'} \mathbb{E} \left[ \prod_{s=1}^t \frac{p(Y_s | \theta, A_s)}{p(Y_s | \theta_0, A_s)} \right] \end{aligned}$$

603 By Lemma 25, we have that  $\exp \left( L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta) \right) = \prod_{s=1}^t \frac{p(Y_s | \theta, A_s)}{p(Y_s | \theta_0, A_s)}$  is a non-negative supermartingale with respect to the filtration  $\mathcal{F}'_t = \sigma(\mathcal{F}_{t-1}, A_t)$ . That implies that each term in the sum is upper bounded by 1. Hence,

$$\mathbb{E} \left[ L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_{MLE,t}(\Theta')) \right] \leq \log \sum_{\theta \in \Theta'} 1 = \log |\Theta'|,$$

606 which proves the claim.  $\square$

607 To extend the previous bound to the full parameter space, we use a covering argument. A subset  $\Theta' \subset \Theta$  is said to be a valid  $\rho$ -covering of  $\Theta$  with respect to the  $\ell_1$  norm if for every  $\theta \in \Theta$ , there exists a  $\theta' \in \Theta'$  such that  $\|\theta - \theta'\|_1 \leq \rho$ . We denote by  $\mathcal{N}(\Theta, \|\cdot\|_1, \rho)$  the smallest possible cardinality of a valid  $\rho$  covering. We have the following bound on this quantity.

611 **Lemma 16.** For every  $\rho > 0$ ,

$$\log \mathcal{N}(\Theta, \|\cdot\|_1, \rho) \leq \log \binom{d}{s} \left(1 + \frac{2}{\rho}\right)^s \leq s \log \frac{ed(1 + 2/\rho)}{s}.$$

612

613 *Proof.* For each subset  $S \subset [d]$  of cardinality  $|S| = s$ , there is a surjective isometric embedding from  $(\Theta_S, \|\cdot\|_1)$  to  $(\mathbb{B}_1^s(1), \|\cdot\|_1)$ . In particular, to embed  $\theta \in \Theta_S$  into  $\mathbb{B}_1^s(1)$ , one can simply remove all the components of  $\theta$  corresponding to indices not in  $S$ . Therefore, for every  $\rho > 0$ ,  $\mathcal{N}(\Theta_S, \|\cdot\|_1, \rho) \leq \mathcal{N}(\mathbb{B}_1^s(1), \|\cdot\|_1, \rho)$ . Moreover, via a standard argument, we have  $\mathcal{N}(\mathbb{B}_1^s(1), \|\cdot\|_1, \rho) \leq (1 + \frac{2}{\rho})^s$  (see, e.g., Lemma 5.7 in [Wainwright, 2019](#)). Now, let  $\Theta_{S,\rho}$  denote any minimal  $\rho$ -covering of  $\Theta_S$  and notice that for an arbitrary  $\theta \in \Theta$  with support  $S$ , there exists a subset  $\tilde{S}$  such that  $S \subseteq \tilde{S}$  and  $|\tilde{S}| = s$ . Therefore, there exists  $\tilde{\theta} \in \Theta_{\tilde{S},\rho}$  such that  $\|\theta - \tilde{\theta}\|_1 \leq \rho$ . Hence,  $\cup_{S \subset [d]: |S|=s} \Theta_{S,\rho}$  forms a valid  $\rho$ -covering of  $\Theta$  and its cardinality is bounded by

$$\mathcal{N}(\Theta, \|\cdot\|_1, \rho) \leq |\cup_{S \subset [d]: |S|=s} \Theta_{S,\rho}| \leq \sum_{S \subset [d]: |S|=s} \left(1 + \frac{2}{\rho}\right)^s = \binom{d}{s} \left(1 + \frac{2}{\rho}\right)^s.$$

621 and we conclude by the elementary inequality  $\binom{d}{s} \leq \left(\frac{de}{s}\right)^s$ .  $\square$

### 622 E.1.1 Proof of Lemma 10

623 We bound the difference between the log-likelihood of the true parameter and that of the maximum likelihood estimator on the full parameter space. To this end, let  $\rho > 0$  and  $\Theta'$  be a minimal valid  $\rho$ -cover of  $\Theta$  as is defined in Lemma 16, and  $\theta' \in \Theta'$  be such that  $\|\theta' - \theta_t\| \leq \rho$ , which exists by

626 definition of a  $\rho$ -covering. Then,

$$\begin{aligned} \mathbb{E} \left[ L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_t) \right] &= \mathbb{E} \left[ L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_{\text{MLE},t}(\Theta')) \right] \\ &\quad + \mathbb{E} \left[ L_t^{(1)}(\theta_{\text{MLE},t}(\Theta')) - L^{(1)}(\theta') \right] \\ &\quad + \mathbb{E} \left[ L_t^{(1)}(\theta') - L^{(1)}(\theta_t) \right] \\ &\leq \log(\mathcal{N}(\Theta, \|\cdot\|_1, \rho)) + 0 + 2\rho t, \end{aligned}$$

627 where the first term is bounded by Lemma 26, the second term is non-positive by definition of  
628 the maximum likelihood estimator because  $\theta' \in \Theta'$  and the third term is bounded because the  
629 mapping  $\theta \mapsto \mathbb{E} \left[ L_t^{(1)}(\theta) \right]$  is  $2t$ -Lipschitz with respect to the 1-norm by Lemma 21. Finally applying  
630 Lemma 16 and setting  $\rho = \frac{2}{t}$  yields the desired bound.  $\square$

## 631 E.2 High-probability bounds

632 We begin with the case where the maximum likelihood estimator is computed over a finite subset of  
633 the parameter space  $\Theta$  and provide a corresponding high-probability bound.

634 **Lemma 17.** *Let  $\Theta'$  be a finite subset of  $\Theta$ , we define  $\theta_{\text{MLE},t}(\Theta') = \arg \min_{\theta \in \Theta'} L_t^{(1)}(\theta)$ . Then*

$$\mathbb{P} \left[ \exists t \geq 1, L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_{\text{MLE},t}(\Theta')) \geq \log \frac{|\Theta'|}{\delta} \right] \leq \delta. \quad (27)$$

635 *Proof.* Fix  $\theta \in \Theta'$ . By Lemma 25, we have that  $\exp \left( L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta) \right) = \prod_{s=1}^t \frac{p(Y_s|\theta, A_s)}{p(Y_s|\theta_0, A_s)}$  is a  
636 non-negative supermartingale with respect to the filtration  $\mathcal{F}'_t = \sigma(\mathcal{F}_{t-1}, A_t)$ , allowing us to invoke  
637 Ville's inequality to get the following guarantee:

$$\mathbb{P} \left[ \exists t \geq 1, \exp(L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta)) \geq \frac{1}{\delta} \right] \leq \delta.$$

638 Taking the logarithm and a union bound on  $\Theta'$  yields the desired result.  $\square$

639 We now provide a bound on the expected product of a bounded random variable with the difference  
640 in log-likelihood between the true parameter and the maximum likelihood estimator.

641 **Lemma 18.** *Let  $B \in \mathbb{R}$  and  $X$  be a random variable satisfying  $0 \leq X \leq B$  almost surely. Then  
642 for any  $t \geq 1$ ,*

$$\begin{aligned} \mathbb{E} \left[ X(L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_t)) \right] &\leq \inf_{\delta, \rho > 0} \left\{ \mathbb{E} \left[ X s \log \frac{ed(1 + \frac{2}{\rho})}{s\delta^{\frac{1}{s}}} \right] + 2B\rho t + B\delta s \log \frac{e^{1+\frac{1}{s}}d(1 + \frac{2}{\rho})}{s\delta^{\frac{1}{s}}} \right\} \\ &\leq 4 + s \log \frac{2e^2 d T^2 B^2}{s} \mathbb{E} \left[ X + \frac{1}{T} \right]. \end{aligned} \quad (28)$$

643 *Proof.* Let  $\delta, \rho > 0$  and  $\Theta'$  be a minimal valid  $\rho$ -cover of  $\Theta$  as defined in Lemma 16,  $N = |\Theta'|$ ,  
644 let  $\theta' = \theta_{\text{MLE},t}(\Theta')$  and let  $\bar{\theta} \in \Theta'$  be such that  $\|\bar{\theta} - \theta_t\| \leq \rho$ , which exists by definition of a valid  
645  $\rho$ -cover. We have the following decomposition:

$$\begin{aligned} \mathbb{E} \left[ X(L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_t)) \right] &\leq \mathbb{E} \left[ X(L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta')) \mathbf{1}_{\{L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta') \leq \log \frac{N}{\delta}\}} \right] \\ &\quad + B \mathbb{E} \left[ (L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta')) \mathbf{1}_{\{L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta') > \log \frac{N}{\delta}\}} \right] \\ &\quad + B \mathbb{E} \left[ (L_t^{(1)}(\bar{\theta}) - L_t^{(1)}(\theta_t)) \right] + B \mathbb{E} \left[ (L_t^{(1)}(\theta') - L_t^{(1)}(\bar{\theta})) \right]. \end{aligned}$$

646 The first term is upper bounded by  $\mathbb{E} \left[ X \log \frac{N}{\delta} \right]$ , the third term is upper bounded by  $2B\rho t$  because  
647  $\mathbb{E} \left[ L_t^{(1)}(\cdot) \right]$  is  $2t$ -Lipschitz by Lemma 21. The fourth term is non-positive because  $\theta'$  minimizes the  
648 negative log likelihood on  $\Theta'$ . Finally, we turn our attention to the second term. To simplify the

649 computations, we let  $Y = L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta')$ , and compute  $\mathbb{E} \left[ Y \mathbf{1}_{\{Y > \log \frac{N}{\delta}\}} \right]$ . Conditioning on  
 650 wheter  $\epsilon$  is larger or smaller than  $\log \frac{N}{\delta}$  yields the following identity

$$\mathbb{P} \left[ Y \mathbf{1}_{\{Y \geq \log \frac{N}{\delta}\}} \geq \epsilon \right] = \begin{cases} \mathbb{P}[Y \geq \epsilon] & \text{if } \epsilon \geq \log \frac{N}{\delta}, \\ \mathbb{P}[Y \geq \log \frac{N}{\delta}] & \text{otherwise.} \end{cases}$$

651 We can now upper bound the expectation as follows

$$\begin{aligned} \mathbb{E} \left[ Y \mathbf{1}_{\{Y \geq \log \frac{N}{\delta}\}} \right] &= \int_0^\infty \mathbb{P} \left[ Y \mathbf{1}_{\{Y \geq \log \frac{N}{\delta}\}} \geq \epsilon \right] d\epsilon \\ &= \log \frac{N}{\delta} \mathbb{P} \left[ Y \geq \log \frac{N}{\delta} \right] + \int_{\log \frac{N}{\delta}}^\infty \mathbb{P}[Y \geq \epsilon] d\epsilon \\ &= \log \frac{N}{\delta} \mathbb{P} \left[ Y \geq \log \frac{N}{\delta} \right] + \int_0^\delta \frac{1}{\delta'} \mathbb{P} \left[ Y \geq \log \frac{N}{\delta'} \right] d\delta' \\ &\leq \delta \log \frac{N}{\delta} + \delta, \end{aligned}$$

652 where we used the change of variable  $\epsilon = \log \frac{N}{\delta'}$  and used  $\mathbb{P} \left[ Y \geq \log \frac{N}{\delta} \right] \leq \delta$  by Lemma 17.  
 653 Finally, putting everything together and using  $N \leq \mathcal{N}(\Theta, \|\cdot\|_1, \rho) \leq \left( \frac{ed(1+\frac{2}{\rho})}{s} \right)^s$ , by Lemma 16,  
 654 we get

$$\mathbb{E} \left[ X(L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_t)) \right] \leq \mathbb{E} \left[ Xs \log \frac{ed(1+\frac{2}{\rho})}{s\delta^{\frac{1}{s}}} \right] + 2B\rho t + B\delta s \log \frac{e^{1+\frac{1}{s}}d(1+\frac{2}{\rho})}{s\delta^{\frac{1}{s}}}.$$

655 To balance the trade-off between the approximation error and the covering complexity, we choose  
 656  $\rho = \frac{2}{BT}$ , and  $\delta = \frac{1}{BT}$  which yields the desired form of the logarithmic factors. Substituting these  
 657 into the bound completes the proof.  $\square$

## 658 E.2.1 Proof of Lemma 14

659 As was noted in the analysis, since  $\lambda_T$  is not used by the algorithm, we can replace  $\lambda_T$  by  $\lambda_{T-1}$  in  
 660 our computations. We have

$$\begin{aligned} &\mathbb{E} \left[ \sum_{t=1}^T \frac{\eta}{\lambda_{t-1}} (L_t^{(1)}(\theta_t) - L_t^{(1)}(\theta_0) + L_{t-1}^{(1)}(\theta_0) - L_{t-1}^{(1)}(\theta_{t-1})) + \frac{\eta}{\lambda_T} (L_T^{(1)}(\theta_0) - L_T^{(1)}(\theta_T)) \right] \\ &= \mathbb{E} \left[ \sum_{t=1}^T \frac{\eta}{\lambda_{t-1}} (L_t^{(1)}(\theta_t) - L_t^{(1)}(\theta_0)) - \sum_{t=1}^T \frac{\eta}{\lambda_t} (L_t^{(1)}(\theta_t) - L_t^{(1)}(\theta_0)) \right] \\ &= \eta \cdot \sum_{t=1}^T \mathbb{E} \left[ (L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_t)) \left( \frac{1}{\lambda_t} - \frac{1}{\lambda_{t-1}} \right) \right]. \end{aligned}$$

661 Let  $C_{1,T}$  be a deterministic upper bound on  $\left( \frac{1}{\lambda_{t+1}} - \frac{1}{\lambda_t} \right)$ . Applying Lemma 28 to  $X =$   
 662  $\left( \frac{1}{\lambda_{t+1}} - \frac{1}{\lambda_t} \right)$  and telescoping, we get

$$\begin{aligned} &\eta \cdot \sum_{t=1}^T \mathbb{E} \left[ (L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_t)) \left( \frac{1}{\lambda_t} - \frac{1}{\lambda_{t-1}} \right) \right] \\ &\leq \eta \left( 4 + s \log \frac{2e^2 dt^2 C_{1,T}^2}{s} \right) \sum_{t=1}^T \mathbb{E} \left[ \left( \frac{1}{\lambda_t} - \frac{1}{\lambda_{t-1}} \right) + \frac{1}{T} \right] \\ &\leq \eta \left( 4 + s \log \frac{2e^2 dt^2 C_{1,T}^2}{s} \right) \mathbb{E} \left[ \left( \frac{1}{\lambda_T} + 1 \right) \right] \\ &\leq \mathbb{E} \left[ \frac{\eta(12 + 3s \log \frac{2e^2 dt^2 C_{1,T}^2}{s})}{2\lambda_{T-1}} \right], \end{aligned}$$

663 where in the last step, we used  $1 \leq \frac{1}{2\lambda_T}$  which implies  $\frac{1}{\lambda_T} + 1 \leq \frac{3}{2\lambda_T}$ . This finishes the proof.  $\square$

## 664 **F Bounding the surrogate information ratio**

### 665 **F.1 Proof of Lemma 6**

666 The surrogate regret of a policy is directly related to its 2- and 3-information ratio by definition

$$\widehat{\Delta}_t(\pi) = \sqrt{\overline{\text{IG}}_t(\pi)\overline{\text{IR}}_t^{(2)}(\pi)} = \left(\overline{\text{IG}}_t(\pi)\overline{\text{IR}}_t^{(3)}(\pi)\right)^{\frac{1}{3}}.$$

667 By the AM-GM inequality, we have that for any  $\lambda > 0$ , the surrogate regret is controlled as follows

$$\widehat{\Delta}_t(\pi) \leq \frac{\overline{\text{IG}}_t(\pi)}{\lambda} + \frac{\lambda}{4}\overline{\text{IR}}_t^{(2)}(\pi).$$

668 Similarly, by Lemma 27 which generalizes the AM-GM inequality, we can obtain the following  
669 regret bound

$$\widehat{\Delta}_t(\pi) \leq \frac{\overline{\text{IG}}_t(\pi)}{\lambda} + c_3^* \sqrt{\lambda \overline{\text{IR}}_t^{(3)}(\pi)},$$

670 where  $c_3^* < 2$  is an absolute constant defined in Lemma 27. This concludes the proof.  $\square$

### 671 **F.2 Proof of Lemma 1**

672 The proof of Lemma 1 is essentially the same as the proof of Lemma 5.6 in Hao et al. [2021], but we  
673 state it here for completeness. Throughout this proof, we use  $\langle p, f \rangle = \sum_{a \in \mathcal{A}} p(a)f(a)$  to denote  
674 the inner product between a signed measure  $p$  on  $\mathcal{A}$  and a function  $f : \mathcal{A} \rightarrow \mathbb{R}$ . Using this notation,  
675 we can, for example, write the generalized surrogate information ratio as  $\overline{\text{IR}}_t^{(\gamma)}(\pi) = \langle \pi, \overline{\text{IR}}_t^{(\gamma)} \rangle$ .

676 We define  $\pi_t^{(\gamma)} \in \arg \min_{\pi \in \Delta(\mathcal{A})} \overline{\text{IR}}_t^{(\gamma)}(\pi)$  to be any minimizer of the generalized surrogate infor-  
677 mation ratio with parameter  $\gamma \geq 2$ . First, we observe that

$$\nabla_{\pi} \overline{\text{IR}}_t^{(2)}(\pi) = \frac{2\langle \pi, \widehat{\Delta}_t \rangle \widehat{\Delta}_t}{\langle \pi, \overline{\text{IG}}_t \rangle} - \frac{(\langle \pi, \widehat{\Delta}_t \rangle)^2 \overline{\text{IG}}_t}{(\langle \pi, \overline{\text{IG}}_t \rangle)^2}.$$

678 Therefore, from the first-order optimality condition for convex constrained minimization (and the  
679 fact that  $\overline{\text{IR}}_t^{(2)}$  is convex on  $\Delta(\mathcal{A})$ ), we have

$$\forall \pi \in \Delta(\mathcal{A}), 0 \leq \langle \pi - \pi_t^{(\text{SOIDS})}, \nabla_{\pi} \overline{\text{IR}}_t^{(2)}(\pi_t^{(\text{SOIDS})}) \rangle.$$

680 In particular,

$$0 \leq \frac{2\langle \pi_t^{(\text{SOIDS})}, \widehat{\Delta}_t \rangle \langle \pi_t^{(\gamma)} - \pi_t^{(\text{SOIDS})}, \widehat{\Delta}_t \rangle}{\langle \pi_t^{(\text{SOIDS})}, \overline{\text{IG}}_t \rangle} - \frac{(\langle \pi_t^{(\text{SOIDS})}, \widehat{\Delta}_t \rangle)^2 \langle \pi_t^{(\gamma)} - \pi_t^{(\text{SOIDS})}, \overline{\text{IG}}_t \rangle}{(\langle \pi_t^{(\text{SOIDS})}, \overline{\text{IG}}_t \rangle)^2}.$$

681 This inequality is equivalent to

$$2\langle \pi_t^{(\gamma)}, \widehat{\Delta}_t \rangle \geq \langle \pi_t^{(\text{SOIDS})}, \widehat{\Delta}_t \rangle \left( 1 + \frac{\langle \pi_t^{(\gamma)}, \overline{\text{IG}}_t \rangle}{\langle \pi_t^{(\text{SOIDS})}, \overline{\text{IG}}_t \rangle} \right) \geq \langle \pi_t^{(\text{SOIDS})}, \widehat{\Delta}_t \rangle.$$

682 From this inequality, we obtain

$$\begin{aligned} \frac{(\langle \pi_t^{(\text{SOIDS})}, \widehat{\Delta}_t \rangle)^{\gamma}}{\langle \pi_t^{(\text{SOIDS})}, \overline{\text{IG}}_t \rangle} &= \frac{(\langle \pi_t^{(\text{SOIDS})}, \widehat{\Delta}_t \rangle)^2 (\langle \pi_t^{(\text{SOIDS})}, \widehat{\Delta}_t \rangle)^{\gamma-2}}{\langle \pi_t^{(\text{SOIDS})}, \overline{\text{IG}}_t \rangle} \\ &\leq \frac{(\langle \pi_t^{(\gamma)}, \widehat{\Delta}_t \rangle)^2 (\langle \pi_t^{(\text{SOIDS})}, \widehat{\Delta}_t \rangle)^{\gamma-2}}{\langle \pi_t^{(\gamma)}, \overline{\text{IG}}_t \rangle} \\ &\leq 2^{\gamma-2} \frac{(\langle \pi_t^{(\gamma)}, \widehat{\Delta}_t \rangle)^{\gamma}}{\langle \pi_t^{(\gamma)}, \overline{\text{IG}}_t \rangle} = 2^{\gamma-2} \min_{\pi \in \Delta(\mathcal{A})} \overline{\text{IR}}_t^{(\gamma)}(\pi), \end{aligned}$$

683 thus proving the claim.  $\square$

684 **F.3 Proof of Lemma 7**

685 This section is focused on bounding the information ratios of the sparse optimistic information  
 686 directed sampling policy. As is widely done in the information directed sampling literature, we will  
 687 introduce a “forerunner” algorithm with controlled surrogate information ratio. By Lemma 1, the  
 688 sOIDS policy will then automatically inherit the bound of the forerunner.

689 As one of our forerunners, we will make use of the “Feel-Good Thompson Sampling” first intro-  
 690 duced by Zhang [2022]. Letting  $\tilde{\theta}_t \sim Q_t^+$ , the FGTS policy is defined as

$$\pi_t^{(\text{FGTS})}(a) = \mathbb{P}_t \left[ a^*(\tilde{\theta}_t) = a \right]. \quad (29)$$

691 Which can be seen as the policy obtained by sampling a parameter  $\tilde{\theta}_t \sim Q_t^+$  and then picking the  
 692 optimal action under this parameter. Compared to the usual Thompson Sampling policy, this boils  
 693 down to replacing the Bayesian posterior by the optimistic posterior. Whenever the optimal action  
 694 for  $\theta$  is non-unique, we define  $a^*(\theta)$  to be any optimal action with minimal 0-norm. If there are  
 695 multiple optimal actions with minimal 0-norm, ties can be broken arbitrarily.

696 For the bound on the surrogate 3-information ratio, we assume that the prior  $Q_1^+$  and the action set  
 697  $\mathcal{A}$  are such that for all  $\theta$  in the support of the prior, there exists  $a' \in \arg \max_{a \in \mathcal{A}} r(a, \theta)$  such that  
 698  $\|a'\|_0 \leq s$ . We refer to this as the sparse optimal action property. Since the support of our prior  $Q_1^+$   
 699 only contains  $s$ -sparse vectors, the sparse optimal action property is satisfied whenever the action  
 700 set is a unit  $\ell_p$  ball. Note also that the hard instances in both the  $\sqrt{sdT}$  lower bound in Theorem  
 701 24.3 of Lattimore and Szepesvári [2020] and the  $s^{2/3}T^{2/3}$  lower bound in Theorem 5 of Jang et al.  
 702 [2022] satisfy the sparse optimal action property<sup>2</sup>. Therefore, even with this additional assumption,  
 703 the lower bounds for both the data-rich and data-poor regimes remain meaningful. Whenever the  
 704 optimal action for  $\theta$  is non-unique, we define  $a^*(\theta)$  to be any optimal action with minimal 0-norm,  
 705 with ties broken arbitrarily.

706 **F.3.1 Bounding the two information ratio**

707 We will now prove the first part of lemma 7, by showing that the information ratio of the FGTS  
 708 policy is bounded by the dimension. The proof is exactly the same as in the Bayesian setting as  
 709 is done in Proposition 5 of Russo and Roy [2016], Lemma 7 of Lemma 7 in Neu et al. [2022] or  
 710 in Lemma 5.7 of Hao et al. [2021], except the Bayesian posterior is replaced with the optimistic  
 711 posterior. We provide the proof here for completeness.

712 Since we defined the surrogate information gain in terms of the model  $\theta$ , as opposed to the optimal  
 713 action  $a^*(\theta)$ , we follow the proof of Lemma 7 in Neu et al. [2022]. For brevity, we let  $\alpha_a =$   
 714  $\pi_t^{(\text{FGTS})}(a) = \mathbb{P}_t \left[ a^*(\tilde{\theta}_t) = a \right]$ . We define the  $|\mathcal{A}| \times |\mathcal{A}|$  matrix  $M$  by

$$M_{a,a'} = \sqrt{\alpha_a \alpha_{a'}} (\mathbb{E}_t[r(a, \tilde{\theta}_t) | a^*(\tilde{\theta}_t) = a'] - r(a, \bar{\theta}(Q_t^+))).$$

715 Next, we relate the surrogate information gain and the surrogate regret to the Frobenius norm and  
 716 the trace of  $M$ . First, we can lower bound the surrogate information gain of FGTS as

$$\begin{aligned} \overline{\text{IG}}_t(\pi_t^{(\text{FGTS})}) &= \frac{1}{2} \sum_{a \in \mathcal{A}} \alpha_a \int_{\Theta} (r(a, \bar{\theta}(Q_t^+)) - r(a, \theta))^2 dQ_t^+(\theta) \\ &= \frac{1}{2} \sum_{a \in \mathcal{A}} \alpha_a \int_{\Theta} \sum_{a' \in \mathcal{A}} \mathbf{1}_{\{a^*(\theta) = a'\}} (r(a, \bar{\theta}(Q_t^+)) - r(a, \theta))^2 dQ_t^+(\theta) \\ &= \frac{1}{2} \sum_{a \in \mathcal{A}} \sum_{a' \in \mathcal{A}} \alpha_a \int_{\Theta} \mathbf{1}_{\{a^*(\theta) = a'\}} dQ_t^+(\theta) \mathbb{E}_t[(r(a, \bar{\theta}(Q_t^+)) - r(a, \tilde{\theta}_t) | a^*(\tilde{\theta}_t) = a')] \\ &\geq \frac{1}{2} \sum_{a \in \mathcal{A}} \sum_{a' \in \mathcal{A}} \alpha_a \alpha_{a'} \left( r(a, \bar{\theta}(Q_t^+)) - \mathbb{E}_t[r(a, \tilde{\theta}_t) | a^*(\tilde{\theta}_t) = a'] \right)^2 \\ &= \frac{1}{2} \sum_{a \in \mathcal{A}} \sum_{a' \in \mathcal{A}} M_{a,a'}^2 = \frac{1}{2} \|M\|_F^2. \end{aligned}$$

<sup>2</sup>The optimal actions in the hard instance used to prove Theorem 5 in Jang et al. [2022] are  $2s$ -sparse, which still allows us to prove the same bound on the surrogate 3-information ratio, up to constant factors.

717 Next, we can re-write the surrogate regret of FGTS as

$$\begin{aligned}
\widehat{\Delta}_t(\pi_t^{(\text{FGTS})}) &= \int_{\Theta} r(a^*(\theta), \theta) dQ_t^+(\theta) - \sum_{a \in \mathcal{A}} \alpha_a \int_{\Theta} r(a, \theta) dQ_t^+ \\
&= \int_{\Theta} \sum_{a \in \mathcal{A}} \mathbf{1}_{\{a^*(\theta)=a\}} r(a^*(\theta), \theta) dQ_t^+(\theta) - \sum_{a \in \mathcal{A}} \alpha_a r(a, \bar{\theta}(Q_t^+)) \\
&= \sum_{a \in \mathcal{A}} \alpha_a \mathbb{E}_t[r(a, \tilde{\theta}_t) | a^*(\tilde{\theta}_t) = a] - \sum_{a \in \mathcal{A}} \alpha_a r(a, \bar{\theta}(Q_t^+)) \\
&= \text{tr}(M).
\end{aligned} \tag{30}$$

718 Using Fact 10 from [Russo and Roy \[2016\]](#), we bound  $\overline{\text{IR}}_t^{(2)}(\pi_t^{(\text{FGTS})})$  as

$$\overline{\text{IR}}_t^{(2)}(\pi_t^{(\text{FGTS})}) = \frac{(\widehat{\Delta}_t(\pi_t^{(\text{FGTS})}))^2}{\overline{\text{IG}}_t(\pi_t^{(\text{FGTS})})} \leq \frac{2(\text{tr}(M))^2}{\|M\|_F^2} \leq 2 \cdot \text{rank}(M).$$

719 All the remains is to show that  $M$  has rank at most  $d$ . Enumerate the actions as  $\mathcal{A} = \{a_1, \dots, a_{|\mathcal{A}|}\}$ ,  
720 and let  $\mu_i = \mathbb{E}_t[\tilde{\theta}_t | a^*(\tilde{\theta}_t) = a_i]$ . By linearity of expectation (and of the reward function), we can  
721 write

$$M_{i,j} = \sqrt{\alpha_i \alpha_j} \langle \mu_i - \bar{\theta}(Q_t^+), a_j \rangle.$$

722 Therefore,  $M$  can be factorised as

$$M = \begin{bmatrix} \sqrt{\alpha_1}(\mu_1 - \bar{\theta}(Q_t^+))^\top \\ \vdots \\ \sqrt{\alpha_{|\mathcal{A}|}}(\mu_{|\mathcal{A}|} - \bar{\theta}(Q_t^+))^\top \end{bmatrix} \begin{bmatrix} \sqrt{\alpha_1} a_1 & \cdots & \sqrt{\alpha_{|\mathcal{A}|}} a_{|\mathcal{A}|} \end{bmatrix}.$$

723 Since  $M$  is the product of a  $K \times d$  matrix and a  $d \times K$  matrix, it must have rank at most  $\min(K, d)$ .

### 724 F.3.2 Bounding the three information ratio

725 To bound the 3 information ratio we follow [Hao et al. \[2021\]](#) and we introduce the exploratory policy

$$\mu = \arg \max_{\pi \in \Delta(\mathcal{A})} \sigma_{\min} \left( \sum_{a \in \mathcal{A}} \pi(a) a a^\top \right). \tag{31}$$

726 We define the mixture policy  $\pi_t^{(\text{mix})} = (1 - \gamma)\pi_t^{(\text{FGTS})} + \gamma\mu$  where  $\gamma \geq 0$  will be determined later.  
727 First, we lower bound the surrogate information gain of the mixture policy in the same way that we  
728 lower bounded the surrogate information gain of the FGTS policy previously. This time, we obtain  
729 the lower bound

$$\begin{aligned}
\overline{\text{IG}}_t(\pi_t^{(\text{mix})}) &\geq \frac{1}{2} \sum_{a \in \mathcal{A}} \pi_t^{(\text{mix})}(a) \sum_{a' \in \mathcal{A}} \mathbb{P}_t(a^*(\tilde{\theta}_t) = a') (r(a, \bar{\theta}(Q_t^+)) - \mathbb{E}_t[r(a, \tilde{\theta}_t) | a^*(\tilde{\theta}_t) = a'])^2 \\
&= \frac{1}{2} \sum_{a \in \mathcal{A}} \pi_t^{(\text{mix})}(a) \sum_{a' \in \mathcal{A}} \mathbb{P}_t(a^*(\tilde{\theta}_t) = a') \langle \mu_{a'} - \bar{\theta}(Q_t^+), a \rangle^2,
\end{aligned}$$

730 where  $\mu_{a'} = \mathbb{E}_t[\tilde{\theta}_t | a^*(\tilde{\theta}_t) = a']$ . From the inequality  $\pi_t^{(\text{mix})}(a) \geq \gamma\mu(a)$ , and the definition of  
731  $C_{\min}$ , we have

$$\begin{aligned}
\overline{\text{IG}}_t(\pi_t^{(\text{mix})}) &\geq \frac{\gamma}{2} \sum_{a' \in \mathcal{A}} \mathbb{P}_t(a^*(\tilde{\theta}_t) = a') \sum_{a \in \mathcal{A}} \mu(a) (\mu_{a'} - \bar{\theta}(Q_t^+))^\top a a^\top (\mu_{a'} - \bar{\theta}(Q_t^+)) \\
&\geq \frac{\gamma}{2} \sum_{a' \in \mathcal{A}} \mathbb{P}_t(a^*(\tilde{\theta}_t) = a') C_{\min} \|\mu_{a'} - \bar{\theta}(Q_t^+)\|_2^2.
\end{aligned}$$

732 Using the expression for the surrogate regret of FGTS in (30), we obtain

$$\begin{aligned}
\widehat{\Delta}_t(\pi_t^{(\text{FGTS})}) &= \sum_{a \in \mathcal{A}} \mathbb{P}_t(a^*(\tilde{\theta}_t) = a) (\mathbb{E}_t[\langle \tilde{\theta}_t, a \rangle | a^*(\tilde{\theta}_t) = a] - \langle \bar{\theta}(Q_t^+), a \rangle) \\
&\leq \sqrt{\sum_{a \in \mathcal{A}} \mathbb{P}_t(a^*(\tilde{\theta}_t) = a) (\mathbb{E}_t[\langle \tilde{\theta}_t, a \rangle | a^*(\tilde{\theta}_t) = a] - \langle \bar{\theta}(Q_t^+), a \rangle)^2},
\end{aligned}$$

733 where in the last line we used the Cathy-Schwarz inequality. Due to the sparse optimal action  
734 property, all actions for which  $\mathbb{P}_t(a^*(\tilde{\theta}_t) = a) > 0$  have at most  $s$  non-zero elements. Therefore,

$$\sum_{a \in \mathcal{A}} \mathbb{P}_t(a^*(\tilde{\theta}_t) = a) (\mathbb{E}_t[\langle \tilde{\theta}_t, a \rangle | a^*(\tilde{\theta}_t) = a] - \langle \bar{\theta}(Q_t^+), a \rangle)^2 \leq \sum_{a \in \mathcal{A}} \mathbb{P}_t(a^*(\tilde{\theta}_t) = a) s \|\mu_a - \bar{\theta}(Q_t^+)\|_2^2.$$

735 This, combined with the lower bound on  $\overline{\text{IG}}_t(\pi_t^{(\text{mix})})$  means that

$$\begin{aligned} \widehat{\Delta}_t(\pi_t^{(\text{FGTS})}) &\leq \sqrt{\sum_{a \in \mathcal{A}} \mathbb{P}_t(a^*(\tilde{\theta}_t) = a) s \|\mu_a - \bar{\theta}(Q_t^+)\|_2^2} \\ &= \sqrt{\frac{2s}{\gamma C_{\min}} \frac{\gamma}{2} \sum_{a \in \mathcal{A}} \mathbb{P}_t(a^*(\tilde{\theta}_t) = a) C_{\min} \|\mu_a - \bar{\theta}(Q_t^+)\|_2^2} \\ &\leq \sqrt{\frac{2s}{\gamma C_{\min}} \overline{\text{IG}}_t(\pi_t^{(\text{mix})})}. \end{aligned}$$

736 Choosing  $\gamma = 1$ , this tells us that

$$(\widehat{\Delta}_t(\pi_t^{(\text{FGTS})}))^2 \leq \frac{2s}{C_{\min}} \overline{\text{IG}}_t(\mu).$$

737 We bound the information ratio in three cases. First, suppose that  $\widehat{\Delta}_t(\mu) \leq \widehat{\Delta}_t(\pi_t^{(\text{FGTS})})$ . In this  
738 case,

$$\overline{\text{IR}}_t^{(3)}(\mu) = \frac{\widehat{\Delta}_t(\mu)(\widehat{\Delta}_t(\mu))^2}{\overline{\text{IG}}_t(\mu)} \leq \frac{2(\widehat{\Delta}_t(\pi_t^{(\text{FGTS})}))^2}{\overline{\text{IG}}_t(\mu)} \leq \frac{4s}{C_{\min}}.$$

739 Next, we consider the case where  $\widehat{\Delta}_t(\mu) > \widehat{\Delta}_t(\pi_t^{(\text{FGTS})})$ . For any  $\gamma \in (0, 1]$ ,

$$\overline{\text{IR}}_t^{(3)}(\pi_t^{(\text{mix})}) = \frac{((1-\gamma)\widehat{\Delta}_t(\pi_t^{(\text{FGTS})}) + \gamma\widehat{\Delta}_t(\mu))^3}{(1-\gamma)\overline{\text{IG}}_t(\pi_t^{(\text{FGTS})}) + \gamma\overline{\text{IG}}_t(\mu)} \leq \frac{((1-\gamma)\widehat{\Delta}_t(\pi_t^{(\text{FGTS})}) + \gamma\widehat{\Delta}_t(\mu))^3}{\gamma\overline{\text{IG}}_t(\mu)}.$$

740 We define  $f(\gamma) = ((1-\gamma)\widehat{\Delta}_t(\pi_t^{(\text{FGTS})}) + \gamma\widehat{\Delta}_t(\mu))^3 / (\gamma\overline{\text{IG}}_t(\mu))$  to be the RHS of the previous  
741 equation. One can verify that the derivative of  $f(\gamma)$  is

$$f'(\gamma) = \frac{((1-\gamma)\widehat{\Delta}_t(\pi_t^{(\text{FGTS})}) + \gamma\widehat{\Delta}_t(\mu))^2}{\gamma^2\overline{\text{IG}}_t(\mu)} \left[ 2\gamma(\widehat{\Delta}_t(\mu) - \widehat{\Delta}_t(\pi_t^{(\text{FGTS})})) - \widehat{\Delta}_t(\pi_t^{(\text{FGTS})}) \right],$$

742 and that  $f(\gamma)$  is minimised w.r.t.  $\gamma > 0$  at  $\widehat{\gamma}$ , where  $\widehat{\gamma}$  is the positive solution of  $f'(\widehat{\gamma}) = 0$ , which is

$$\widehat{\gamma} = \frac{\widehat{\Delta}_t(\pi_t^{(\text{FGTS})})}{2(\widehat{\Delta}_t(\mu) - \widehat{\Delta}_t(\pi_t^{(\text{FGTS})}))}.$$

743 That  $\widehat{\gamma}$  is always positive follows from the fact that  $\widehat{\Delta}_t(\mu) > \widehat{\Delta}_t(\pi_t^{(\text{FGTS})})$ . If  $\widehat{\gamma} \leq 1$ , then we can  
744 take the forerunner to be the mixture policy with  $\gamma = \widehat{\gamma}$ . In this case,

$$\begin{aligned} \overline{\text{IR}}_t^{(3)}(\pi_t^{(\text{mix})}) &= \frac{(\frac{3}{2})^3 2(\widehat{\Delta}_t(\mu) - \widehat{\Delta}_t(\pi_t^{(\text{FGTS})})) \widehat{\Delta}_t(\pi_t^{(\text{FGTS})})^2}{\overline{\text{IG}}_t(\mu)} \\ &\leq \frac{(\frac{3}{2})^3 8s}{C_{\min}} = \frac{27s}{C_{\min}}. \end{aligned}$$

745 Otherwise, if  $\widehat{\gamma} > 1$ , then

$$\widehat{\Delta}_t(\mu) \leq \frac{3}{2} \widehat{\Delta}_t(\pi_t^{(\text{FGTS})}).$$

746 In this case, we can take the forerunner to be  $\mu$ . The surrogate 3-information ratio can then be upper  
747 bounded as

$$\overline{\text{IR}}_t^{(3)}(\mu) = \frac{\widehat{\Delta}_t(\mu)(\widehat{\Delta}_t(\mu))^2}{\overline{\text{IG}}_t(\mu)} \leq \frac{2(\frac{3}{2})^2 (\widehat{\Delta}_t(\pi_t^{(\text{FGTS})}))^2}{\overline{\text{IG}}_t(\mu)} \leq \frac{(\frac{3}{2})^2 4s}{C_{\min}} = \frac{9s}{C_{\min}}.$$

748 Therefore, one can always find a value of  $\gamma \in (0, 1]$  such that

$$\overline{\text{IR}}_t^{(3)}(\pi_t^{(\text{mix})}) \leq \frac{27s}{C_{\min}}.$$

749 **G Choosing the learning rates**

750 This section is focused on the choice of the learning rates required to obtain the bound of Theorem 2.

751 **G.1 Technical tools**

752 We start by a collection of technical results to help with choosing a time-dependent learning rate.

753 **Lemma 19.** *Let  $a_i \geq 0$  and  $f : [0, \infty) \rightarrow [0, \infty)$  be a nonincreasing function. Then*

$$\sum_{t=1}^T a_t f\left(\sum_{i=0}^t a_i\right) \leq \int_{a_0}^{\sum_{t=0}^T a_t} f(x) dx. \quad (32)$$

754 The proof follows from elementary manipulations comparing sums and integrals. The result is taken  
755 from Lemma 4.13 of [Orabona \[2019\]](#), where a complete proof is also supplied. The following  
756 lemma ensures that the learning rates are non-increasing.

757 **Lemma 20.** *Let  $C_1 > e, C_2 > 0$  and define  $\lambda_t = \frac{\log(C_1 t)}{C_2 t}$ , then  $\lambda_t$  is a non-decreasing sequence.*

758 *Proof.* Let  $t > 0$ , we have

$$\frac{\log(C_1(t+1))}{\log(C_1 t)} = \frac{\log(C_1 t \left(\frac{t+1}{t}\right))}{\log(C_1 t)} = \frac{\log(C_1 t) + \log\left(\frac{t+1}{t}\right)}{\log(C_1 t)} \leq 1 + \frac{1}{t \log(C_1 t)} \leq 1 + \frac{1}{t},$$

759 where the first inequality uses  $\log(1+x) \leq x$  for any  $x > -1$  and the second inequality uses  
760  $\log(C_1 t) \geq \log(C_1) \geq 1$  because we assumed  $C_1 \geq e$ . Since  $\frac{C_2(t+1)}{C_2 t} = 1 + \frac{1}{t}$ , we can conclude  
761 that the sequence  $\lambda_t$  is non-increasing.  $\square$

762 **G.2 Data-rich regime: Proof of Lemma 8**

763 We start by focusing on the data rich regime, and we bound the following part of the regret bound  
764 given in Equation (12):

$$\frac{C_T}{\lambda_{T-1}} + \frac{32}{3} \sum_{t=1}^T \lambda_{t-1} \overline{\mathbf{R}}_t^{(2)}(\pi_t).$$

765 Here,  $C_T = 5 + 2s \log \frac{edT}{s}$ . To proceed, we let  $\lambda_t = \alpha \sqrt{\frac{C_{t+1}}{d(t+1)}}$ , where  $\alpha > 0$  is a constant that we  
766 will optimize later. Because  $t \rightarrow C_t$  is increasing, we get that  $\lambda_{t-1} \leq \alpha \sqrt{\frac{C_t}{dt}}$ . By Lemma 7, we  
767 know that for all  $t \geq 1$ ,  $\overline{\mathbf{R}}_t^{(2)}(\pi_t) \leq 2d$ , hence

$$\begin{aligned} \frac{C_T}{\lambda_{T-1}} + \frac{32}{3} \sum_{t=1}^T \lambda_{t-1} \overline{\mathbf{R}}_t^{(2)}(\pi_t) &\leq \frac{1}{\alpha} \sqrt{C_T d T} + \frac{64}{3} \alpha \sqrt{C_T} \sum_{t=1}^T \frac{d}{\sqrt{dt}} \\ &\leq \frac{1}{\alpha} \sqrt{C_T d T} + \frac{128}{3} \alpha \sqrt{C_T d T} \\ &\leq \left( \frac{1}{\alpha} + \frac{128}{3} \alpha \right) \sqrt{C_T d T} \\ &\leq 16 \sqrt{\frac{2}{3} C_T d T}, \end{aligned}$$

768 where the second line uses the standard inequality  $\sum_{t=1}^T \frac{1}{\sqrt{t}} \leq 2\sqrt{T}$ , and the last line is obtained by  
769 optimizing the expression  $\left(\frac{1}{\alpha} + \frac{128}{3} \alpha\right)$  with the optimal choice  $\alpha = \sqrt{\frac{3}{128}}$  which yields the value  
770  $16\sqrt{\frac{2}{3}}$ . This concludes the proof of the claim.  $\square$

771 **G.3 Data-poor regime: proof of Lemma 8**

772 We now focus on the data-poor regime and specifically on bounding the following part of the bound  
 773 given in Equation (12):

$$\frac{C_T}{\lambda_{T-1}} + \frac{16}{3} c_3^* \sum_{t=1}^T \sqrt{3\lambda_{t-1} \overline{\text{IR}}_t^{(3)}(\pi_t)}.$$

774 Here,  $C_T = 5 + 2s \log \frac{edT}{s}$ . Now, we let  $\lambda_t = \alpha \left( \frac{C_{t+1} \sqrt{C_{\min}}}{(t+1)\sqrt{s}} \right)^{\frac{2}{3}}$ , where  $\alpha > 0$  is a constant that  
 775 we will optimize later. Because  $t \rightarrow C_t$  is increasing, we get that  $\lambda_{t-1} \leq \alpha \left( \frac{C_T \sqrt{C_{\min}}}{t s} \right)^{\frac{2}{3}}$ . By  
 776 Lemma 7, the 3-surrogate-information ratio is bounded for all  $t \geq 1$  as  $\overline{\text{IR}}_t^{(3)}(\pi_t) \leq \frac{54s}{C_{\min}}$ . Hence,  
 777 the following holds:

$$\begin{aligned} \frac{C_T}{\lambda_{T-1}} + \frac{16}{3} c_3^* \sum_{t=1}^T \sqrt{3\lambda_{t-1} \overline{\text{IR}}_t^{(3)}(\pi_t)} &\leq \frac{1}{\alpha} (C_T)^{\frac{1}{3}} \left( \frac{T\sqrt{s}}{\sqrt{C_{\min}}} \right)^{\frac{2}{3}} + 48c_3^* \sqrt{2\alpha} (C_T)^{\frac{1}{3}} \left( \frac{\sqrt{s}}{\sqrt{C_{\min}}} \right)^{\frac{2}{3}} \sum_{t=1}^T \frac{1}{t^{\frac{1}{3}}} \\ &\leq \frac{1}{\alpha} (C_T)^{\frac{1}{3}} \left( \frac{T\sqrt{s}}{\sqrt{C_{\min}}} \right)^{\frac{2}{3}} + 72c_3^* \sqrt{2\alpha} (C_T)^{\frac{1}{3}} \left( \frac{T\sqrt{s}}{\sqrt{C_{\min}}} \right)^{\frac{2}{3}} \\ &\leq \left( \frac{1}{\alpha} + 72c_3^* \sqrt{2\alpha} \right) (C_T)^{\frac{1}{3}} \left( \frac{T\sqrt{s}}{\sqrt{C_{\min}}} \right)^{\frac{2}{3}} \\ &\leq 12 \cdot 6^{\frac{1}{3}} (C_T)^{\frac{1}{3}} \left( \frac{T\sqrt{s}}{\sqrt{C_{\min}}} \right)^{\frac{2}{3}}. \end{aligned}$$

778 Here, we have applied Lemma 19 with the function  $f(x) = x^{\frac{1}{3}}$  and  $a_i = 1$  to bound  $\sum_{t=1}^T t^{-1/3} \leq$   
 779  $\frac{3}{2} T^{\frac{2}{3}}$  in the second line, the last line comes from the choice  $\alpha = \frac{1}{4 \cdot 6^{\frac{1}{3}}}$  which optimizes the constant  
 780  $\left( \frac{1}{\alpha} + 144c_3^* \sqrt{2\alpha} \right)$  (as per Lemma 27). This proves the statement.  $\square$

781 **G.4 Joint learning rates, end of the proof of Theorem 2**

782 In the section below, we present the technical derivation related to choosing the choice of learning  
 783 rate  $\lambda_t = \min(\frac{1}{2}, \max(\lambda_t^{(2)}, \lambda_t^{(3)}))$ , where  $\lambda_t^{(2)} = \sqrt{\frac{3C_{t+1}}{128d(t+1)}}$  and  $\lambda_t^{(3)} = \frac{1}{4 \cdot 6^{\frac{1}{3}}} \left( \frac{C_{t+1} \sqrt{C_{\min}}}{(t+1)\sqrt{s}} \right)^{\frac{2}{3}}$ ,  
 784 with  $C_t = 5 + 2s \log \frac{edt}{s}$ . This choice interpolates between the data-rich and data-poor regimes. As  
 785 a first step, we start by confirming via Lemma 20 that both  $\lambda_t^{(2)}$  and  $\lambda_t^{(3)}$  are non-increasing and the  
 786 bound of Theorem 1 holds with our choice of  $\lambda_t$ .

787 First, note that our choice of learning rates ensures that  $\lambda_t \leq \frac{1}{2}$  holds as long as  $T$  is larger than  
 788 an absolute constant, and thus we focus on this case here (and relegate the complete details of  
 789 establishing this absolute constant to Appendix G.5). To proceed, we define the (constant-free)  
 790 regret rates  $R_t^{(2)} = \sqrt{C_t dt}$  and  $R_t^{(3)} = \left( t \sqrt{s \frac{C_t}{C_{\min}}} \right)^{\frac{2}{3}}$  and note that they correspond to the regret  
 791 bounds obtained when using the respective learning rates  $\lambda_t^{(2)}$  and  $\lambda_t^{(3)}$ , as per Lemma 8.

792 We now consider the last time that the learning rates  $\lambda_t^{(3)}$  and  $\lambda_t^{(2)}$  have been used. More specifically,  
 793 we denote  $T_2 = \max\{t \leq T, \lambda_{t-1}^{(2)} \geq \lambda_{t-1}^{(3)}\}$ , and  $T_3 = \max\{t \leq T, \lambda_{t-1}^{(3)} \geq \lambda_{t-1}^{(2)}\}$ . Combining the  
 794 bound of Equation 12 and using the definition  $\lambda_t = \min(\frac{1}{2}, \max(\lambda_t^{(2)}, \lambda_t^{(3)}))$ , the following bound

795 holds

$$\begin{aligned}
& R_T \\
& \leq \mathbb{E} \left[ \frac{C_T}{\lambda_{T-1}} + \sum_{t=1}^T \min \left( \frac{32}{3} \lambda_{t-1} \overline{\mathbf{R}}_t^{(2)}(\pi_t), \frac{16}{3} c_3^* \sqrt{3 \lambda_{t-1} \overline{\mathbf{R}}_t^{(3)}(\pi_t)} \right) \right] \\
& = \mathbb{E} \left[ \frac{C_T}{\min(\frac{1}{2}, \max(\lambda_{T-1}^{(2)}, \lambda_{T-1}^{(3)}))} \right. \\
& \quad \left. + \sum_{t=1}^T \min \left( \frac{32}{3} \min(\frac{1}{2}, \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)})) \overline{\mathbf{R}}_t^{(2)}(\pi_t), \frac{16}{3} c_3^* \sqrt{3 \min(\frac{1}{2}, \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)})) \overline{\mathbf{R}}_t^{(3)}(\pi_t)} \right) \right] \\
& \leq \mathbb{E} \left[ C_T \min \left( \frac{1}{\lambda_{T-1}^{(2)}}, \frac{1}{\lambda_{T-1}^{(3)}} \right) + \sum_{t=1}^T \min \left( \frac{32}{3} \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)}) \overline{\mathbf{R}}_t^{(2)}(\pi_t), \frac{16}{3} c_3^* \sqrt{3 \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)}) \overline{\mathbf{R}}_t^{(3)}(\pi_t)} \right) \right].
\end{aligned}$$

796 We can now separate the sum obtained at the last line based on which learning rate was used at time  
797 t.

$$\begin{aligned}
& \sum_{t=1}^T \min \left( \frac{32}{3} \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)}) \overline{\mathbf{R}}_t^{(2)}(\pi_t), \frac{16}{3} c_3^* \sqrt{3 \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)}) \overline{\mathbf{R}}_t^{(3)}(\pi_t)} \right) \\
& \leq \sum_{\lambda_t^{(2)} \geq \lambda_t^{(3)}} \frac{32}{3} \lambda_{t-1}^{(2)} \overline{\mathbf{R}}_t^{(2)}(\pi_t) + \sum_{\lambda_t^{(3)} \geq \lambda_t^{(2)}} \frac{16}{3} c_3^* \sqrt{3 \lambda_{t-1}^{(3)} \overline{\mathbf{R}}_t^{(3)}(\pi_t)} \\
& \leq \sum_{t=1}^{T_2} \frac{32}{3} \lambda_{t-1}^{(2)} \overline{\mathbf{R}}_t^{(2)}(\pi_t) + \sum_{t=1}^{T_3} \frac{16}{3} c_3^* \sqrt{3 \lambda_{t-1}^{(3)} \overline{\mathbf{R}}_t^{(3)}(\pi_t)}.
\end{aligned}$$

798 Following exactly the same step as in the proof of Lemma 8, we further bound  
799  $\sum_{t=1}^{T_2} \frac{32}{3} \lambda_{t-1}^{(2)} \overline{\mathbf{R}}_t^{(2)}(\pi_t) \leq 8 \sqrt{\frac{2}{3}} R_{T_2}^{(2)}$  and  $\sum_{t=1}^{T_3} \frac{16}{3} c_3^* \sqrt{3 \lambda_{t-1}^{(3)} \overline{\mathbf{R}}_t^{(3)}(\pi_t)} \leq 8 \cdot 6^{\frac{1}{3}} R_{T_3}^{(3)}$ .

800 The crucial observation is that which of  $\lambda_T^{(3)}$  or  $\lambda_T^{(2)}$  is bigger will determine whether  $R_T^{(2)}$  or  $R_T^{(3)}$   
801 is the term of leading order (up to some constants). More specifically, Let  $T$  be such that  $\lambda_{T-1}^{(2)} \geq$   
802  $\lambda_{T-1}^{(3)}$  which means that  $\sqrt{\frac{3C_T}{128dT}} \geq \frac{1}{4 \cdot 6^{\frac{1}{3}}} \left( \frac{C_T \sqrt{C_{\min}}}{T \sqrt{s}} \right)^{\frac{2}{3}}$ . Rearranging, this implies that  $\sqrt{C_T d T} \leq$   
803  $\frac{6^{\frac{5}{6}}}{4} \left( T \sqrt{s \frac{C_T}{C_{\min}}} \right)^{\frac{2}{3}}$ , which means that  $R_T^{(2)} \leq \frac{6^{\frac{5}{6}}}{4} R_T^{(3)}$ . Following the exact same steps, we also  
804 have that  $\lambda_{T-1}^{(3)} \geq \lambda_{T-1}^{(2)}$  implies that  $R_T^{(3)} \leq \frac{4}{6^{\frac{5}{6}}} R_T^{(2)}$ . We apply this to the time  $T_2$  in which  
805  $\lambda_{T_2-1}^{(2)} \geq \lambda_{T_2-1}^{(3)}$  by definition. we have that  $R_{T_2}^{(2)} \leq \frac{6^{\frac{5}{6}}}{4} R_{T_2}^{(3)}$  and putting this together with the  
806 previous bound, we have

$$\begin{aligned}
R_T & \leq \frac{C_T}{\lambda_{T-1}^{(3)}} + 8 \sqrt{\frac{2}{3}} R_{T_2}^{(2)} + 8 \cdot 6^{\frac{1}{3}} R_{T_3}^{(3)} \\
& \leq 4 \cdot 6^{\frac{1}{3}} R_T^{(3)} + 8 \sqrt{\frac{2}{3}} \cdot \frac{6^{\frac{5}{6}}}{4} R_{T_2}^{(2)} + 8 \cdot 6^{\frac{1}{3}} R_{T_3}^{(3)} \\
& \leq 4 \cdot 6^{\frac{1}{3}} R_T^{(3)} + 4 \cdot 6^{\frac{1}{3}} R_{T_2}^{(3)} + 8 \cdot 6^{\frac{1}{3}} R_{T_3}^{(3)} \\
& \leq 4 \cdot 6^{\frac{1}{3}} R_T^{(3)} + 4 \cdot 6^{\frac{1}{3}} R_T^{(3)} + 8 \cdot 6^{\frac{1}{3}} R_{T_3}^{(3)} \\
& \leq 16 \cdot 6^{\frac{1}{3}} R_T^{(3)},
\end{aligned}$$

807 where we use the fact that  $T \rightarrow R_T^{(3)}$  is increasing and  $T_2 \leq T, T_3 \leq T$ .

808 Using the same argument as before, we have that  $\lambda_{T_3-1}^{(3)} \geq \lambda_{T_3-1}^{(2)}$ , and we can conclude that  $R_{T_3}^{(3)} \leq$   
809  $\frac{4}{6^{\frac{5}{6}}} R_{T_3}^{(2)}$ .

810 Putting this together, with the previous bound, we have

$$\begin{aligned}
R_T &\leq \frac{C_T}{\lambda_{T-1}^{(2)}} + 8\sqrt{\frac{2}{3}}R_{T_2}^{(2)} + 8 \cdot 6^{\frac{1}{3}}R_{T_3}^{(3)} \\
&\leq 8\sqrt{\frac{2}{3}}R_T^{(2)} + 8\sqrt{\frac{2}{3}}R_{T_2}^{(2)} + 8 \cdot 6^{\frac{1}{3}} \cdot \frac{4}{6^{\frac{5}{6}}}R_{T_3}^{(3)} \\
&\leq 8\sqrt{\frac{2}{3}}R_T^{(2)} + 8\sqrt{\frac{2}{3}}R_{T_2}^{(2)} + 16\sqrt{\frac{2}{3}}R_{T_3}^{(2)} \\
&\leq 8\sqrt{\frac{2}{3}}R_T^{(2)} + 8\sqrt{\frac{2}{3}}R_T^{(2)} + 16\sqrt{\frac{2}{3}}R_T^{(2)} \\
&\leq 32\sqrt{\frac{2}{3}}R_T^{(2)},
\end{aligned}$$

811 where we use the fact that  $T \rightarrow R_T^{(3)}$  is increasing and  $T_2 \leq T, T_3 \leq T$ . Evaluating the constants  
812 numerically yields  $16 \cdot 6^{\frac{1}{3}} \approx 29.07 \leq 30$  and  $32\sqrt{\frac{2}{3}} \approx 26.13 \leq 27$ .

### 813 G.5 Upper bound on the learning rates

814 We now consider the case where the learning rates exceed  $\frac{1}{2}$ , and show that this only holds for small  
815 values of  $T$ . First, we have that  $\lambda_{T-1}^{(2)} \leq \frac{1}{2}$  if

$$\sqrt{\frac{3C_T}{128dT}} \leq \frac{1}{2}.$$

816 Rearranging the inequality and recalling  $C_T = 5 + 2s \log \frac{edT}{s}$ , this is equivalent to

$$T \geq \frac{15}{32d} + \frac{3s}{16d} \log \frac{edT}{s}.$$

817 Using the loose inequality  $\log \frac{edT}{s} \leq \frac{dT}{s}$ , we get that this condition is satisfied for any  $T \geq 1$ .

818 Similarly, we have that  $\lambda_{T-1}^{(3)} \leq \frac{1}{2}$  if

$$\frac{1}{4 \cdot 6^{\frac{1}{3}}} \left( \frac{C_T \sqrt{C_{\min}}}{T \sqrt{s}} \right)^{\frac{2}{3}} \leq \frac{1}{2}.$$

819 We note that

$$C_{\min} = \max_{\mu \in \Delta(A)} \sigma_{\min}(\mathbb{E}_{A \sim \mu} [AA^T]) \leq \max_{\mu \in \Delta(A)} \frac{\text{Tr}(\mathbb{E}_{A \sim \mu} [AA^T])}{d} \leq 1,$$

820 where the first inequality uses that the trace of a matrix is always bigger than  $d$ -times its smallest  
821 eigenvalue and the second inequality uses the fact that for any matrix  $A$ , we have  $\text{Tr}(AA^T) =$   
822  $\sum_{i=1}^d a_i^2 \leq d \max_i |a_i| \leq d$  because we assumed that all the actions are bounded in infinity norm.  
823 Hence the previous inequality will be satisfied if

$$\frac{1}{4 \cdot 6^{\frac{1}{3}}} \left( \frac{C_T}{T \sqrt{s}} \right)^{\frac{2}{3}} \leq \frac{1}{2}.$$

824 Rearranging the inequality, this is equivalent to

$$T \geq 4\sqrt{\frac{3}{s}}C_t = 8\sqrt{3s} \log(eT) + \sqrt{3s} \left( \frac{20}{s} + 8 \log \frac{d}{s} \right).$$

825 Applying Lemma 24 with  $a = 8\sqrt{3s}$  and  $b = \sqrt{3s} \left( \frac{20}{s} + 8 \log \left( \frac{d}{s} \right) \right)$ , we find that the previous  
826 inequality is satisfied for all

$$T \geq 2a \log ea + 2b = 40\sqrt{\frac{3}{s}} + 16\sqrt{3s} \log \frac{8e\sqrt{3d}}{\sqrt{s}}.$$

827 Thus, letting  $T_{\min} = 40\sqrt{\frac{3}{s}} + 16\sqrt{3s} \log \frac{8e\sqrt{3d}}{\sqrt{s}}$  be the constant given above, both learning rates  
828 stay upper bounded by  $\frac{1}{2}$  for all  $T \geq T_{\min}$  and the upper bound on the regret given the previous  
829 subsection holds. Otherwise, we upper bound the instantaneous regret by 2 and this leads to an  
830 additional  $2T_{\min} = \mathcal{O}(\sqrt{s} \log \frac{d}{\sqrt{s}})$  in the regret. Putting this together with the bound proved in the  
831 previous section, we thus have that the following regret bound is valid for any  $T \geq 1$ :

$$R_T \leq \min \left( 27\sqrt{\left(5 + 2s \log \frac{edT}{s}\right) dT}, 30 \left(5 + 2s \log \frac{edT}{s}\right)^{\frac{1}{3}} \left(\frac{T\sqrt{s}}{\sqrt{C_{\min}}}\right)^{\frac{2}{3}} \right) + \mathcal{O}\left(\sqrt{s} \log \frac{d}{\sqrt{s}}\right).$$

832 This concludes the proof of Theorem 2.  $\square$

## 833 I Technical Results

834 In this section, we state and prove the remaining technical results.

835 **Lemma 21.** *Let  $\pi \in \Delta(\mathcal{A})$ , the function  $\theta \rightarrow \Delta(\pi, \theta)$  is 2-Lipschitz with respect to the 1 norm. Let*  
836  *$t \geq 1$ , the function  $\theta \rightarrow \mathbb{E} \left[ \log \left( \frac{1}{p_t(Y_t|\theta, A_t)} \right) \right]$  is 2-Lipschitz with respect to the 1 norm.*

837 *Proof.* Let  $\theta, \theta' \in \Theta$ , we have

$$\begin{aligned} |r(\pi, \theta) - r(\pi, \theta')| &= \left| \sum_{a \in \mathcal{A}} \pi(a) \langle \theta - \theta', a \rangle \right| \\ &\leq \sum_{a \in \mathcal{A}} \pi(a) |\langle \theta - \theta', a \rangle| \\ &\leq \sum_{a \in \mathcal{A}} \pi(a) \|\theta - \theta'\|_1 \|a\|_{\infty} \\ &\leq \|\theta - \theta'\|_1. \end{aligned}$$

838 Similarly,

$$|r^*(\theta) - r^*(\theta')| = \left| \max_{a \in \mathcal{A}} r(a, \theta) - \max_{a \in \mathcal{A}} r(a, \theta') \right| \leq \max_{a \in \mathcal{A}} |r(a, \theta) - r(a, \theta')| \leq \|\theta - \theta'\|_1.$$

839 Finally

$$|\Delta(\pi, \theta) - \Delta(\pi, \theta')| = |r^*(\theta) - r^*(\theta') + r(\pi, \theta') - r(\pi, \theta)| \leq 2\|\theta - \theta'\|_1.$$

840 For the negative log-likelihood, for simplicity, we let  $r = \langle \theta, A_t \rangle$ ,  $r' = \langle \theta', A_t \rangle$  and  $r_0 = \langle \theta_0, A_t \rangle$ ,

$$\begin{aligned} \mathbb{E} \left[ \log \left( \frac{1}{p(Y_t|\theta, A_t)} \right) - \log \left( \frac{1}{p(Y_t|\theta', A_t)} \right) \right] &= \frac{1}{2} \mathbb{E} [(\langle \theta, A_t \rangle - Y_t)^2 - (\langle \theta', A_t \rangle - Y_t)^2] \\ &= \frac{1}{2} \mathbb{E} [(r - Y_t)^2 - (r' - Y_t)^2] \\ &= \frac{1}{2} \mathbb{E} [(r - r')(r + r' - 2Y_t)] \\ &= \frac{1}{2} \mathbb{E} [(r - r')(r + r' - 2r_0)] \\ &\leq 2\|\theta - \theta'\|_1. \end{aligned}$$

841  $\square$

842 **Lemma 22.** (Hoeffding's Lemma) *Let  $X$  be a bounded real random variable such that  $X \in [a, b]$*   
843 *almost surely. Let  $\eta \neq 0$ , then we have*

$$\frac{1}{\eta} \log \mathbb{E} [\exp(\eta X)] \leq \mathbb{E}[X] + \frac{\eta(b-a)^2}{8}. \quad (33)$$

844 *Proof.* See for instance Chapter 2 in [Boucheron et al. \[2013\]](#).  $\square$

845 We now provide a data dependent version of Hoeffding's lemma that is used in the analysis of the  
 846 gaps in the optimistic posterior.

847 **Lemma 23.** (A data dependent version of Hoeffding's Lemma) Let  $X$  be a real random variable  
 848 and  $\eta \neq 0$  be such that  $\eta X \leq 1$  almost surely, then we have

$$\frac{1}{\eta} \log \mathbb{E} [\exp (\eta X)] \leq \mathbb{E} [X] + \eta \mathbb{E} \left[ X^2 \right] \leq 2 \mathbb{E} [X]. \quad (34)$$

849 *Proof.* Using the elementary inequalities  $\log (x) \leq x - 1$  for  $x > 0$  and  $e^x \leq 1 + x + x^2$  for  $x \leq 1$ ,  
 850 we get that

$$\begin{aligned} \frac{1}{\eta} \log \mathbb{E} [\exp (\eta X)] &\leq \frac{1}{\eta} \mathbb{E} [\exp (\eta X) - 1] \\ &\leq \frac{1}{\eta} \mathbb{E} [\eta X + \eta^2 X^2] \\ &\leq \mathbb{E} [X] + \eta \mathbb{E} \left[ X^2 \right]. \end{aligned}$$

851

□

852 The following lemmas help us to analyze when the learning rates are smaller or bigger than  $\frac{1}{2}$ .

853 **Lemma 24.** Let  $a \geq 1, b \geq 0$ , then, the equation  $t \geq a \log et + b$  is verified for any  $t \geq 2a \log ea + 2b$   
 854 .

855 *Proof.* We let  $f(t) = t - a \log et - b$ , we have that  $f'(t) \geq 0$  on  $[a, +\infty)$  and  $f(a) \leq 0$ . Hence  
 856  $f(t) = 0$  has a unique solution  $\alpha$  on  $[a, \infty)$  such that  $f(t) \geq 0$  if  $t \geq \alpha$ . We now focus on upper  
 857 bounding  $\alpha$ . The equation  $f(\alpha) = 0$  is equivalent to

$$\log \alpha = \frac{\alpha - b}{a} - 1.$$

858 Now taking the exponential and reordering this is also equivalent to

$$\frac{-\alpha}{a} \exp \left( \frac{-\alpha}{a} \right) = \frac{\exp \left( -\frac{a+b}{a} \right)}{a}.$$

859 Let

$$\begin{aligned} g : (-\infty, -1] &\longrightarrow \left[ -\frac{1}{e}, 0 \right) \\ x &\longmapsto xe^x. \end{aligned}$$

860 The previous equation can be rewritten  $g \left( \frac{-\alpha}{a} \right) = -\frac{\exp \left( -\frac{a+b}{a} \right)}{a}$ .

861 We define  $W_{-1} : \left[ -\frac{1}{e}, 0 \right) \longrightarrow (-\infty, -1]$  as the (functional) inverse of  $g$ .  $g$  is the  $-1$  branch of the  
 862 Lambert W function.

863 We have that for any  $x \leq -1$ ,  $W_{-1}(xe^x) = x$  and that for any  $y \geq e$ ,  $-W_{-1}\left(-\frac{1}{y}\right) \leq 2 \log(y)$ .  
 864 Since  $g$  is decreasing on its domain,  $W_{-1}$  is well-defined and decreasing. Moreover, for any  $x \leq -1$

865 ,  $W_{-1}(g(x)) = x$ . In particular, we have that  $\alpha = aW_{-1} \left( -\frac{\exp \left( -\frac{a+b}{a} \right)}{a} \right)$ . We will use that

866 formulation to find an upper bound on  $\alpha$ .

867 We fix some  $y \geq e$ . We have  $-2 \log(y) \leq -1$  hence  $W_{-1} \left( -2 \log(y) e^{-2 \log(y)} \right) = -2 \log(y)$ ,  
 868 which means that  $2 \log(y) = -W_{-1}\left(-\frac{1}{y^*}\right)$  where  $y^* = \frac{e^{(2 \log(y))}}{2 \log(y)} = \frac{y^2}{2 \log(y)}$ .

869 Because of the elementary inequality  $2 \log(x) \leq x$  for  $x > 0$ , we conclude that  $y \leq y^*$ . Since  
 870  $y \longrightarrow -W_{-1}\left(-\frac{1}{y}\right)$  is an increasing function we finally have that for any  $y \geq e$

$$W_{-1} \left( -\frac{1}{y} \right) \leq W_{-1} \left( -\frac{1}{y^*} \right) = 2 \log(y).$$

871 Applying this to  $y = a \exp\left(\frac{a+b}{a}\right) \geq e$ , we get

$$\alpha = W_{-1}\left(\frac{-1}{y}\right) \leq 2 \log(y) = 2a \log ea + 2b.$$

872 Since any  $t \geq \alpha$  will satisfy  $f(t) \geq 0$ , this concludes our proof. □

873

874 **Lemma 25.** Let  $\theta \in \Theta$ , then  $M_t = \exp(L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta)) = \prod_{s=1}^t \frac{p(Y_t|\theta, A_t)}{p(Y_t|\theta_0, A_t)}$  is a supermartin-  
875 gale with respect to the filtration  $\mathcal{F}_t$ .

876 *Proof.* We have

$$\begin{aligned} \mathbb{E}\left[\frac{p(Y_t|\theta, A_t)}{p(Y_t|\theta_0, A_t)} \middle| \mathcal{F}_{t-1}, A_t\right] &= \mathbb{E}\left[\exp\left(\frac{\langle \theta_0, A_t \rangle - Y_t^2 - (\langle \theta, A_t \rangle - Y_t^2)}{2}\right) \middle| \mathcal{F}_{t-1}, A_t\right] \\ &= \mathbb{E}\left[\exp\left(\frac{\epsilon_t^2 - (\langle \theta - \theta_0, A_t \rangle - \epsilon_t)^2}{2}\right) \middle| \mathcal{F}_{t-1}, A_t\right] \\ &= \exp\left(-\frac{\langle \theta - \theta_0, A_t \rangle^2}{2}\right) \mathbb{E}[\exp(\epsilon_t \langle \theta - \theta_0, A_t \rangle) | \mathcal{F}_{t-1}, A_t] \\ &\leq \exp\left(-\frac{\langle \theta - \theta_0, A_t \rangle^2}{2}\right) \cdot \exp\left(\frac{\langle \theta - \theta_0, A_t \rangle^2}{2}\right) \\ &= 1, \end{aligned}$$

877 where the inequality comes from the conditional subgaussianity of  $\epsilon_t$ . Finally, by the tower rule of  
878 conditional expectations

$$\mathbb{E}\left[\frac{p(Y_t|\theta, A_t)}{p(Y_t|\theta_0, A_t)} \middle| \mathcal{F}_{t-1}\right] = \mathbb{E}\left[\mathbb{E}\left[\frac{p(Y_t|\theta, A_t)}{p(Y_t|\theta_0, A_t)} \middle| \mathcal{F}_{t-1}, A_t\right] \middle| \mathcal{F}_{t-1}\right] \leq 1.$$

879 □

### 880 I.1 Proof of Proposition 1

881 This is coming from the fact that the mean is the constant minimizing the mean squared error. We  
882 remind the reader of the definition of the surrogate information gain and the true information gain  
883 for a policy  $\pi \in \Delta(\mathcal{A})$

$$\overline{\text{IG}}_t(\pi) = \frac{1}{2} \sum_{a \in \mathcal{A}} \pi(a) \int_{\Theta} (\langle \theta - \bar{\theta}(Q_t^+), a \rangle)^2 dQ(\theta), \quad (35)$$

884 where  $\bar{\theta}(Q_t^+) = \mathbb{E}_{\theta \sim Q_t^+}[\theta]$  is the mean parameter under the optimistic posterior  $Q_t^+$ .

$$\text{IG}_t(\pi) = \frac{1}{2} \sum_{a \in \mathcal{A}} \pi(a) \int_{\Theta} (\langle \theta, a \rangle - \langle \theta_0, a \rangle)^2 dQ_t^+(\theta), \quad (36)$$

885 Let's fix  $a \in \mathcal{A}$ , we have that

$$\begin{aligned} (\langle \theta - \theta_0, a \rangle)^2 &= (\langle \theta - \bar{\theta}(Q_t^+) + \bar{\theta}(Q_t^+) - \theta_0, a \rangle)^2 \\ &= (\langle \theta - \bar{\theta}(Q_t^+), a \rangle)^2 + 2\langle \theta - \bar{\theta}(Q_t^+), a \rangle \langle \bar{\theta}(Q_t^+) - \theta_0, a \rangle + (\langle \bar{\theta}(Q_t^+) - \theta_0, a \rangle)^2 \\ &\geq (\langle \theta - \bar{\theta}(Q_t^+), a \rangle)^2 + 2\langle \theta - \bar{\theta}(Q_t^+), a \rangle \langle \bar{\theta}(Q_t^+) - \theta_0, a \rangle \end{aligned}$$

886 Now using that  $\bar{\theta}(Q_t^+) = \int_{\Theta} \theta dQ_t^+(\theta)$  and integrating, we get

$$\int_{\Theta} (\langle \theta - \theta_0, a \rangle)^2 dQ_t^+(\theta) \geq \int_{\Theta} (\langle \theta - \bar{\theta}(Q_t^+), a \rangle)^2 dQ_t^+(\theta).$$

887 Multiplying by  $\pi(a)$  and summing over actions, we get the claim of the lemma.

888 **I.2 Generalization of the AM-GM inequality**

889 Dealing with the generalized information ratio requires bounding the cubic root of products. While  
 890 one could use Hölder's inequality to deal directly with products, we find it more flexible to use a  
 891 variational form of this inequality. In all that follows, we let  $p > 1$  be a real number and  $q$  be such  
 892 that  $\frac{1}{p} + \frac{1}{q} = 1$ . It is not hard to check that  $q = \frac{p}{p-1}$ . We start by stating a direct consequence of the  
 893 Fenchel-Young Inequality which can be seen as an extension of the AM-GM inequality.

894 **Lemma 26.** *Let  $x, y \geq 0$ , then*

$$xy \leq \frac{x^p}{p} + \frac{y^q}{q}. \quad (37)$$

895 *With equality if and only if  $px^{p-1} = y$*

896 *Proof.* One can check that the Fenchel dual of the function

$$\begin{aligned} f : \mathbb{R}^+ &\longrightarrow \mathbb{R} \\ x &\longmapsto \frac{x^p}{p} \end{aligned}$$

897 is exactly  $f^*(y) = \frac{1}{q}|y|^q \text{sgn}(y)$ . Then the Lemma is a direct consequence of the Fenchel Young  
 898 inequality and of its equality case. □

899 Refining a bit this Lemma, we get the following variational form of the previous inequality :

900 **Lemma 27.** *Let  $x, y \geq 0, \lambda > 0$ , then*

$$\sqrt[p]{xy} \leq \frac{x}{\lambda} + c_p^*(\lambda y)^{\frac{1}{p-1}} \quad (38)$$

901 *where  $c_p^* = (p-1)\frac{1}{p} \frac{p-1}{p}$  with equality if and only if  $x = y = 0$  or  $\lambda = p \frac{x^{\frac{p-1}{p}}}{y^{\frac{1}{p}}}$ .*

902 *Proof.* We apply the previous lemma to  $\sqrt[p]{\frac{px}{\lambda}}$  and  $\sqrt[p]{\frac{\lambda y}{p}}$ . □

903 In order to go from the variational form to the product form, we may use the following result.

904 **Lemma 28.** *Let  $\alpha, \beta > 0$ , then*

$$\inf_{\lambda > 0} \frac{\alpha}{\lambda} + \beta \lambda^{\frac{1}{p-1}} = c_p \alpha^{\frac{1}{p}} \beta^{\frac{p-1}{p}}, \quad (39)$$

905 *where  $c_p = p \frac{1}{p-1} \frac{p-1}{p}$  satisfies  $c_p \cdot c_p^* \frac{p-1}{p} = 1$ , and the minimum is reached at  $\lambda^* = (p-1) \frac{p-1}{p} \frac{\alpha^{\frac{p-1}{p}}}{\beta^{\frac{p-1}{p}}}$ .*

906 *Proof.* Applying the previous Lemma to  $x = \alpha$  and  $y = c_p^{\frac{p}{p-1}} \beta^{p-1}$  yields the result. □

907 **Remark** An alternative is to pick  $\lambda$  to make both terms equals resulting in the same result but with  
 908 2 as a leading constant. Now

$$\begin{aligned} c_p &= p^{\frac{1}{p}} \frac{p}{p-1} \frac{p-1}{p} \\ &= \exp\left(\frac{1}{p} \log p + \frac{p-1}{p} \log \frac{p}{p-1}\right) \\ &\leq \frac{1}{p} \cdot p + \frac{p-1}{p} \cdot \frac{p}{p-1} \\ &= 2. \end{aligned}$$

909 With equality if and only if  $p = 2$ . So, the choice of  $c_p$  always yields a better leading constant.  
 910 However,  $c_3 \simeq 1.88$  so one could argue that the gain is small. Since we will usually use Lemma 27,  
 911  $c_p^*$  will naturally appear and  $c_p$  will cancel it, ultimately making the leading constant as simple as  
 912 possible.

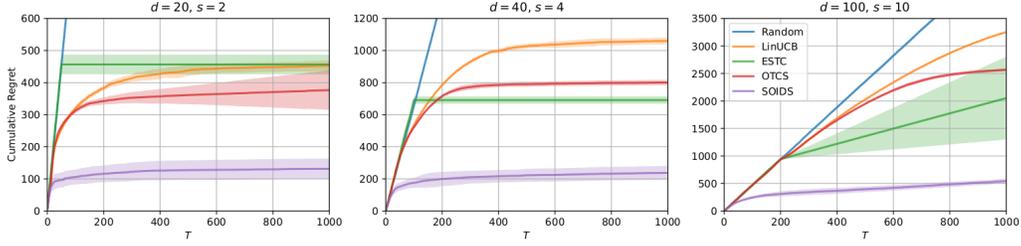


Figure 1: Cumulative regret for  $d = 20$  (left) 40 (middle) and 100 (right). We plot the mean  $\pm$  standard deviation over 10 repetitions.

## 913 J Experimental details

914 We aim to verify that, in both the data-rich and data-poor regimes simultaneously, the regret of  
 915 SOIDS is comparable with the regret of existing algorithms that achieve near optimal worst-case  
 916 regret in either the data-rich or the data-poor regime. Our baseline for the data-rich regime is the  
 917 online-to-confidence-set (OTCS) method proposed by [Abbasi-Yadkori et al. \[2012\]](#), which has worst  
 918 case regret of the order  $\sqrt{sdT}$ . For a tougher comparison, we run this method with the confidence  
 919 sets from Theorem 4.7 of [Clerico et al. \[2025\]](#), which have much smaller constant factors than  
 920 those used by [Abbasi-Yadkori et al. \[2012\]](#). Our baseline for the data-poor regime is the Explore  
 921 the Sparsity Then Commit (ESTC) algorithm proposed by [Hao et al. \[2020\]](#), which has worst-case  
 922 regret of the order  $(sT)^{2/3}$ . For reference, we also compare with LinUCB [Abbasi-Yadkori et al.](#)  
 923 [\[2011\]](#), which does not adapt to sparsity.

924 It is generally difficult to run the SOIDS algorithm exactly because the surrogate information ratio  
 925 contains expectations w.r.t. the optimistic posterior. In our implementation of SOIDS, we use  
 926 the empirical Bayesian sparse sampling procedure of [Hao et al. \[2021\]](#) to draw approximate sam-  
 927 ples from the optimistic posterior, and then approximate the surrogate information ratio via sample  
 928 averages.

929 For each  $d \in \{20, 40, 100\}$ ,  $\theta_0$  is the  $s$ -sparse vector in  $\mathbb{R}^d$ , with  $s = d/10$ , in which first  $s$  com-  
 930 ponents are  $10/s$  and the remaining components are zero. The action set consists of 200 random  
 931 draws from the uniform distribution on  $[-1, 1]^d$ . The noise variance is 1 and we run each method  
 932 10 times. In Figure 1, we report the cumulative regret over  $T = 1000$  steps. As  $d$  is varied from 20  
 933 to 100, we appear to transition from the data-rich regime to the data-poor regime: for  $d = 20$ , the  
 934 OTCS method is the best performing baseline, whereas for  $d = 100$ , ETCS is the best performing  
 935 baseline. As our theoretical results would suggest, SOIDS performs well in both regimes.

936 To run the SOIDS algorithm, one must minimise  $\bar{\mathbb{R}}_t^{(2)}(\pi)$  w.r.t.  $\pi$  in each round  $t$ . This is not  
 937 straightforward, because  $\bar{\mathbb{R}}_t^{(2)}(\pi)$  contains expectations w.r.t. the optimistic posterior  $Q_t^+$ . When  
 938 we use the Spike-and-Slab prior in Appendix B.2, we are not aware of any efficient method that can  
 939 be used to maximise  $\bar{\mathbb{R}}_t^{(2)}(\pi)$ . Instead, we draw (approximate) samples  $\theta^{(1)}, \dots, \theta^{(M)}$  from  $Q_t^+$   
 940 to produce the estimates  $\tilde{\Delta}_t(\pi)$  and  $\tilde{\mathbb{I}}_t(\pi)$  for the surrogate regret and the surrogate information  
 941 ratio respectively, where

$$\tilde{\Delta}_t(\pi) = \sum_{a \in \mathcal{A}} \pi(a) \frac{1}{M} \sum_{i=1}^M \Delta(a, \theta^{(i)}), \quad \tilde{\mathbb{I}}_t(\pi) = \frac{1}{2} \sum_{a \in \mathcal{A}} \pi(a) \frac{1}{M} \sum_{i=1}^M ((\theta^{(i)} - \bar{\theta}_M, a))^2.$$

942 Here,  $\bar{\theta}_M$  is the sample mean  $\frac{1}{M} \sum_{i=1}^M \theta^{(i)}$ . We then maximise the approximate surrogate infor-  
 943 mation ratio  $\tilde{\mathbb{R}}_t^{(2)}(\pi)$ , where

$$\tilde{\mathbb{R}}_t^{(2)}(\pi) = \frac{(\tilde{\Delta}_t(\pi))^2}{\tilde{\mathbb{I}}_t(\pi)}.$$

944 To draw the samples  $\theta^{(1)}, \dots, \theta^{(M)}$ , we use the empirical Bayesian sparse sampling procedure pro-  
 945 posed by [Hao et al. \[2021\]](#), which is designed to draw samples from the Bayesian posterior. To  
 946 sample from the optimistic posterior, we incorporate the optimistic adjustment into the likelihood.

947 This method replaces the theoretically sound spike-and-slab prior with a relaxation in which the  
 948 “spikes” are Laplace distributions with small variance, and the “slabs” are Gaussian distributions  
 949 with large variance. In particular, the density of this prior is

$$q_1(\theta) = \sum_{\gamma \in \{0,1\}^d} p(\gamma) \prod_{j=1}^d [\gamma_j \psi_1(\theta_j) + (1 - \gamma_j) \psi_0(\theta_j)].$$

950 Here,  $\psi_1(\theta)$  is the density function of a univariate Gaussian distribution, with mean 0 and vari-  
 951 ance  $\rho_1$ , and  $\psi_0$  is the density function of a univariate Laplace distribution, with mean 0 and scale  
 952 parameter  $\rho_0$ .  $p(\gamma)$  is a product of Bernoulli distributions with mean  $\beta$ . In our experiments, we  
 953 always use  $\rho_1 = 10$ ,  $\rho_0 = 0.1$  and  $\beta = 0.1$ . Also, we set the learning rates to  $\eta = 1/2$  and

$$954 \lambda_t = \min\left(\frac{1}{2}, \frac{1}{10} \max\left(\sqrt{\frac{s \log(edt/s)}{dt}}, \left(\frac{\log(edt/s)}{t}\right)^{2/3}\right)\right).$$

955 Implementing the OTCS baseline exactly would require us to compute the means of the distributions  
 956 played by an exponentially weighted average forecaster with a sparsity prior. These distributions are  
 957 the same as the optimistic posterior, except  $\lambda_t = 0$  (i.e. there is no optimistic adjustment). In our  
 958 implementation of the OTCS baseline, we draw samples using the same empirical Bayesian sparse  
 959 sampling procedure, and then replace the exact means with the sample means. We use the same  
 960 choices for the parameters  $\eta$ ,  $\rho_1$ ,  $\rho_0$  and  $\beta$ . We set the radii of the confidence sets to the values given  
 961 in Theorem 4.7 of Clerico et al. [2025]

962 For the LinUCB baseline, we set the radii of the confidence sets to the values given in Theorem 2 of  
 963 Abbasi-Yadkori et al. [2011]. For the ESTC baseline, we set the exploration length  $T_1$  to 50 when  
 964  $d = 20$ , 100 when  $d = 40$  and  $d = 100$ . These values were chosen based on a small amount of trial  
 965 and error. The theoretically motivated values in Theorem 4.2 of Hao et al. [2020] are much larger  
 966 than these values. Also for ESTC, we set the LASSO regularisation parameter to  $\lambda = 4\sqrt{\log(d)/T_1}$ ,  
 967 which is the value given in Theorem 4.2 of Hao et al. [2020].