Sparse Optimistic Information Directed Sampling

Ludovic Schwartz Hamish Flynn Gergely Neu

Universitat Pompeu Fabra, Barcelona, Spain {ludovic.schwartz, hamishedward.flynn, gergely.neu}@upf.edu

Abstract

Many high-dimensional online decision-making problems can be modeled as stochastic sparse linear bandits. Most existing algorithms are designed to achieve optimal worst-case regret in either the data-rich regime, where polynomial dependence on the ambient dimension is unavoidable, or the data-poor regime, where dimension-independence is possible at the cost of worse dependence on the number of rounds. In contrast, the sparse Information Directed Sampling (IDS) algorithm satisfies a Bayesian regret bound that has the optimal rate in both regimes simultaneously. In this work, we explore the use of Sparse Optimistic Information Directed Sampling (SOIDS) to achieve the same adaptivity in the worst-case setting, without Bayesian assumptions. Through a novel analysis that enables the use of a time-dependent learning rate, we show that SOIDS can optimally balance information and regret. Our results extend the theoretical guarantees of IDS, providing the first algorithm that simultaneously achieves optimal worst-case regret in both the data-rich and data-poor regimes. We empirically demonstrate the good performance of SOIDS.

1 Introduction

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

IDS. These results hold on average, assuming that the problem instance is drawn at random from a known prior distribution.

In this work, we lift this assumption and develop an algorithm that can adapt to both regimes in a "frequentist" sense: we assume that the true parameter is fixed and unknown to the learner, and provide guarantees that hold for any given instance. The algorithm is an adaptation of the recently proposed Optimistic Information Directed Sampling (OIDS) algorithm of Neu, Papini, and Schwartz [2024], which itself is an adaptation of the classic Bayesian IDS algorithm originally proposed by Russo and Van Roy [2017]. Within the Bayesian setting, it has been shown that IDS can exploit various types of problem structure, and adapt to the hardness of the given instance [Hao and Lattimore, 2022, Hao et al., 2022]. These results have been complemented by the recent work of Neu, Papini, and Schwartz [2024], which showed that similar improvements can be achieved without Bayesian assumptions, via a simple adjustment of the standard IDS method. In this paper, we continue this line of work. We show that a sparse version of OIDS enjoys a worst-case regret bound that matches the optimal rate in both regimes simultaneously.

Our contribution is as follows:

- We extend the analysis of the optimistic posterior to allow the use of time-dependent learning rates and history-dependent learning rates. Time-dependent learning rates allow us to drop the requirement that the horizon is known in advance, and are essential for worst-case regret bounds that can adapt to both regimes. History-dependent learning rates allow us to update the learning rate based on data observed by the agent instead of some loose theoretical constant, a necessity for efficient algorithms.
- We demonstrate that the Sparse Optimistic Information Directed Sampling (SOIDS) algorithm recovers almost optimal rates in both the data-poor and data-rich regimes. This is the first algorithm to do so in a frequentist setting.

2 Preliminaries

Sparse linear bandits. We consider the following decision-making game, in which a learning agent interacts with an environment over a sequence of T rounds. At the start of each round t, the 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, the environment generates a stochastic reward $Y_t = r(A_t) + \epsilon_t$, where $r : \mathcal{A} \to \mathbb{R}$ is a fixed reward function and ϵ_t is zero-mean, conditionally 1-sub-Gaussian noise. We assume that the action set \mathcal{A} is finite, and that the reward function can be written as

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

where $\theta_0 \in \mathbb{R}^d$ is an unknown parameter vector. We make the mild boundedness assumptions 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 parameter vector θ_0 is s-sparse in the sense that at most $s \ll d$ of its components are non-zero. In other words, we assume that θ_0 belongs to the following sparse parameter space:

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

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

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

where the expectation is taken with respect to both the random choices of the agent and the random noise in the observed rewards. We note that the regret is implicitly a function of the true parameter θ_0 . Our focus is on proving regret bounds that hold for arbitrary choices of $\theta_0 \in \Theta$.

The data-rich and data-poor regimes. As mentioned in the introduction, it is known there exist algorithms for sparse linear bandits with worst-case regret of the order $\mathcal{O}(\sqrt{sdT})$ [Abbasi-Yadkori et al., 2012]. This regret bound is only meaningful when the dimension d is smaller than the number of rounds T, a situation referred to as the data-rich regime. Under the assumption that there exists an exploratory policy, Hao et al. [2020] showed that there is a simple algorithm that satisfies a

problem-dependent regret bound, which can be meaningful in the so-called data-poor regime, where d is much larger than T. Formally, we say that there exists an exploratory policy if the action set $\mathcal A$ is such that

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

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

Adapting to both regimes. Recently, Hao et al. [2021] showed that the sparse Information Directed Sampling (IDS) algorithm performs well in both regimes. Under the sparse optimal action condition (Definition 1), IDS satisfies a regret bound of the order $\mathcal{O}(\min(\sqrt{dT\Delta},(sT)^{2/3}\Delta^{1/3}C_{\min}^{-1/3}))$, where $\Delta \propto \min(\log(|\mathcal{A}|),s\log(dT/s))$. This is simultaneously optimal in both the data-rich and data-poor regimes. However, this result is limited to the Bayesian setting. This is because IDS uses the Bayesian posterior to quantify uncertainty, which is only meaningful if θ_0 really is a random draw from the prior.

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

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

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

Notation. We conclude this section by introducing some additional notation that will be used in the subsequent sections. For any candidate parameter vector (or model) $\theta \in \mathbb{R}^d$, we let $r(a,\theta) = \langle \theta,a \rangle$ denote the corresponding linear reward function. In addition, we define $a^*(\theta) = \arg\max_{a \in \mathcal{A}} r(a,\theta)$ (with ties broken arbitrarily) and $r^*(\theta) = r(a^*(\theta),\theta)$ to be the optimal action and maximum reward for the model θ . The gap of an action a for a model θ is $\Delta(a,\theta) = r^*(\theta) - r(a,\theta)$. Similarly, the gap for a policy $\pi \in \Delta(\mathcal{A})$ and a model distribution $Q \in \Delta(\Theta)$ is $\Delta(\pi,Q) = \int_{\mathcal{A} \times \Theta} \Delta(a,\theta) \, d\pi \otimes Q(a,\theta)$, and we let $\Delta_t = \Delta(\pi_t,\theta_0)$ denote the gap of the policy played by the agent in round t under the true model θ_0 . Using this notation, the regret can be written as $R_T = \mathbb{E}[\sum_{t=1}^T \Delta_t]$. We define the unnormalized Gaussian likelihood function $p(y|\theta,a) = \exp(-\frac{(y-\langle \theta,a\rangle)^2}{2})$. Finally, we let $\mathcal{F}_t = \sigma(A_1,Y_1,\ldots,A_t,Y_t)$ denote the σ -algebra generated by the interaction between the agent and the environment up to the end of round t.

3 Sparse Optimistic Information Directed Sampling

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

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

¹The 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.

Here, η is a positive constant that should be thought of as "large", and $(\lambda_t)_t$ is a decreasing sequence of positive real numbers that decays to 0, and should be thought of as "small". We allow λ_t to be computed by the algorithm at the end of the round t. In other words, any \mathcal{F}_t -measurable λ_t is admissible. Note that when $\eta = 1$ and $\lambda_t = 0$, the optimistic posterior coincides with the Bayesian posterior. While this construction is closely related to the optimistic posterior update described in Zhang [2022] and Neu, Papini, and Schwartz [2024], there are a few important differences. First, the $\Delta(A_s, \theta)$ term appearing in the adjustment serves as an alternative to their proposal of using $r^*(\theta)$ for the same purpose. Intuitively this serves to "overestimate" the true gaps with the optimistic posterior, driving exploration towards parameters that promise rewards much higher than whatever would have been accrued by the agent. In contrast, the adjustment of Zhang [2022] drives exploration towards parameters θ with high optimal reward regardless of how well the agent would have performed under the same θ —meaning that it unduly assigns mass to uninteresting parameter choices, where any policy is guaranteed to work well anyway. Intuition aside, this adjustment greatly simplifies our analysis of the optimistic posterior as compared to the analysis of Zhang [2022] and Neu, Papini, and Schwartz [2024]. An important additional novelty is that our update features a time-dependent exploration parameter λ_t , which is crucial for the adaptive regret bounds that we seek in this work. To describe the OIDS algorithm, we must first define the surrogate information gain and the surrogate regret. For any round t and any policy $\pi \in \Delta(\mathcal{A})$, the surrogate information gain is defined as

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

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

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

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

$$\overline{\operatorname{IR}}_{t}^{(\gamma)}(\pi) = \frac{(\widehat{\Delta}_{t}(\pi))^{\gamma}}{\overline{\operatorname{IG}}_{t}(\pi)} = 2 \cdot \frac{\left(\sum_{a \in \mathcal{A}} \pi(a) \int_{\Theta} \langle \theta, a^{*}(\theta) - a \rangle dQ_{t}^{+}(\theta)\right)^{\gamma}}{\sum_{a \in \mathcal{A}} \pi(a) \int_{\Theta} (\langle \theta - \overline{\theta}(Q_{t}^{+}), a \rangle)^{2} dQ_{t}^{+}(\theta)}.$$
(3)

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

$$\pi_t^{(\mathbf{SOIDS})} = \mathop{\arg\min}_{\pi \in \Delta(\mathcal{A})} \overline{\mathrm{IR}}_t^{(2)}(\pi) \,. \tag{4}$$

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

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

$$\overline{IR}_t^{(\gamma)}(\pi_t^{(\mathbf{SOIDS})}) \le 2^{\gamma-2} \min_{\pi \in \Delta(\mathcal{A})} \overline{IR}_t^{(\gamma)}(\pi) \,.$$

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

4 Main results

In this section, we state our main results. First, we relate the true regret of any policy sequence to the surrogate regret of the same policy sequence. We then use the fact that the surrogate regret is controlled by both the 2- and 3-information ratios. This, combined with Lemma 1, allows us to show that with properly tuned parameters, SOIDS has optimal worst-case regret in both the data-poor and data-rich regimes. Finally, we show that SOIDS can be tuned in a data-dependent manner, such that its regret bound scales with the cumulative observed information ratio instead of the time horizon.

4.1 General bound for the optimistic posterior

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

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

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

Otherwise, if the learning rates depend on the history, let $C_{1,T}$ be a deterministic upper bound on $\frac{1}{\lambda_t} - \frac{1}{\lambda_{t-1}}$ 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 sequence of policies π_t satisfies

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

4.2 Adapting to both regimes

We show that the SOIDS algorithm with properly tuned parameters achieves the optimal regret rate in both the data-rich and data-poor regimes.

Theorem 2. Assume that the sparse optimal action condition in Definition 1 is satisfied. Let $\lambda_t^{(2)} = \sqrt{\frac{3C_{t+1}}{128d(t+1)}}$ and $\lambda_t^{(3)} = \frac{1}{4\cdot6^{\frac{1}{3}}}\left(\frac{C_{t+1}\sqrt{C_{\min}}}{(t+1)\sqrt{s}}\right)^{\frac{2}{3}}$, with $C_t = 5 + 2s\log\frac{edt}{s}$. If $\lambda_t = \min(\frac{1}{2}, \max(\lambda_t^{(2)}, \lambda_t^{(3)}))$, then the regret of SOIDS satisfies

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

$$= \min\left(\mathcal{O}\left(\sqrt{sdT\log\frac{edT}{s}}\right), \mathcal{O}\left((sT)^{\frac{2}{3}}\left(\log\frac{edT}{s}\right)^{\frac{1}{3}}\right)\right),$$
(7)

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

We observe that our algorithm enjoys both the $\mathcal{O}(\sqrt{sdT})$ and the $\mathcal{O}((sT)^{\frac{2}{3}})$ regret rates. Unlike the Bayesian regret bound for the sparse IDS algorithm of Hao et al. [2021], our regret bound holds in a "worst-case" sense for any value of $\theta_0 \in \Theta$. To our knowledge, this makes our method the first algorithm to achieve optimal worst-case regret in both the data-poor and data-rich regimes

4.3 Instance-dependent guarantees

The bounds presented in the previous sections are minimax in nature, meaning they hold uniformly over all problem instances. We present a bound in which the scaling with respect to the horizon T is replaced with the cumulative surrogate-information ratio, which could be much smaller than T in "easier" instances, leading to better guarantees.

Theorem 3. Assume that the sparse optimal action condition in Definition 1 is satisfied, and

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

 $\lambda_t = \max(\lambda_t^{(3)}, \lambda_t^{(2)})$, then the regret of SOIDS satisfies

$$R_{T} \leq \left(\frac{2}{s} + \frac{80}{3} + 5\log\frac{edT}{s}\right) \min\left(\sqrt{s\left(2d + \sum_{t=1}^{T-1} \overline{R}_{t}^{(2)}(\pi_{t})\right)}, s^{\frac{1}{3}}\left(\frac{3\sqrt{6}s}{\sqrt{C_{\min}}} + \sum_{t=1}^{T} \sqrt{\overline{R}_{t}^{(3)}(\pi_{t})}\right)^{\frac{2}{3}}\right)$$

$$= \mathcal{O}\left(\log\frac{edT}{s}\min\left(\sqrt{s\left(2d + \sum_{t=1}^{T-1} \overline{R}_{t}^{(2)}(\pi_{t})\right)}, s^{\frac{1}{3}}\left(\frac{3\sqrt{6}s}{\sqrt{C_{\min}}} + \sum_{t=1}^{T} \sqrt{\overline{R}_{t}^{(3)}(\pi_{t})}\right)^{\frac{2}{3}}\right)\right).$$
(8)

This type of result is only possible because our novel analysis of the optimistic posterior (cf. Theorem 1) can handle history-dependent learning rates. A full proof is provided in Appendix D. This result shows that (with appropriate choices of the learning rates) SOIDS is fully adaptive to which of the two regimes is best. Because our analysis requires decreasing learning rates, we are forced to leave the $\log(T)$ terms out of the learning rates, and our logarithmic term has a worse power than in the bound of Theorem 2. An interesting open question is whether it is possible to improve the dependency on this logarithmic term while still using data-dependent learning rates.

5 Analysis

We now provide an outline of the proofs of the main results.

5.1 Proof of Theorem 1

A key observation is that the optimistic posterior can be interpreted as a learner playing an auxiliary online learning game over distributions $\Delta(\Theta)$. The loss of that game is a weighted sum of negative log-likelihood and estimation error losses. We define

$$L_t^{(1)}(\theta) = \sum_{s=1}^t \log\left(\frac{1}{p(Y_s|\theta, A_s)}\right) = \sum_{s=1}^t \frac{1}{2} (\langle \theta, A_s \rangle - Y_s)^2$$

to be the *cumulative negative log-likelihood loss* of θ and

$$L_t^{(2)}(\theta) = \sum_{s=1}^t -\Delta(A_s, \theta)$$

to be the *cumulative estimation error loss* of θ . In addition, we define the regularizer $\Phi: \Delta(\Theta) \to \mathbb{R}$ by the mapping $P \mapsto \mathcal{D}_{KL}\left(P \middle\| Q_1^+\right)$, which is the KL-divergence with respect to the prior Q_1^+ . With those notations, the optimistic posterior corresponds to an instance of the Follow the Regularized Leader (FTRL) algorithm introduced by Hazan and Kale [2010] and Abernethy et al. [2008]. FTRL is a standard method in online convex optimization that balances cumulative loss minimization with a regularization term to enforce stability and guarantee controlled regret. The update can be reframed as

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

This formulation enables the application of tools from convex analysis and online learning, such as Fenchel duality, to derive regret bounds for this auxiliary online learning game and to understand the interplay between the two losses under the learning rates η and λ_t . We now focus on the case in which the learning rates λ_t don't depend on the history and relegate the analysis of history-dependent learning rates to Appendix C. The following lemma provides a bound on the average regret when the model θ_0 is drawn from an arbitrary comparator distribution P.

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

$$\sum_{t=1}^{T} \Delta(A_t, P) \le \frac{\Phi(P) + \Phi^*(\eta(L_T^{(1)}(\theta_T) - L_T^{(1)}(\cdot)) - \lambda_T L_T^{(2)}(\cdot))}{\lambda_T} + \frac{\eta}{\lambda_T} (\langle P, L_T^{(1)} \rangle - L_T^{(1)}(\theta_T)). \tag{9}$$

Here $\theta_t = \arg\min_{\theta \in \Theta} L_t^{(1)}(\theta)$ denotes the maximum likelihood estimator at time t, and $\Phi^*(L) = \log \int_{\Theta} \exp(L(\theta)) dQ_1^+(\theta)$ is the Fenchel dual of the regularizer Φ . A complete proof of this result is provided in appendix B.1.1. We aim to chose a comparator P and the prior Q_1^+ such that P is concentrated around θ_0 and the KL divergence $\mathcal{D}_{\mathrm{KL}}\left(P \middle\| Q_1^+\right)$ is controlled. If the parameter space were finite, the natural choice would be to take P as a Dirac on θ_0 and Q_1^+ as a uniform distribution on the whole parameter space; more care is necessary here. Choosing Q_1^+ as a subset-selection prior and P as a uniform distribution on a sparse neighborhood of θ_0 satisfies both requirements.

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

$$\forall \theta' \in \text{supp}(P_{\theta}), \ \|\theta - \theta'\|_{1} \leq \epsilon \quad and \quad \mathcal{D}_{KL}\left(P_{\theta} \| Q_{1}^{+}\right) \leq s \log \frac{2ed}{\epsilon s}.$$

The proof of this lemma, as well as the exact choice of the prior Q_1^+ and the comparator $P(\theta_0)$, are provided in Appendix B.2. In Appendix I (cf. Lemma 21), we establish that both $L_T^{(2)}(\cdot)$ and $\mathbb{E}\left[L_T^{(1)}(\cdot)\right]$ are 2T-Lipschitz with respect to the ℓ_1 -norm. Hence,

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

Combining these with Lemma 2, we obtain the following bound on the cumulative regret:

$$R_{T} \leq \mathbb{E}\left[\frac{s\log\frac{2ed}{\epsilon s} + 2T(\lambda_{T} + \eta)\epsilon}{\lambda_{T}} + \frac{\Phi^{*}(-\eta(L_{T}^{(1)}(\cdot) - L_{T}^{(1)}(\theta_{T})) - \lambda_{T}L_{T}^{(2)}(\cdot))}{\lambda_{T}}\right] + \mathbb{E}\left[\frac{\eta}{\lambda_{T}}(L_{T}^{(1)}(\theta_{0}) - L_{T}^{(1)}(\theta_{T}))\right].$$

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

Lemma 4.

$$\mathbb{E}\left[\frac{\Phi^{*}(\eta(L_{T}^{(1)}(\theta_{T}) - L_{T}^{(1)}(\cdot)) - \lambda_{T}L_{T}^{(2)}(\cdot))}{\lambda_{T}}\right] \\
\leq \mathbb{E}\left[\sum_{t=1}^{T} \frac{\Phi^{*}(\eta(L_{t}^{(1)}(\theta_{0})L_{t}^{(1)}(\cdot)) - \lambda_{t-1}L_{t}^{(2)}(\cdot))}{\lambda_{t-1}} - \frac{\Phi^{*}(\eta(L_{t-1}^{(1)}(\theta_{0}) - L_{t-1}^{(1)}(\cdot)) - \lambda_{t-1}L_{t-1}^{(2)}(\cdot))}{\lambda_{t-1}}\right] \\
+ \frac{\eta(6 + s\log\frac{edT}{s})}{\lambda_{T}}.$$
(11)

A detailed proof of this result is provided in Appendix B.1.2. Finally, the remaining sum can be handled by looking at the explicit formula for Φ^* . The terms related to the likelihood and the gap estimates can be separated using Hölder's inequality, as is done in Zhang [2022] and Neu, Papini, and Schwartz [2024]. More explicitly, by now choosing $\eta = \frac{1}{4}$, we obtain the following lemma.

Lemma 5

$$\mathbb{E}\left[\sum_{t=1}^{T} \frac{\Phi^{*}(\eta(L_{t}^{(1)}(\theta_{0}) - L_{t}^{(1)}(\cdot)) - \lambda_{t-1}L_{t}^{(2)}(\cdot))}{\lambda_{t-1}} - \frac{\Phi^{*}(\eta(L_{t-1}^{(1)}(\theta_{0}) - L_{t-1}^{(1)}(\cdot)) - \lambda_{t-1}L_{t-1}^{(2)}(\cdot))}{\lambda_{t-1}}\right] \\
\leq \mathbb{E}\left[-\sum_{t=1}^{T} \frac{3\overline{IG}_{t}(\pi_{t})}{32\lambda_{t-1}} + 2\sum_{t=1}^{T} \widehat{\Delta}(\pi_{t})\right]. \tag{12}$$

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

5.2 Proof of Theorem 2

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

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

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

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

This is a consequence of a simple generalization of the AM-GM inequality and is proved in Appendix F.1. Combining the previous lemma with $\lambda = \frac{64}{3}\lambda_{t-1}$ and Theorem 1, we can further upper bound the regret of a sequence of policies $(\pi_t)_t$ as

$$R_{T} \leq \mathbb{E}\left[\frac{5 + 2s\log\frac{edT}{s}}{\lambda_{T-1}} - \sum_{t=1}^{T} \frac{3\overline{\mathrm{IG}}_{t}(\pi_{t})}{32\lambda_{t-1}} + 2\sum_{t=1}^{T} \widehat{\Delta}_{t}(\pi_{t})\right]$$

$$\leq \mathbb{E}\left[\frac{C_{T}}{\lambda_{T-1}} + \sum_{t=1}^{T} \min\left(\frac{32}{3}\lambda_{t-1}\overline{\mathrm{IR}}_{t}^{(2)}(\pi_{t}), \frac{16}{3}c_{3}^{*}\sqrt{3\lambda_{t-1}\overline{\mathrm{IR}}_{t}^{(3)}(\pi_{t})}\right)\right], \tag{13}$$

where $C_T = 5 + 2s \log \frac{edT}{s}$. Usually, bounds on the 2-information ratio can be converted to $\mathcal{O}(\sqrt{T})$ bounds and bounds on the 3-information ratio can be converted to $\mathcal{O}(T^{\frac{2}{3}})$ bounds. Hence we will use the 2-information ratio to control the regret in the data-rich regime and the 3-information ratio to control the regret in the data-poor regime. Due to Lemma 1, the SOIDS policy minimizes the 2-information ratio and approximately minimizes the 3-information ratio. As a result, if there exists a "forerunner" algorithm with bounded 2-information ratio or 3-information ratio, SOIDS inherits these bounds automatically. In particular, we can use a different forerunner for each regime and SOIDS will match the regret guarantees of the best forerunner in each regime.

This forerunner-based technique is widely used to analyze IDS based algorithms and has been applied to a variety of Bayesian settings [Russo and Van Roy, 2017, Hao et al., 2021, Hao and Lattimore, 2022] and some frequentist settings [Kirschner and Krause, 2018, Kirschner et al., 2020, 2021]. An advantage of the OIDS framework is that since the surrogate quantities are defined with respect to the optimistic posterior, the analysis of the surrogate information ratio is virtually identical to the corresponding analysis of the information ratio in the Bayesian setting.

The forerunner we consider for the 2-information ratio is the *Feel-Good Thompson Sampling* (FGTS) algorithm of Zhang [2022]. For the 3-information ratio, we consider a mixture of the FGTS policy and an exploratory policy. The following lemma provides bounds on the surrogate information ratios of the SOIDS algorithm.

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

$$\overline{IR}_{t}^{(2)}(\pi_{t}^{(\mathbf{SOIDS})}) \leq \overline{IR}_{t}^{(2)}(\pi_{t}^{(\mathbf{FGTS})}) \leq 2d \tag{14}$$

and

$$\overline{IR}_{t}^{(3)}(\pi_{t}^{(\mathbf{SOIDS})}) \le 2\overline{IR}_{t}^{(3)}(\pi_{t}^{(\mathbf{mix})}) \le \frac{54s}{C_{\min}}.$$
(15)

The explicit definition of both forerunner algorithms, as well as the proof of this lemma, are deferred to Appendix F.3. Finally, it remains to pick the learning rate λ_t . The following lemma describes the appropriate learning rate for the data-poor and the data-rich regimes separately.

Lemma 8. The choice of learning rate $\lambda_t^{(2)} = \sqrt{\frac{3C_{t+1}}{128d(t+1)}}$ guarantees

$$\frac{C_T}{\lambda_{T-1}^{(2)}} + \frac{32}{3} \sum_{t=1}^{T} \lambda_{t-1}^{(2)} \overline{IR}_t^{(2)} (\pi_t^{(\mathbf{SOIDS})}) \le 16\sqrt{\frac{2}{3}C_T dT}.$$

The choice of learning rate $\lambda_t^{(3)} = \frac{1}{4 \cdot 6^{\frac{1}{3}}} \left(\frac{C_{t+1} \sqrt{C_{\min}}}{(t+1)\sqrt{s}} \right)^{\frac{2}{3}}$ guarantees

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

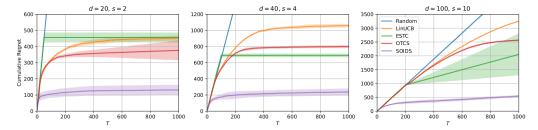


Figure 1: Cumulative regret for d=20 (left) 40 (middle) and 100 (right). We plot the mean \pm standard deviation over 10 repetitions.

The proof is deferred to Appendix G.2. The expression for the constant $c_3^* \leq 2$ can be found in Appendix I.2 (cf. Lemma 27). It remains to analyze what happens when the learning rate $\lambda_t = \min(\frac{1}{2}, \max(\lambda_t^{(2)}, \lambda_t^{(3)}))$ is chosen. We defer this to Appendix G.4.

6 Experiments

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

It is generally difficult to run the SOIDS algorithm exactly because the surrogate information ratio contains expectations w.r.t. the optimistic posterior. In our implementation of SOIDS, we use the empirical Bayesian sparse sampling procedure of Hao et al. [2021] to draw approximate samples from the optimistic posterior, and then approximate the surrogate information ratio via sample averages. We provide further details regarding the implementations of each method in Appendix J.

For each $d \in \{20, 40, 100\}$, θ_0 is the s-sparse vector in \mathbb{R}^d , with s = d/10, in which first s components are 10/s and the remaining components are zero. The action set consists of 200 random draws from the uniform distribution on $[-1,1]^d$. The noise variance is 1 and we run each method 10 times. In Figure 1, we report the cumulative regret over T = 1000 steps. As d is varied from 20 to 100, we appear to transition from the data-rich regime to the data-poor regime: for d = 20, the OTCS method is the best performing baseline, whereas for d = 100, ETCS is the best performing baseline. As our theoretical results would suggest, SOIDS performs well in both regimes.

7 Conclusion

There remain several interesting questions that our work leaves open for future research, such as the possibility of improving the logarithmic terms in the instance-dependent regret bound (as mentioned earlier in Section 4). We highlight another question below.

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

Acknowledgements. The authors wish to thank Johannes Kirschner for thought-provoking discussions about information directed sampling before the preparation of this manuscript. This project has received funding from the European Research Council (ERC), under the European Union's Horizon 2020 research and innovation programme (Grant agreement No. 950180).

References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24, 2011.
- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Online-to-confidence-set conversions and application to sparse stochastic bandits. 2012.
- Naoki Abe and Philip M Long. Associative reinforcement learning using linear probabilistic concepts. In *ICML*, pages 3–11. Citeseer, 1999.
- Jacob D. Abernethy, Elad Hazan, and Alexander Rakhlin. Competing in the dark: An efficient algorithm for bandit linear optimization. pages 263–274, 2008. URL http://colt2008.cs.helsinki.fi/papers/127-Abernethy.pdf.
- Pierre Alquier and Karim Lounici. Pac-bayesian theorems for sparse regression estimation with exponential weights. *Electronic Journal of Statistics*, 5:127–145, 2011.
- Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002.
- Hamsa Bastani and Mohsen Bayati. Online decision making with high-dimensional covariates. *Operations Research*, 68(1):276–294, 2020.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. Concentration Inequalities A Nonasymptotic Theory of Independence. 2013. ISBN 978-0-19-953525-5. doi: 10.1093/ACPROF:OSO/9780199535255.001.0001. URL https://doi.org/10.1093/acprof:oso/9780199535255.001.0001.
- Sébastien Bubeck and Mark Sellke. First-Order Bayesian Regret Analysis of Thompson Sampling, 2022. URL http://arxiv.org/abs/1902.00681.
- Sunrit Chakraborty, Saptarshi Roy, and Ambuj Tewari. Thompson sampling for high-dimensional sparse linear contextual bandits. In *International Conference on Machine Learning*, pages 3979–4008. PMLR, 2023.
- Wei Chu, Lihong Li, Lev Reyzin, and Robert E. Schapire. Contextual bandits with linear payoff functions. volume 15 of *JMLR Proceedings*, pages 208–214, 2011. URL http://proceedings.mlr.press/v15/chu11a/chu11a.pdf.
- Eugenio Clerico, Hamish Flynn, Wojciech Kotłowski, and Gergely Neu. Confidence sequences for generalized linear models via regret analysis, 2025. URL https://arxiv.org/abs/2504.16555.
- Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit feedback. In *COLT*, volume 2, page 3, 2008.
- Dylan J. Foster and Alexander Rakhlin. Beyond UCB: Optimal and Efficient Contextual Bandits with Regression Oracles, 2020. URL http://arxiv.org/abs/2002.04926.
- Dylan J. Foster, Noah Golowich, Jian Qian, Alexander Rakhlin, and Ayush Sekhari. A Note on Model-Free Reinforcement Learning with the Decision-Estimation Coefficient, 2022a. URL http://arxiv.org/abs/2211.14250.
- Dylan J. Foster, Sham M. Kakade, Jian Qian, and Alexander Rakhlin. The Statistical Complexity of Interactive Decision Making, 2022b. URL http://arxiv.org/abs/2112.13487.
- Dylan J. Foster, Alexander Rakhlin, Ayush Sekhari, and Karthik Sridharan. On the Complexity of Adversarial Decision Making, 2022c. URL http://arxiv.org/abs/2206.13063.
- Sébastien Gerchinovitz. Sparsity regret bounds for individual sequences in online linear regression. *The Journal of Machine Learning Research*, 14(1):729–769, 2013.
- Botao Hao and Tor Lattimore. Regret bounds for information-directed reinforcement learning. 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/b733cdd80ed2ae7e3156d8c33108c5d5-Abstract-Conference.html.

- Botao Hao, Tor Lattimore, and Mengdi Wang. High-dimensional sparse linear bandits. 2020.
- Botao Hao, Tor Lattimore, and Wei Deng. Information directed sampling for sparse linear bandits. 2021.
- Botao Hao, Tor Lattimore, and Chao Qin. Contextual Information-Directed Sampling, 2022. URL http://arxiv.org/abs/2205.10895.
- Elad Hazan and Satyen Kale. Extracting certainty from uncertainty: Regret bounded by variation in costs. *Machine Learning*, 80(2-3):165–188, 2010. doi: 10.1007/S10994-010-5175-X. URL https://doi.org/10.1007/s10994-010-5175-x.
- Kyoungseok Jang, Chicheng Zhang, and Kwang-Sung Jun. Popart: Efficient sparse regression and experimental design for optimal sparse linear bandits. *Advances in Neural Information Processing Systems*, 35:2102–2114, 2022.
- Gi-Soo Kim and Myunghee Cho Paik. Doubly-robust lasso bandit. *Advances in Neural Information Processing Systems*, 32, 2019.
- Johannes Kirschner and Andreas Krause. Information Directed Sampling and Bandits with Heteroscedastic Noise, 2018. URL http://arxiv.org/abs/1801.09667.
- Johannes Kirschner, Tor Lattimore, and Andreas Krause. Information directed sampling for linear partial monitoring. volume 125 of *Proceedings of Machine Learning Research*, pages 2328–2369, 2020. URL http://proceedings.mlr.press/v125/kirschner20a.html.
- Johannes Kirschner, Tor Lattimore, Claire Vernade, and Csaba Szepesvári. Asymptotically optimal information-directed sampling. volume 134 of *Proceedings of Machine Learning Research*, pages 2777–2821, 2021. URL http://proceedings.mlr.press/v134/kirschner21a.html.
- Johannes Kirschner, Seyed Alireza Bakhtiari, Kushagra Chandak, Volodymyr Tkachuk, and Csaba Szepesvári. Regret minimization via saddle point optimization. 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/6eaf8c729af4fbeb18006dc2e6a41d9b-Abstract-Conference.html.
- Tor Lattimore and András György. Mirror Descent and the Information Ratio. volume 134 of *Proceedings of Machine Learning Research*, pages 2965–2992, 2021. URL http://proceedings.mlr.press/v134/lattimore21b.html.
- Tor Lattimore and Csaba Szepesvári. Bandit algorithms. Cambridge University Press, 2020.
- Gergely Neu, Julia Olkhovskaya, Matteo Papini, and Ludovic Schwartz. Lifting the Information Ratio: An Information-Theoretic Analysis of Thompson Sampling for Contextual Bandits, 2022. URL http://arxiv.org/abs/2205.13924.
- Gergely Neu, Matteo Papini, and Ludovic Schwartz. Optimistic information directed sampling. 2024.
- Min-hwan Oh, Garud Iyengar, and Assaf Zeevi. Sparsity-agnostic lasso bandit. In *International Conference on Machine Learning*, pages 8271–8280. PMLR, 2021.
- Francesco Orabona. A modern introduction to online learning. *CoRR*, abs/1912.13213, 2019. URL http://arxiv.org/abs/1912.13213.
- Paat Rusmevichientong and John N Tsitsiklis. Linearly parameterized bandits. *Mathematics of Operations Research*, 35(2):395–411, 2010.
- Daniel Russo and Benjamin Van Roy. An information-theoretic analysis of thompson sampling. *Journal of Machine Learning Research*, 17:68:1–68:30, 2016. URL https://jmlr.org/papers/v17/14-087.html.
- Daniel Russo and Benjamin Van Roy. Learning to Optimize via Information-Directed Sampling, 2017. URL http://arxiv.org/abs/1403.5556.

- Michal Valko, Rémi Munos, Branislav Kveton, and Tomáš Kocák. Spectral bandits for smooth graph functions. volume 32 of *JMLR Workshop and Conference Proceedings*, pages 46–54, 2014. URL http://proceedings.mlr.press/v32/valko14.html.
- M.J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. 2019. ISBN 978-1-108-49802-9. URL https://books.google.es/books?id=8C8nuQEACAAJ.
- Xue Wang, Mingcheng Wei, and Tao Yao. Minimax concave penalized multi-armed bandit model with high-dimensional covariates. In *International Conference on Machine Learning*, pages 5200–5208. PMLR, 2018.
- Tong Zhang. Feel-good thompson sampling for contextual bandits and reinforcement learning. *SIAM Journal on Mathematics of Data Science*, 4(2):834–857, 2022. doi: 10.1137/21M140924X. URL https://doi.org/10.1137/21m140924x.

NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

1 Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We provide theoretical proof and experiments supporting the claims made in the abstract.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We discuss the limitations of the work and some remaining open questions. We comment on the computational efficiency of the proposed algorithm.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.

- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: The complete proof and assumptions of every single result is either in the main text or the appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The code and data are provided and the experimental setting is described.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.

- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The code is available online.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The full code is available.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental
 material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: See the figures included in the paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: All of the experiments can be run on the CPU of a laptop.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We respect NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
 deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: Due the theoretical nature of the work, societal impacts are hard to foresee.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]
Justification:
Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.

- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
 not require this, but we encourage authors to take this into account and make a best
 faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All of the code has been written by the authors of the paper.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]
Justification:
Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]
Justification:
Guidelines:

• The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]
Justification:
Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]
Justification:
Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

A Related work

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

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

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

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

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

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

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

Optimistic posterior approaches for frequentist guarantees: The approach most closely related to our work modifies the Bayesian posterior to provide frequentist guarantees. Introduced by Zhang [2022], the optimistic posterior is a modification of the Bayesian posterior which enables frequentist regret bounds for a variant of Thompson Sampling. Subsequently, Neu et al. [2024] studied the optimistic posterior framework in greater depth, defining a frequentist analog of the information

ratio to extend IDS to frequentist settings. A notable limitation of these works is their restriction to constant learning rates in the optimistic posterior, which limits adaptivity, an issue that we address in this paper.

B Analysis of the Optimistic posterior

This section provides further details about the prior underlying the optimistic posterior and guarantees on the posterior updates.

B.1 Follow the regularized leader analysis

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

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

B.1.1 Proof of Lemma 2

As is usual in the analysis of the follow the regularized leader algorithm, we introduce the Fenchel conjugate $\Phi^*: \mathbb{R}^\Theta \to \mathbb{R}$ of the regularization function $\Phi = \mathcal{D}_{KL}\left(\cdot \middle\| Q_1^+\right)$, which takes values $\Phi^*(L) = \sup_{P \in \Delta(\Theta)} \left\{ \langle P, L \rangle - \Phi(P) \right\}$. The Fenchel–Young inequality guarantees that for any $P \in \Delta(\Theta)$ and $L \in \mathbb{R}^\Theta$,

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

We introduce the maximum likelihood estimator $\theta_t = \arg\min_{\theta \in \Theta} L_t^{(1)}(\theta)$ and the function

$$L(\cdot) = -\eta (L_T^{(1)}(\cdot) - L_T^{(1)}(\theta_T)) - \lambda_T L_T^{(2)}(\cdot).$$

Since λ_T is never used by the algorithm, we can assume that $\lambda_T = \lambda_{T-1}$. The role of the maximum likelihood estimator is to ensure that the term $L_t^{(1)}(\theta) - L_t^{(1)}(\theta_t)$ is always non-negative. The Fenchel–Young inequality tells us that

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

Noticing that $\langle P, L_T^{(2)} \rangle = -\sum_{t=1}^T \Delta(P, A_t)$ and rearranging the terms concludes the proof.

B.1.2 Proof of Lemma 4

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

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

In the usual follow-the-regularized-leader analysis, we use the fact that $\lambda \mapsto \frac{\phi^*(\lambda L)}{\lambda}$ is non-decreasing for any $L \in \mathbb{R}^\Theta$. Here however, only one component of the linear loss is scaled by λ_t , and so the standard FTRL analysis fails. Crucially, because we introduced the maximum likelihood estimator θ_t , we have that $L_t^{(1)}(\cdot) - L_t^{(1)}(\theta_t) \geq 0$. We can therefore use the following lemma that guarantees that a scaled and shifted dual is monotonic.

Lemma 9. Let $\Phi \geq 0$ and Φ^* be a convex function and its dual as defined previously. For $L_1, L_2 \in \mathbb{R}^\Theta$ such that $L_1 \geq 0$, the mapping $\lambda \mapsto \frac{\Phi^*(-L_1 + \lambda L_2)}{\lambda}$ is non-decreasing on \mathbb{R}^{+*} .

Proof. By definition, we have

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

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 non-decreasing with respect to lambda. Since the supremum of non-decreasing functions is also non-decreasing, this concludes the proof.

Applying the previous lemma, we upper bound the previous sum by replacing each λ_t factor by λ_{t-1} (using the convention $\lambda_0=1/2$), and then we replace the maximum likelihood estimator θ_t by θ_0 inside Φ^* to obtain

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

It remains to bound the difference of the negative log likelihood of the true parameter and the maximum likelihood estimator. This is done via the following result (whose proof we relegate to appendix E.1.1).

Lemma 10. For any $t \ge 1$, we have

$$0 \le \mathbb{E}\left[L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_t)\right] \le \inf_{\rho} \left\{2\rho t + s\log\frac{ed(1+2/\rho)}{s}\right\} \le 6 + s\log\frac{edt}{s}$$
 (16)

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

$$\mathbb{E}\left[\sum_{t=1}^{T} \frac{\eta}{\lambda_{t-1}} (L_{t}^{(1)}(\theta_{t}) - L_{t}^{(1)}(\theta_{0}) + L_{t-1}^{(1)}(\theta_{0}) - L_{t-1}^{(1)}(\theta_{t-1})) + \frac{\eta}{\lambda_{T}} (L_{T}^{(1)}(\theta_{0}) - L_{T}^{(1)}(\theta_{T}))\right]$$

$$= \mathbb{E}\left[\sum_{t=1}^{T} \frac{\eta}{\lambda_{t-1}} (L_{t}^{(1)}(\theta_{t}) - L_{t}^{(1)}(\theta_{0})) - \sum_{t=1}^{T} \frac{\eta}{\lambda_{t}} (L_{t}^{(1)}(\theta_{t}) - L_{t}^{(1)}(\theta_{0}))\right]$$

$$\leq \eta \cdot \sum_{t=1}^{T} \mathbb{E}\left[(L_{t}^{(1)}(\theta_{0}) - L_{t}^{(1)}(\theta_{t})) \right] \left(\frac{1}{\lambda_{t}} - \frac{1}{\lambda_{t-1}}\right)$$

$$\leq \frac{\eta(6 + s \log \frac{edT}{s})}{\lambda_{T}}.$$

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

B.1.3 Controlling the losses separately

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

For this analysis, we define the true information gain as

$$IG_t(\pi) = \frac{1}{2} \sum_{a \in \mathcal{A}} \pi(a) \int_{\Theta} (\langle \theta - \theta_0, a \rangle)^2 dQ_t^+(\theta), \tag{17}$$

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

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

$$\overline{IG}_t(\pi) \le IG_t(\pi) \tag{18}$$

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

Lemma 11. Assume that the noise ϵ_t is conditionnally 1-sub-Gaussian, then for any $t \ge 1, \eta, \alpha \ge 0$ such that $\gamma = \frac{\eta \alpha}{2} (1 - \eta \alpha) > 0$, the following inequality holds

$$\mathbb{E}\left[\log \int_{\Theta} \left(\frac{p(Y_t|\theta, A_t)}{p(Y_t|\theta_0, A_t)}\right)^{\eta \alpha} dQ_t^+(\theta)\right] \le -2\gamma (1 - 2\gamma) \mathbb{E}\left[IG_t(\pi_t)\right]$$
(19)

$$\leq -2\gamma(1-2\gamma)\mathbb{E}\left[\overline{IG}_t(\pi_t)\right].$$
 (20)

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

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

$$\begin{split} \mathbb{E}\left[\log\int_{\Theta}\left(\frac{p(Y_{t}|\theta,A_{t})}{p(Y_{t}|\theta_{0},A_{t})}\right)^{\eta\alpha}\mathrm{d}Q_{t}^{+}(\theta)\right] &= \mathbb{E}\left[\mathbb{E}\left[\log\int_{\Theta}\left(\frac{p(Y_{t}|\theta,A_{t})}{p(Y_{t}|\theta_{0},A_{t})}\right)^{\eta\alpha}\mathrm{d}Q_{t}^{+}(\theta)\middle|\mathcal{F}_{t-1},A_{t}\right]\right] \\ &\leq \mathbb{E}\left[\log\mathbb{E}\left[\int_{\Theta}\left(\frac{p(Y_{t}|\theta,A_{t})}{p(Y_{t}|\theta_{0},A_{t})}\right)^{\eta\alpha}\mathrm{d}Q_{t}^{+}(\theta)\middle|\mathcal{F}_{t-1},A_{t}\right]\right] \\ &= \mathbb{E}\left[\log\int_{\Theta}\mathbb{E}\left[\exp\left(-\eta\alpha\left(\frac{(Y_{t}-\langle\theta,A_{t}\rangle)^{2}}{2}-\frac{(Y_{t}-\langle\theta_{0},A_{t}\rangle)^{2}}{2}\right)\right)\middle|\mathcal{F}_{t-1},A_{t}\right]\mathrm{d}Q_{t}^{+}(\theta)\right]. \end{split}$$

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 some elementary manipulations and the conditional sub-Gaussian tail behaviour of ϵ_t and $Y_t = r_0 + \epsilon_t$ which implies that for any (\mathcal{F}_{t-1}, A_t) -measurable ζ_t , $\mathbb{E}\left[\exp\left(Y_t\zeta_t\right)|\mathcal{F}_{t-1}, A_t\right] = \exp(r_0\zeta_t)\mathbb{E}\left[\exp\left(\epsilon_t\zeta_t\right)|\mathcal{F}_{t-1}, A_t\right] \leq \exp(r_0\zeta_t)\exp\left(\frac{\zeta_t^2}{2}\right)$, we have

$$\mathbb{E}\left[\exp\left(-\eta\alpha\left(\frac{(Y_t-r)^2}{2} - \frac{(Y_t-r_0)^2}{2}\right)\right)\middle|\mathcal{F}_{t-1}, A_t\right]$$

$$= \mathbb{E}\left[\exp\left(-\frac{\eta\alpha}{2}(2Y_t-r-r_0)(r_0-r)\right)\middle|\mathcal{F}_{t-1}, A_t\right]$$

$$= \exp\left(\eta\alpha\frac{r_0^2-r^2}{2}\right)\mathbb{E}\left[\exp\left(\eta\alpha Y_t(r-r_0)\right)\middle|\mathcal{F}_{t-1}, A_t\right]$$

$$\leq \exp\left(\eta\alpha\frac{r_0^2-r^2}{2}\right) \cdot \exp\left(\eta\alpha r_0(r-r_0)\right) \exp\left(\frac{\eta^2\alpha^2}{2}(r-r_0)^2\right)$$

$$= \exp\left(-(r-r_0)^2 \cdot \frac{\eta\alpha}{2}(1-\eta\alpha)\right).$$

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

$$\mathbb{E}\left[\exp\left(-\eta\alpha\left(\frac{(Y_t - r)^2}{2} - \frac{(Y_t - r_0)^2}{2}\right)\right) \middle| \mathcal{F}_{t-1}, A_t\right]$$

$$\leq \exp(-\gamma(r - r_0)^2)$$

$$\leq 1 - \gamma(r - r_0)^2 + \frac{\gamma^2}{2}(r - r_0)^4$$

$$\leq 1 - \gamma(r - r_0)^2 + 2\gamma^2(r - r_0)^2$$

$$\leq 1 - \gamma(1 - 2\gamma)(r - r_0)^2.$$

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

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

Rearranging and combining the result with Proposition 1 yields the claim of the lemma.

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

Lemma 12. For any $t \ge 1, \beta > 1$, if $\beta \lambda_{t-1} \le 1$, we have

$$\mathbb{E}\left[\frac{1}{\beta\lambda_{t-1}}\log\int_{\Theta}\exp(\beta\lambda_{t-1}\Delta(A_t,\theta))\,\mathrm{d}Q_t^+(\theta)\right] \le \mathbb{E}\left[2\widehat{\Delta}_t(\pi_t)\right]. \tag{21}$$

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

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

$$\mathbb{E}\left[\sum_{t=1}^{T} \frac{\Phi^{*}(-\eta(L_{t}^{(1)}(\cdot) - L_{t}^{(1)}(\theta_{0})) - \lambda_{t-1}L_{t}^{(2)}(\cdot))}{\lambda_{t-1}} - \frac{\Phi^{*}(-\eta(L_{t-1}^{(1)}(\cdot) - L_{t-1}^{(1)}(\theta_{0})) - \lambda_{t-1}L_{t-1}^{(2)}(\cdot))}{\lambda_{t-1}}\right]$$

$$= \mathbb{E}\left[\sum_{t=1}^{T} \frac{1}{\lambda_{t-1}} \log \frac{\int_{\Theta} \left(\frac{p(Y_{t}|\theta, a_{t})}{p(Y_{t}|\theta_{0}, A_{t})}\right)^{\eta} \exp\left(\lambda_{t-1}\Delta(A_{t}, \theta)\right) \exp\left(-\eta L_{t-1}^{(1)}(\theta) - \lambda_{t-1}L_{t-1}^{(2)}(\theta)\right) dQ_{1}^{+}(\theta)}{\int_{\Theta} \exp\left(-\eta L_{t-1}^{(1)}(\theta) - \lambda_{t-1}L_{t-1}^{(2)}(\theta)\right) dQ_{1}^{+}(\theta)}\right]$$

$$= \mathbb{E}\left[\sum_{t=1}^{T} \frac{1}{\lambda_{t-1}} \log \int_{\Theta} \left(\frac{p(Y_{t}|\theta, A_{t})}{p(Y_{t}|\theta_{0}, A_{t})}\right)^{\eta} \exp\left(\lambda_{t-1}\Delta(A_{t}, \theta)\right) dQ_{t}^{+}(\theta)\right]$$

$$\leq \mathbb{E}\left[\sum_{t=1}^{T} \frac{1}{\alpha\lambda_{t-1}} \log \int_{\Theta} \left(\frac{p(Y_{t}|\theta, A_{t})}{p(Y_{t}|\theta_{0}, A_{t})}\right)^{\eta\alpha} + \frac{1}{\beta\lambda_{t-1}} \log \int_{\Theta} \exp\left(\beta\lambda_{t-1}\Delta(A_{t}, \theta)\right) dQ_{t}^{+}(\theta)\right],$$

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

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

In order to construct the prior Q_1^+ and the comparator P for the regret analysis, we need to take into account two criteria: that $\mathcal{D}_{\mathrm{KL}}\left(P\big\|Q_1^+\right)$ be controlled and that $|\langle P,L\rangle-L(\theta_0)|$ be small. The comparator will be a function of the unknown parameter θ_0 , and thus we denote it by P_{θ_0} .

As for the prior, it should take into account the sparsity level of the unknown θ_0 , but should have no access to its support. We first design a distribution Π over the set of all subsets of $[d] = \{1, \ldots, d\}$, which have cardinality at most s. We choose the distribution such that: a) the probability assigned to each subset depends only on its cardinality; b) the probability assigned to the set of all subsets of size k is proportional to 2^{-k} , where $1 \le k \le s$. In other words, we prefer smaller subsets and have no preference over which indices in [d] are included. The distribution that satisfies these requirements is

$$\Pi(S) = \frac{2^{-|S|}}{\binom{d}{|S|} \sum_{k=1}^{s} 2^{-k}}.$$
(22)

For $S=\emptyset$, we set $\Pi(S)=0$. Doing so only complicates matters if the support of θ_0 is empty (i.e., $\theta_0=0$). However, in this case, the reward function is 0 everywhere, which means any algorithm would have 0 regret. We therefore continue under the assumption that $\theta_0\neq 0$. The most important property of this distribution, which we will use later, is that for any subset S of cardinality s, $\log(1/\Pi(S)) \leq s\log(2ed/s)$. For each subset S, we define S0 to be the uniform distribution on S1. The prior is defined to be

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

As for the comparator distribution P_{θ_0} , we would ideally like to take a Dirac measure on θ_0 , but this would make the KL divergence appearing in the bound blow up. Thus, we pick a comparator P_{θ_0} which dilutes its mass around θ_0 . For any $\bar{\theta} \in \Theta$, with support \bar{S} , and any $\epsilon \in (0,1)$, we define the set $(1-\epsilon)\bar{\theta}+\epsilon\Theta_{\bar{S}}=\{(1-\epsilon)\bar{\theta}+\epsilon\theta':\theta'\in\Theta_{\bar{S}}\}\subset\Theta_{\bar{S}}$. We choose P_{θ_0} to be the uniform distribution on $(1-\epsilon)\theta_0+\epsilon\Theta_{S_0}$, where S_0 is the support of θ_0 . We now bound $\Phi(P_{\theta_0})=\mathcal{D}_{\mathrm{KL}}\left(P_{\theta_0}\|Q_1^+\right)$ for this choice of P_{θ_0} in the following lemma, from which the claim of Lemma 3 then directly follows.

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

Proof. We notice that $(1-\epsilon)\bar{\theta}+\epsilon\Theta_{\bar{S}}$ is an s-dimensional ℓ_1 -ball of radius ϵ , which is contained in $\Theta_{\bar{S}}$. Therefore, on the support of $P_{\bar{\theta}}$, $\frac{\mathrm{d}P_{\bar{\theta}}}{\mathrm{d}Q_{\bar{S}}}$ is equal to the ratio of the volumes of a unit ℓ_1 -ball and an ℓ_1 -ball of radius ϵ , which is $(1/\epsilon)^s$. Thus,

$$\mathcal{D}_{\mathrm{KL}}\left(P_{\bar{\theta}} \middle\| Q_{1}^{+}\right) = \int \log \frac{\mathrm{d}P_{\bar{\theta}}}{\sum_{S} \Pi(S) \mathrm{d}Q_{S}} \mathrm{d}P_{\bar{\theta}} \leq \int \log \frac{\mathrm{d}P_{\bar{\theta}}}{\Pi(\bar{S}) \mathrm{d}Q_{\bar{S}}} \mathrm{d}P_{\bar{\theta}} \leq s \log \frac{1}{\epsilon} + \log \frac{1}{\Pi(\bar{S})} \ .$$

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

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

Combining everything, we obtain

$$\mathcal{D}_{KL}\left(P_{\bar{\theta}} \| Q_1^+\right) \le s \log \frac{1}{\epsilon} + s \log \frac{2ed}{s} = s \log \frac{2ed}{\epsilon s}, \tag{23}$$

as advertised.

C Proof of the history-dependent part of Theorem 1

Following the original analysis, we arrive again at (9).

$$\sum_{t=1}^{T} \Delta(A_t, P) \leq \frac{\Phi(P)}{\lambda_T} + \frac{\Phi^*(-\eta L_T^{(1)}(\cdot) + \eta L_T^{(1)}(\theta_T) - \lambda_T L_T^{(2)}(\cdot))}{\lambda_T} + \frac{\eta}{\lambda_T} (\langle P, L_T^{(1)} \rangle - L_T^{(1)}(\theta_T)),$$

where $P \in \Delta(\Theta)$ can be any comparator distribution. Lemma 3 is still valid and we can chose the same prior as before. We can still choose a comparator distribution supported on an ϵ -ball around θ_0 . However, because λ_t depends on the history, we can no longer upper bound $\mathbb{E}\left[\frac{|P \cdot L_T^{(1)} - L_T^{(1)}(\theta_0)|}{\lambda_{T-1}}\right]$

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-Lipschitz. Hence,

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

where we used $C_{2,T}$, a deterministic upper bound on $\frac{1}{\lambda_{T-1}}$. Exactly the same telescoping of Φ^* can be done, however because the learning rate is history-dependent, the difference between the negative log likelihood of θ_0 and θ_t must be treated with more care. We have the following lemma

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

$$\mathbb{E}\left[\sum_{t=1}^{T} \frac{\eta}{\lambda_{t-1}} \left(L_{t}^{(1)}(\theta_{t}) - L_{t}^{(1)}(\theta_{0}) + L_{t-1}^{(1)}(\theta_{0}) - L_{t-1}^{(1)}(\theta_{t-1})\right) + \frac{\eta}{\lambda_{T}} \left(L_{T}^{(1)}(\theta_{0}) - L_{T}^{(1)}(\theta_{T})\right)\right] \\
\leq \mathbb{E}\left[\frac{\eta(15 + 3s \log \frac{2e^{2}dT^{2}C_{1,T}^{2}}{s})}{2\lambda_{T-1}}\right].$$
(24)

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

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

D Proof of Theorem 3

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

$$R_T \le \mathbb{E}\left[\frac{C_T}{\lambda_{T-1}} + \min\left(\sum_{t=1}^T \frac{32}{3} \lambda_{t-1} \overline{\mathbf{IR}}_t^{(2)}(\pi_t), \frac{16}{3} c_3^* \sqrt{3\lambda_{t-1} \overline{\mathbf{IR}}_t^{(3)}(\pi_t)}\right)\right],\tag{25}$$

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

We let
$$\lambda_t^{(2)} = \sqrt{\frac{s}{2d + \sum_{s=1}^t \overline{\mathbb{R}}_s^{(2)}(\pi_s)}}$$
 and $\lambda_t^{(3)} = \left(\frac{s}{\frac{3\sqrt{6}s}{\sqrt{C_{\min}}} + \sum_{s=1}^t \sqrt{\overline{\mathbb{R}}_s^{(3)}(\pi_s)}}\right)^{\frac{2}{3}}$, and verify that $\lambda_t = \frac{s}{\sqrt{C_{\min}}}$

 $\max(\lambda_t^{(2)}, \lambda_t^{(3)})$ is decreasing and always smaller than $\frac{1}{2}$. We also verify that $C_{1,T} = C_{2,T} = \sqrt{\frac{dT}{s}}$ are valid upper bounds. As a result, we have the following upper bound

$$C_T = 2 + s \log \frac{4e^3d^2T^3C_{1,T}^2C_{2,T}}{s^2} \le 2 + s \log 4e^3T^{4.5} \left(\frac{d}{s}\right)^{3.5} \le 2 + 5s \log(\frac{edT}{s}).$$
 (26)

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

$$\sum_{t=1}^{T} \lambda_{t-1}^{(2)} \overline{\mathbb{R}}_{t}^{(2)}(\pi) = \sqrt{s} \sum_{t=1}^{T} \frac{\overline{\mathbb{R}}_{t}^{(2)}(\pi_{t})}{\sqrt{2d + \sum_{s=1}^{t-1} \overline{\mathbb{R}}_{s}^{(2)}(\pi_{s})}}$$

$$\leq \sqrt{s} \sum_{t=1}^{T} \frac{\overline{\mathbb{R}}_{t}^{(2)}(\pi_{t})}{\sqrt{\sum_{s=1}^{t} \overline{\mathbb{R}}_{s}^{(2)}(\pi_{s})}}$$

$$\leq 2\sqrt{s} \sum_{t=1}^{T} \overline{\mathbb{R}}_{t}^{(2)}(\pi_{t})$$

$$\leq 2\sqrt{s} \sum_{t=1}^{T} \overline{\mathbb{R}}_{t}^{(2)}(\pi_{t})}$$

$$\leq 2\sqrt{s} \left(2d + \sum_{t=1}^{T-1} \overline{\mathbb{R}}_{t}^{(2)}(\pi_{t})\right),$$

where we applied Lemma 19 with the function $f(x)=\frac{1}{\sqrt{x}}$ and $a_i=\overline{\rm IR}_i^{(2)}(\pi_i)$ to get the second 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 now define $R_T^{(2)}=\sqrt{s\left(2d+\sum_{t=1}^{T-1}\overline{\rm IR}_t^{(2)}(\pi_t)\right)}$, the data-dependent regret rate associated to the 2-surrogate-information ratio.

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

$$\sum_{t=1}^{T} \sqrt{\lambda_{t-1}^{(3)} \overline{\mathbb{R}}_{t}^{(3)}(\pi_{t})} = s^{\frac{1}{3}} \sum_{t=1}^{T} \frac{\sqrt{\overline{\mathbb{R}}_{t}^{(3)}(\pi_{t})}}{\left(\frac{3\sqrt{6}s}{\sqrt{C_{min}}} + \sum_{s=1}^{t-1} \sqrt{\overline{\mathbb{R}}_{s}^{(3)}(\pi_{s})}\right)^{\frac{1}{3}}}$$

$$\leq s^{\frac{1}{3}} \sum_{t=1}^{T} \frac{\sqrt{\overline{\mathbb{R}}_{t}^{(3)}(\pi_{t})}}{\left(\sum_{s=1}^{t} \sqrt{\overline{\mathbb{R}}_{s}^{(3)}(\pi_{s})}\right)^{\frac{1}{3}}}$$

$$\leq \frac{3}{2} s^{\frac{1}{3}} \left(\sum_{t=1}^{T} \sqrt{\overline{\mathbb{R}}_{t}^{(3)}(\pi_{t})}\right)^{\frac{2}{3}}$$

$$\leq \frac{3}{2} s^{\frac{1}{3}} \left(\frac{3\sqrt{6}s}{\sqrt{C_{min}}} + \sum_{t=1}^{T-1} \sqrt{\overline{\mathbb{R}}_{t}^{(3)}(\pi_{t})}\right),$$

where we applied Lemma 19 with the function $f(x) = \frac{1}{x^{\frac{1}{3}}}$ and $a_i = \sqrt{\overline{IR}_i^{(3)}}(\pi_i)$ to get the second inequality. This can be seen as a generalization of the usual $\sum_{t=1}^T \frac{1}{t^{\frac{1}{3}}} \leq \frac{3}{2}T^{\frac{2}{3}}$. We now

define
$$R_T^{(3)} = s^{\frac{1}{3}} \left(\frac{3\sqrt{6}s}{\sqrt{C_{min}}} + \sum_{t=1}^{T-1} \sqrt{\overline{\mathrm{IR}}_t^{(3)}(\pi_t)} \right)^{\frac{2}{3}}$$
, the data-dependent regret rate associated to

the 3-surrogate-information ratio. We now consider the last time that the learning rates $\lambda_t^{(3)}$ 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 $T_3 = \max\{t \leq T, \lambda_{t-1}^{(3)} \geq \lambda_{t-1}^{(2)}\}$. Coming back to the bound of Equation 25 and using the definition $\lambda_t = \max(\lambda_t^{(2)}, \lambda_t^{(3)})$, the following bound holds

$$\begin{split} &R_{T} \\ &\leq \mathbb{E}\left[\frac{C_{T}}{\lambda_{T-1}} + \sum_{t=1}^{T} \min\left(\frac{32}{3}\lambda_{t-1}\overline{\mathrm{IR}}_{t}^{(2)}(\pi_{t}), \frac{16}{3}c_{3}^{*}\sqrt{3\lambda_{t-1}}\overline{\mathrm{IR}}_{t}^{(3)}(\pi_{t})\right)\right] \\ &\leq \mathbb{E}\left[C_{T} \min\left(\frac{1}{\lambda_{T-1}^{(2)}}, \frac{1}{\lambda_{T-1}^{(3)}}\right) + \sum_{t=1}^{T} \min\left(\frac{32}{3}\max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)})\overline{\mathrm{IR}}_{t}^{(2)}(\pi_{t}), \frac{16}{3}c_{3}^{*}\sqrt{3\max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)})\overline{\mathrm{IR}}_{t}^{(3)}(\pi_{t})}\right)\right]. \end{split}$$

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

$$\sum_{t=1}^{T} \min \left(\frac{32}{3} \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)}) \overline{\mathbb{IR}}_{t}^{(2)}(\pi_{t}), \frac{16}{3} c_{3}^{*} \sqrt{3 \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)}) \overline{\mathbb{IR}}_{t}^{(3)}(\pi_{t})} \right) \\
\leq \sum_{\lambda_{t-1}^{(2)} \geq \lambda_{t-1}^{(3)}} \frac{32}{3} \lambda_{t-1}^{(2)} \overline{\mathbb{IR}}_{t}^{(2)}(\pi_{t}) + \sum_{\lambda_{t-1}^{(3)} \geq \lambda_{t-1}^{(2)}} \frac{16}{3} c_{3}^{*} \sqrt{3 \lambda_{t-1}^{(3)} \overline{\mathbb{IR}}_{t}^{(3)}(\pi_{t})} \\
\leq \sum_{t=1}^{T_{2}} \frac{32}{3} \lambda_{t-1}^{(2)} \overline{\mathbb{IR}}_{t}^{(2)}(\pi_{t}) + \sum_{t=1}^{T_{3}} \frac{16}{3} c_{3}^{*} \sqrt{3 \lambda_{t-1}^{(3)} \overline{\mathbb{IR}}_{t}^{(3)}(\pi_{t})}.$$

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

The crucial observation is that which of $\lambda_T^{(3)}$ or $\lambda_T^{(2)}$ is bigger will determine whether $R_T^{(2)}$ or $R_T^{(3)}$ is the term of leading order (up to some constants). More specifically, Let T be such that $\lambda_{T-1}^{(2)} \geq \lambda_{T-1}^{(3)}$

which means that $\sqrt{\frac{s}{2d+\sum_{t=1}^{T-1}\overline{\mathrm{IR}}_t^{(2)}(\pi_t)}} \geq \left(\frac{s}{\frac{3\sqrt{6}s}{\sqrt{C_{\min}}} + \sum_{t=1}^{T-1}\sqrt{\overline{\mathrm{IR}}_t^{(3)}(\pi_t)}}\right)^{\frac{2}{3}}$. Rearranging, this implies that $\sqrt{s\left(2d+\sum_{s=1}^{T-1}\overline{\mathrm{IR}}_t^{(2)}(\pi_t)\right)} \leq s^{\frac{2}{3}}\left(\frac{3\sqrt{6}s}{\sqrt{C_{\min}}} + \sum_{t=1}^{T-1}\sqrt{\overline{\mathrm{IR}}_t^{(3)}}(\pi_t)\right)^{\frac{2}{3}}$, which 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 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 $R_{T_2}^{(2)} \leq R_{T_2}^{(3)}$ and putting this together with the previous bound, we have

$$R_{T} \leq \mathbb{E}\left[\frac{C_{T}}{\lambda_{T-1}^{(3)}} + \frac{64}{3}R_{T_{2}}^{(2)} + \frac{16}{3}R_{T_{3}}^{(3)}\right]$$

$$\leq \mathbb{E}\left[\frac{C_{T}}{s}R_{T}^{(3)} + \frac{64}{3}R_{T_{2}}^{(2)} + \frac{16}{3}R_{T_{3}}^{(3)}\right]$$

$$\leq \mathbb{E}\left[\frac{C_{T}}{s}R_{T}^{(3)} + \frac{64}{3}R_{T_{2}}^{(3)} + \frac{16}{3}R_{T_{3}}^{(3)}\right]$$

$$\leq \mathbb{E}\left[\frac{C_{T}}{s}R_{T}^{(3)} + \frac{64}{3}R_{T}^{(3)} + \frac{16}{3}R_{T}^{(3)}\right]$$

$$\leq \mathbb{E}\left[\left(\frac{C_{T}}{s} + \frac{80}{3}\right)R_{T}^{(3)}\right],$$

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

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

where we use the fact that $T \to R_T^{(2)}$ and $T \to R_T^{(3)}$ are non-decreasing and $T_2 \le T, T_3 \le T$. Putting both of those bounds together with Equation 26 yields the claim of the Theorem.

E Maximum likelihood estimation

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

E.1 Bound in expectation

We start with the case in which the maximum likelihood estimator is computed on a finite subset of the parameter space Θ .

Lemma 15. Let $t \ge 1$, and Θ' be a finite subset of Θ , we define the MLE over Θ' as

$$\theta_{\mathit{MLE},t}(\Theta') = \operatorname*{arg\,min}_{\theta \in \Theta'} L_t^{(1)}(\theta).$$

Then,

$$\mathbb{E}\left[L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_{MLE,t}(\Theta'))\right] \le \log|\Theta'| \tag{27}$$

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

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

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,

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

which proves the claim.

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 ρ -covering of Θ with respect to the ℓ_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 ρ covering. We have the following bound on this quantity.

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

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

Proof. For each subset $S \subset [d]$ of cardinality |S| = s, there is a surjective isometric embedding from $(\Theta_S, \|\cdot\|_1)$ to $(\mathbb{B}^s_1(1), \|\cdot\|_1)$. In particular, to embed $\theta \in \Theta_S$ into $\mathbb{B}^s_1(1)$, one can simply remove all the components of θ corresponding to indices not in S. Therefore, for every $\rho > 0$, $\mathcal{N}(\Theta_S, \|\cdot\|_1, \rho) \leq \mathcal{N}(\mathbb{B}^s_1(1), \|\cdot\|_1, \rho)$. Moreover, via a standard argument, we have $\mathcal{N}(\mathbb{B}^s_1(1), \|\cdot\|_1, \rho) \leq (1 + \frac{2}{\rho})^s$ (see, e.g., Lemma 5.7 in Wainwright, 2019). Now, let $\Theta_{S,\rho}$ denote any minimal ρ -covering of Θ_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, $\bigcup_{S \subset [d]: |S| = s} \Theta_{S,\rho}$ forms a valid ρ -covering of Θ and its cardinality is bounded by

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

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

E.1.1 Proof of Lemma 10

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 Θ' be a minimal valid ρ -cover of Θ as is defined in Lemma 16, and $\theta' \in \Theta'$ be such that $\|\theta' - \theta_t\| \le \rho$, which exists by

definition of a ρ -covering. Then,

$$\mathbb{E}\left[L_{t}^{(1)}(\theta_{0}) - L_{t}^{(1)}(\theta_{t})\right] = \mathbb{E}\left[L_{t}^{(1)}(\theta_{0}) - L_{t}^{(1)}(\theta_{\text{MLE},t}(\Theta'))\right] \\ + \mathbb{E}\left[L_{t}^{(1)}(\theta_{\text{MLE},t}(\Theta')) - L_{t}^{(1)}(\theta')\right] \\ + \mathbb{E}\left[L_{t}^{(1)}(\theta') - L_{t}^{(1)}(\theta_{t})\right] \\ \leq \log(\mathcal{N}(\Theta, \|\cdot\|_{1}, \rho)) + 0 + 2\rho t,$$

where the first term is bounded by Lemma 15, the second term is non-positive by definition of the maximum likelihood estimator because $\theta' \in \Theta'$ and the third term is bounded because the 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 Lemma 16 and setting $\rho = \frac{2}{t}$ yields the desired bound.

E.2 High-probability bounds

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

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

$$\mathbb{P}\left[\exists t \ge 1, L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_{MLE,t}(\Theta')) \ge \log \frac{|\Theta'|}{\delta}\right] \le \delta.$$
 (28)

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 non-negative supermartingale with respect to the filtration $(\mathcal{F}_t')_t$, where $\mathcal{F}_t' = \sigma(\mathcal{F}_t,A_{t+1})$, allowing us to invoke Ville's inequality to get the following guarantee:

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

Taking the logarithm and a union bound on Θ' yields the desired result.

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

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

$$\mathbb{E}\left[X(L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_t))\right] \leq \inf_{\delta, \rho > 0} \left\{ \mathbb{E}\left[Xs\log\frac{ed(1+\frac{2}{\rho})}{s\delta^{\frac{1}{s}}}\right] + \frac{5}{2}B\rho t + B\delta s\log\frac{e^{1+\frac{1}{s}}d(1+\frac{2}{\rho})}{s\delta^{\frac{1}{s}}} \right\}$$

$$\leq 5 + s\log\frac{2e^2dT^2B^2}{s} \mathbb{E}\left[X + \frac{1}{T}\right]. \tag{29}$$

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

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

The first term is upper bounded by $\mathbb{E}\left[X\log\frac{N}{\delta}\right]$, the third term is upper bounded by $\frac{5}{2}B\rho t$ because $\|\theta-\theta'\|_1$ is uniformly bounded by ρ and by by Lemma 21. The fourth term is non-positive because θ' minimizes the negative log likelihood on Θ' . Finally, we turn our attention to the second term.

To simplify the 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 whether ϵ is larger or smaller than $\log\frac{N}{\delta}$ yields the following identity

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

We can now upper bound the expectation as follows

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

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. 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, we get

$$\mathbb{E}\left[X(L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_t))\right] \le \mathbb{E}\left[Xs\log\frac{ed(1+\frac{2}{\rho})}{s\delta^{\frac{1}{s}}}\right] + \frac{5}{2}B\rho t + B\delta s\log\frac{e^{1+\frac{1}{s}}d(1+\frac{2}{\rho})}{s\delta^{\frac{1}{s}}}.$$

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

E.2.1 Proof of Lemma 14

As was noted in the analysis, since λ_T is not used by the algorithm, we can replace λ_T by λ_{T-1} in our computations. We have

$$\mathbb{E}\left[\sum_{t=1}^{T} \frac{\eta}{\lambda_{t-1}} (L_{t}^{(1)}(\theta_{t}) - L_{t}^{(1)}(\theta_{0}) + L_{t-1}^{(1)}(\theta_{0}) - L_{t-1}^{(1)}(\theta_{t-1})) + \frac{\eta}{\lambda_{T}} (L_{T}^{(1)}(\theta_{0}) - L_{T}^{(1)}(\theta_{T}))\right]$$

$$= \mathbb{E}\left[\sum_{t=1}^{T} \frac{\eta}{\lambda_{t-1}} (L_{t}^{(1)}(\theta_{t}) - L_{t}^{(1)}(\theta_{0})) - \sum_{t=1}^{T} \frac{\eta}{\lambda_{t}} (L_{t}^{(1)}(\theta_{t}) - L_{t}^{(1)}(\theta_{0}))\right]$$

$$= \eta \cdot \sum_{t=1}^{T} \mathbb{E}\left[(L_{t}^{(1)}(\theta_{0}) - L_{t}^{(1)}(\theta_{t})) \left(\frac{1}{\lambda_{t}} - \frac{1}{\lambda_{t-1}} \right) \right].$$

Let $C_{1,T}$ be a deterministic upper bound on $\left(\frac{1}{\lambda_{t+1}} - \frac{1}{\lambda_t}\right)$. Applying Lemma 18 to $X = \left(\frac{1}{\lambda_{t+1}} - \frac{1}{\lambda_t}\right)$ and telescoping, we get

$$\begin{split} &\eta \cdot \sum_{t=1}^T \mathbb{E}\left[(L_t^{(1)}(\theta_0) - L_t^{(1)}(\theta_t)) \left(\frac{1}{\lambda_t} - \frac{1}{\lambda_{t-1}} \right) \right] \\ &\cdot \leq \eta \left(5 + s \log \frac{2e^2 dt^2 C_{1,T}^2}{s} \right) \sum_{t=1}^T \mathbb{E}\left[\left(\frac{1}{\lambda_t} - \frac{1}{\lambda_{t-1}} \right) + \frac{1}{T} \right] \\ &\leq \eta \left(5 + s \log \frac{2e^2 dt^2 C_{1,T}^2}{s} \right) \mathbb{E}\left[\left(\frac{1}{\lambda_T} + 1 \right) \right] \\ &\leq \mathbb{E}\left[\frac{\eta (15 + 3s \log \frac{2e^2 dt^2 C_{1,T}^2}{s})}{2\lambda_{T-1}} \right], \end{split}$$

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

F Bounding the surrogate information ratio

F.1 Proof of Lemma 6

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

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

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

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

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

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

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

F.2 Proof of Lemma 1

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

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

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

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

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

In particular,

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

This inequality is equivalent to

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

From this inequality, we obtain

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

thus proving the claim.

F.3 Proof of Lemma 7

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

As one of our forerunners, we will make use of the Feel-Good Thompson Sampling (FGTS) algorithm introduced by Zhang [2022]. Letting $\tilde{\theta}_t \sim Q_t^+$, the FGTS policy is defined as

$$\pi_t^{(\mathbf{FGTS})}(a) = \mathbb{P}_t \left[a^*(\widetilde{\theta_t}) = a \right].$$
 (30)

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

F.3.1 Bounding the two information ratio

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

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

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

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

$$\overline{\mathbf{IG}}_{t}(\pi_{t}^{(\mathbf{FGTS})}) = \frac{1}{2} \sum_{a \in \mathcal{A}} \alpha_{a} \int_{\Theta} (r(a, \bar{\theta}(Q_{t}^{+})) - r(a, \theta))^{2} dQ_{t}^{+}(\theta)$$

$$= \frac{1}{2} \sum_{a \in \mathcal{A}} \alpha_{a} \int_{\Theta} \sum_{a' \in \mathcal{A}} \mathbf{1}_{\{a^{*}(\theta) = a'\}} (r(a, \bar{\theta}(Q_{t}^{+})) - r(a, \theta))^{2} dQ_{t}^{+}(\theta)$$

$$= \frac{1}{2} \sum_{a \in \mathcal{A}} \sum_{a' \in \mathcal{A}} \alpha_{a} \int_{\Theta} \mathbf{1}_{\{a^{*}(\theta) = a'\}} dQ_{t}^{+}(\theta) \mathbb{E}_{t} [(r(a, \bar{\theta}(Q_{t}^{+})) - r(a, \tilde{\theta}_{t}) | a^{*}(\tilde{\theta}_{t}) = a']$$

$$\geq \frac{1}{2} \sum_{a \in \mathcal{A}} \sum_{a' \in \mathcal{A}} \alpha_{a} \alpha_{a'} \left(r(a, \bar{\theta}(Q_{t}^{+})) - \mathbb{E}_{t} [r(a, \tilde{\theta}_{t}) | a^{*}(\tilde{\theta}_{t}) = a'] \right)^{2}$$

$$= \frac{1}{2} \sum_{a \in \mathcal{A}} \sum_{a' \in \mathcal{A}} M_{a, a'}^{2} = \frac{1}{2} ||M||_{F}^{2}.$$

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

$$\widehat{\Delta}_{t}(\pi_{t}^{(\mathbf{FGTS})}) = \int_{\Theta} r(a^{*}(\theta), \theta) dQ_{t}^{+}(\theta) - \sum_{a \in \mathcal{A}} \alpha_{a} \int_{\Theta} r(a, \theta) dQ_{t}^{+}$$

$$= \int_{\Theta} \sum_{a \in \mathcal{A}} \mathbf{1}_{\{a^{*}(\theta) = a\}} r(a^{*}(\theta), \theta) dQ_{t}^{+}(\theta) - \sum_{a \in \mathcal{A}} \alpha_{a} r(a, \bar{\theta}(Q_{t}^{+}))$$

$$= \sum_{a \in \mathcal{A}} \alpha_{a} \mathbb{E}_{t} [r(a, \widetilde{\theta}_{t}) | a^{*}(\widetilde{\theta}_{t}) = a] - \sum_{a \in \mathcal{A}} \alpha_{a} r(a, \bar{\theta}(Q_{t}^{+}))$$

$$= \operatorname{tr}(M).$$
(31)

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

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

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

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

Therefore, M can be factorized as

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

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

F.3.2 Bounding the three information ratio

To bound the 3 information ratio we follow Hao et al. [2021] and we introduce the exploratory policy

$$\mu = \underset{\pi \in \Delta(A)}{\arg \max} \, \sigma_{\min} \left(\sum_{a \in A} \pi(a) a a^{\top} \right). \tag{32}$$

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

$$\overline{\mathbf{IG}}_{t}(\pi_{t}^{(\mathbf{mix})}) \geq \frac{1}{2} \sum_{a \in \mathcal{A}} \pi_{t}^{(\mathbf{mix})}(a) \sum_{a' \in \mathcal{A}} \mathbb{P}_{t}(a^{*}(\widetilde{\theta}_{t}) = a')(r(a, \overline{\theta}(Q_{t}^{+})) - \mathbb{E}_{t}[r(a, \widetilde{\theta}_{t}) | a^{*}(\widetilde{\theta}_{t}) = a'])^{2}$$

$$= \frac{1}{2} \sum_{a \in \mathcal{A}} \pi_{t}^{(\mathbf{mix})}(a) \sum_{a' \in \mathcal{A}} \mathbb{P}_{t}(a^{*}(\widetilde{\theta}_{t}) = a') \langle \mu_{a'} - \overline{\theta}(Q_{t}^{+}), a \rangle^{2},$$

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

$$\overline{\operatorname{IG}}_{t}(\pi_{t}^{(\operatorname{\mathbf{mix}})}) \geq \frac{\gamma}{2} \sum_{a' \in \mathcal{A}} \mathbb{P}_{t}(a^{*}(\widetilde{\theta}_{t}) = a') \sum_{a \in \mathcal{A}} \mu(a) (\mu_{a'} - \bar{\theta}(Q_{t}^{+}))^{\top} a a^{\top} (\mu_{a'} - \bar{\theta}(Q_{t}^{+})) \\
\geq \frac{\gamma}{2} \sum_{a' \in \mathcal{A}} \mathbb{P}_{t}(a^{*}(\widetilde{\theta}_{t}) = a') C_{\min} \|\mu_{a'} - \bar{\theta}(Q_{t}^{+})\|_{2}^{2}.$$

Using the expression for the surrogate regret of FGTS in (31), we obtain

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

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

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

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

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

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

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

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

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

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

$$\overline{\operatorname{IR}}_{t}^{(3)}(\pi_{t}^{(\mathbf{mix})}) = \frac{((1-\gamma)\widehat{\Delta}_{t}(\pi_{t}^{(\mathbf{FGTS})}) + \gamma\widehat{\Delta}_{t}(\mu))^{3}}{(1-\gamma)\overline{\operatorname{IG}}_{t}(\pi_{t}^{(\mathbf{FGTS})}) + \gamma\overline{\operatorname{IG}}_{t}(\mu)} \leq \frac{((1-\gamma)\widehat{\Delta}_{t}(\pi_{t}^{(\mathbf{FGTS})}) + \gamma\widehat{\Delta}_{t}(\mu))^{3}}{\gamma\overline{\operatorname{IG}}_{t}(\mu)}.$$

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

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

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

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

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

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

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

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

In this case, we can take the forerunner to be μ . The surrogate 3-information ratio can then be upper bounded as

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

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

$$\overline{\mathsf{IR}}_t^{(3)}(\pi_t^{(\mathbf{mix})}) \le \frac{27s}{C_{\min}}.$$

G Choosing the learning rates

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

G.1 Technical tools

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

Lemma 19. Let $a_i \ge 0$ and $f: [0, \infty) \to [0, \infty)$ be a nonincreasing function. Then

$$\sum_{t=1}^{T} a_t f\left(\sum_{i=0}^{t} a_i\right) \le \int_{a_0}^{\sum_{t=0}^{T} a_t} f(x) \, dx. \tag{33}$$

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

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

Proof. Let t > 0, we have

$$\frac{\log(C_1(t+1))}{\log(C_1t)} = \frac{\log\left(C_1t\left(\frac{t+1}{t}\right)\right)}{\log(C_1t)} = \frac{\log(C_1t) + \log\left(\frac{t+1}{t}\right)}{\log(C_1t)} \le 1 + \frac{1}{t\log(C_1t)} \le 1 + \frac{1}{t},$$

where the first inequality uses $\log(1+x) \leq x$ for any x > -1 and the second inequality uses $\log(C_1t) \geq \log(C_1) \geq 1$ because we assumed $C_1 \geq e$. Since $\frac{C_2(t+1)}{C_2t} = 1 + \frac{1}{t}$, we can conclude that the sequence λ_t is non-increasing.

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

We start by focusing on the data rich regime, and we bound the following part of the regret bound given in Equation (13):

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

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 will optimize later. Because $t\mapsto C_t$ is increasing, we get that $\lambda_{t-1}\leq \alpha\sqrt{\frac{C_T}{dt}}$. By Lemma 7, we know that for all $t\geq 1$, $\overline{\operatorname{IR}}_t^{(2)}(\pi_t)\leq 2d$, hence

$$\frac{C_T}{\lambda_{T-1}} + \frac{32}{3} \sum_{t=1}^T \lambda_{t-1} \overline{\mathbb{IR}}_t^{(2)}(\pi_t) \leq \frac{1}{\alpha} \sqrt{C_T dT} + \frac{64}{3} \alpha \sqrt{C_T} \sum_{t=1}^T \frac{d}{\sqrt{dt}}$$

$$\leq \frac{1}{\alpha} \sqrt{C_T dT} + \frac{128}{3} \alpha \sqrt{C_T dT}$$

$$\leq \left(\frac{1}{\alpha} + \frac{128}{3} \alpha\right) \sqrt{C_T dT}$$

$$\leq 16 \sqrt{\frac{2}{3} C_T dT},$$

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 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 $16\sqrt{\frac{2}{3}}$. This concludes the proof of the claim.

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

We now focus on the data-poor regime and specifically on bounding the following part of the bound given in Equation (13):

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

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 we will optimize later. Because $t \to C_t$ is increasing, we get that $\lambda_{t-1} \le \alpha \left(\frac{C_T \sqrt{C_{\min}}}{ts}\right)^{\frac{2}{3}}$. By Lemma 7, the 3-surrogate-information ratio is bounded for all $t \ge 1$ as $\overline{\mathbb{R}}_t^{(3)}(\pi_t) \le \frac{54s}{C_{min}}$. Hence, the following holds:

$$\frac{C_T}{\lambda_{T-1}} + \frac{16}{3} c_3^* \sum_{t=1}^T \sqrt{3\lambda_{t-1} \overline{\mathbb{R}}_t^{(3)}(\pi_t)} \leq \frac{1}{\alpha} (C_T)^{\frac{1}{3}} \left(\frac{T\sqrt{s}}{\sqrt{C_{\min}}} \right)^{\frac{2}{3}} + 48 c_3^* \sqrt{2\alpha} (C_T)^{\frac{1}{3}} \left(\frac{\sqrt{s}}{\sqrt{C_{\min}}} \right)^{\frac{2}{3}} \sum_{t=1}^T \frac{1}{t^{\frac{1}{3}}} dt^{\frac{1}{3}} dt^{$$

Here, we have applied Lemma 19 with the function $f(x)=x^{\frac{1}{3}}$ and $a_i=1$ to bound $\sum_{t=1}^T t^{-1/3} \leq \frac{3}{2}T^{\frac{2}{3}}$ in the second line, the last line comes from the choice $\alpha=\frac{1}{4\cdot 6^{\frac{1}{3}}}$ which optimizes the constant $\left(\frac{1}{\alpha} + 72c_3^*\sqrt{2\alpha}\right)$ (as per Lemma 27). This proves the statement.

Joint learning rates, end of the proof of Theorem 2

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

First, note that our choice of learning rates ensures that $\lambda_t \leq \frac{1}{2}$ holds as long as T is larger than an absolute constant, and thus we focus on this case here (and relegate the complete details of establishing this absolute constant to Appendix G.5). To proceed, we define the (constant-free) 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

bounds obtained when using the respective learning rates $\lambda_t^{(2)}$ and $\lambda_t^{(3)}$, as per Lemma 8.

We now consider the last time that the learning rates $\lambda_t^{(3)}$ 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 $T_3 = \max\{t \leq T, \lambda_{t-1}^{(3)} \geq \lambda_{t-1}^{(2)}\}$. Combining the bound of Equation 13 and using the definition $\lambda_t = \min(\frac{1}{2}, \max(\lambda_t^{(2)}, \lambda_t^{(3)}))$, the following bound holds

$$\begin{split} &R_{T} \\ &\leq \mathbb{E}\left[\frac{C_{T}}{\lambda_{T-1}} + \sum_{t=1}^{T} \min\left(\frac{32}{3}\lambda_{t-1}\overline{\mathbb{R}}_{t}^{(2)}(\pi_{t}), \frac{16}{3}c_{3}^{*}\sqrt{3\lambda_{t-1}\overline{\mathbb{R}}_{t}^{(3)}(\pi_{t})}\right)\right] \\ &= \mathbb{E}\left[\frac{C_{T}}{\min(\frac{1}{2}, \max(\lambda_{T-1}^{(2)}, \lambda_{T-1}^{(3)}))} \\ &+ \sum_{t=1}^{T} \min\left(\frac{32}{3}\min(\frac{1}{2}, \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)}))\overline{\mathbb{R}}_{t}^{(2)}(\pi_{t}), \frac{16}{3}c_{3}^{*}\sqrt{3\min(\frac{1}{2}, \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)}))}\overline{\mathbb{R}}_{t}^{(3)}(\pi_{t})}\right)\right] \\ &\leq \mathbb{E}\left[C_{T}\min\left(\frac{1}{\lambda_{T-1}^{(2)}}, \frac{1}{\lambda_{T-1}^{(3)}}\right) + \sum_{t=1}^{T}\min\left(\frac{32}{3}\max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)})\overline{\mathbb{R}}_{t}^{(2)}(\pi_{t}), \frac{16}{3}c_{3}^{*}\sqrt{3\max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)})}\overline{\mathbb{R}}_{t}^{(3)}(\pi_{t})}\right)\right]. \end{split}$$

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

$$\sum_{t=1}^{T} \min \left(\frac{32}{3} \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)}) \overline{\mathbb{R}}_{t}^{(2)}(\pi_{t}), \frac{16}{3} c_{3}^{*} \sqrt{3 \max(\lambda_{t-1}^{(2)}, \lambda_{t-1}^{(3)}) \overline{\mathbb{R}}_{t}^{(3)}(\pi_{t})} \right) \\
\leq \sum_{\lambda_{t-1}^{(2)} \geq \lambda_{t-1}^{(3)}} \frac{32}{3} \lambda_{t-1}^{(2)} \overline{\mathbb{R}}_{t}^{(2)}(\pi_{t}) + \sum_{\lambda_{t}^{(3)} \geq \lambda_{t}^{(2)}} \frac{16}{3} c_{3}^{*} \sqrt{3 \lambda_{t-1}^{(3)} \overline{\mathbb{R}}_{t}^{(3)}(\pi_{t})} \\
\leq \sum_{t-1}^{T_{2}} \frac{32}{3} \lambda_{t-1}^{(2)} \overline{\mathbb{R}}_{t}^{(2)}(\pi_{t}) + \sum_{t-1}^{T_{3}} \frac{16}{3} c_{3}^{*} \sqrt{3 \lambda_{t-1}^{(3)} \overline{\mathbb{R}}_{t}^{(3)}(\pi_{t})}.$$

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

The crucial observation is that which of $\lambda_T^{(3)}$ or $\lambda_T^{(2)}$ is bigger will determine whether $R_T^{(2)}$ or $R_T^{(3)}$ is the term of leading order (up to some constants). More specifically, Let T be such that $\lambda_{T-1}^{(2)} \geq \lambda_{T-1}^{(3)}$ 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 $\sqrt{C_T dT} \leq \frac{6^{\frac{5}{6}}}{4} \left(T\sqrt{s\frac{C_T}{C_{\min}}}\right)^{\frac{2}{3}}$, which means that $R_T^{(2)} \leq \frac{6^{\frac{5}{6}}}{4} R_T^{(3)}$. Following the exact same steps, we also 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 $\lambda_{T_2-1}^{(2)} \geq \lambda_{T_2-1}^{(3)}$ 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 previous bound, we have

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

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

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 $R_{T_3}^{(3)} \leq \frac{4}{6\frac{5}{6}}R_{T_3}^{(2)}$.

Putting this together, with the previous bound, we have

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

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

G.5 Upper bound on the learning rates

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

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

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

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

Using the loose inequality $\log \frac{edT}{s} \leq \frac{dT}{s}$, we get that this condition is satisfied for any $T \geq 1$. Similarly, we have that $\lambda_{T-1}^{(3)} \leq \frac{1}{2}$ if

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

We note that

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

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

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

Rearranging the inequality, this is equivalent to

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

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

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

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

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

This concludes the proof of Theorem 2.

I Technical Results

We state and prove the remaining technical results.

Lemma 21. Let $\pi \in \Delta(A)$, the function $\theta \mapsto \Delta(\pi, \theta)$ is 2-Lipschitz with respect to the 1 norm. Let $t \geq 1$, the function $\theta \to \mathbb{E}\left[\log\left(\frac{1}{p_t(Y_t|\theta,A_t)}\right)\right]$ is 2-Lipschitz with respect to the 1 norm. Moreover if θ, θ' are random variables, then

$$\left| \mathbb{E} \left[\log \left(\frac{1}{p_t(Y_t | \theta, A_t)} \right) - \log \left(\frac{1}{p_t(Y_t | \theta', A_t)} \right) \right] \right| \le \frac{5}{2} \sup \|\theta - \theta'\|_1.$$

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

$$|r(\pi, \theta) - r(\pi, \theta')| = \left| \sum_{a \in \mathcal{A}} \pi(a) \langle \theta - \theta', a \rangle \right|$$

$$\leq \sum_{a \in \mathcal{A}} \pi(a) |\langle \theta - \theta', a \rangle|$$

$$\leq \sum_{a \in \mathcal{A}} \pi(a) \|\theta - \theta'\|_1 \|a\|_{\infty}$$

$$\leq \|\theta - \theta'\|_1.$$

Similarly,

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

Finally

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

Let us write $r = \langle \theta, A_t \rangle$, $r' = \langle \theta', A_t \rangle$ and $r_0 = \langle \theta_0, A_t \rangle$. We have that $Y_t = r_0 + \epsilon_t$ where ϵ_t is 1-sub-Gaussian. For the negative log-likelihood, we then have

$$\mathbb{E}\left[\log\left(\frac{1}{p(Y_{t}|\theta, A_{t})}\right) - \log\left(\frac{1}{p(Y_{t}|\theta', A_{t})}\right)\right] = \frac{1}{2}\mathbb{E}\left[(\langle \theta, A_{t} \rangle - Y_{t})^{2} - (\langle \theta', A_{t} \rangle - Y_{t})^{2}\right]$$

$$= \frac{1}{2}\mathbb{E}\left[(r - Y_{t})^{2} - (r' - Y_{t})^{2}\right]$$

$$= \frac{1}{2}\mathbb{E}\left[(r - r')(r + r' - 2Y_{t})\right]$$

$$= \frac{1}{2}\mathbb{E}\left[(r - r')(r + r' - 2r_{0})\right] + \frac{1}{2}\mathbb{E}\left[(r - r')\epsilon_{t}\right]$$

$$\leq 2\mathbb{E}\left[\|\theta - \theta'\|_{1}\right] + \frac{1}{2}\mathbb{E}\left[\|\theta - \theta'\|_{1}\left|\epsilon_{t}\right|\right],$$

where we use the fact that $|r + r' - 2r_0| \le 4$. Now if θ, θ' are fixed, since $\mathbb{E}[\epsilon_t] = 0$, we get

$$\mathbb{E}\left[\log\left(\frac{1}{p(Y_t|\theta, A_t)}\right) - \log\left(\frac{1}{p(Y_t|\theta', A_t)}\right)\right] \le 2\|\theta - \theta'\|_1$$

If θ, θ' are random variables, the subgaussianity of ϵ_t implies that $\mathbb{E}\left[\epsilon_t^2\right] \leq 1$ and by Cauchy-Schwarz we have

$$\mathbb{E}\left[\left\|\theta - \theta'\right\|_{1} |\epsilon_{t}|\right] \leq \sqrt{\mathbb{E}\left[\left\|\theta - \theta'\right\|_{1}^{2}\right] \mathbb{E}\left[\epsilon_{t}^{2}\right]} \leq \sup \left\|\theta - \theta'\right\|_{1}.$$

This yields the second claim of the Lemma.

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

$$\frac{1}{\eta} \log \mathbb{E} \left[\exp \left(\eta X \right) \right] \le \mathbb{E} \left[X \right] + \frac{\eta (b - a)^2}{8}. \tag{34}$$

Proof. See for instance Chapter 2 in Boucheron et al. [2013].

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

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

$$\frac{1}{\eta}\log\mathbb{E}\left[\exp\left(\eta X\right)\right] \le \mathbb{E}\left[X\right] + \eta\mathbb{E}\left[X^2\right] \le 2\mathbb{E}\left[X\right]. \tag{35}$$

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

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

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

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

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

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

Now taking the exponential and reordering this is also equivalent to

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

Let

$$g: (-\infty, -1] \longrightarrow [-\frac{1}{e}, 0)$$

 $x \longmapsto xe^{x}.$

The previous equation can be rewritten $g\left(\frac{-\alpha}{a}\right)=-\frac{\exp\left(-\frac{a+b}{a}\right)}{a}$. We define $W_{-1}:\left[-\frac{1}{e},0\right)\longrightarrow\left(-\infty,1\right]$ as the(functional) inverse of g. The function g is the -1 branch of the Lambert W function. We have that for any $x\leq -1$, $W_{-1}(xe^x)=x$ and that for any $y\geq e$, $-W_{-1}(-\frac{1}{y})\leq 2\log(y)$. Since g is decreasing on its domain, W_{-1} is well-defined and decreasing. Moreover, for any $x\leq -1$

, $W_{-1}(g(x))=x$. In particular, we have that $\alpha=-aW_{-1}\left(-\frac{\exp\left(-\frac{a+b}{a}\right)}{a}\right)$. We will use that formulation to find an upper bound on α .

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

Because of the elementary inequality $2\log(x) \le x$ for x>0, we conclude that $y\le y^*$. Since $y\mapsto -W_{-1}(-\frac{1}{y})$ is an increasing function we finally have that for any $y\ge e$

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

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

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

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

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 supermartingale with respect to the filtration $(\mathcal{F}_t)_t$.

Proof. We have

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

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

$$= \exp\left(-\frac{(\langle\theta - \theta_0, A_t\rangle)^2}{2}\right)\mathbb{E}\left[\exp\left(\epsilon_t\langle\theta - \theta_0, A_t\rangle\right)|\mathcal{F}_{t-1}, A_t\right]$$

$$\leq \exp\left(-\frac{(\langle\theta - \theta_0, A_t\rangle)^2}{2}\right) \cdot \exp\left(\frac{(\langle\theta - \theta_0, A_t\rangle)^2}{2}\right)$$

$$= 1$$

where the inequality comes from the conditional subgaussianity of ϵ_t . Finally, by the tower rule of conditional expectations

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

I.1 Proof of Proposition 1

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

$$\overline{\mathrm{IG}}_t(\pi) = \frac{1}{2} \sum_{a \in \mathcal{A}} \pi(a) \int_{\Theta} (\langle \theta - \overline{\theta}(Q_t^+), a \rangle)^2 \, \mathrm{d}Q_t^+(\theta), \tag{36}$$

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

$$IG_t(\pi) = \frac{1}{2} \sum_{a \in A} \pi(a) \int_{\Theta} (\langle \theta, a \rangle - \langle \theta_0, a \rangle)^2 dQ_t^+(\theta), \tag{37}$$

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

$$(\langle \theta - \theta_0, a \rangle)^2 = (\langle \theta - \bar{\theta}(Q_t^+) + \bar{\theta}(Q_t^+) - \theta_0, a \rangle)^2$$

$$= (\langle \theta - \bar{\theta}(Q_t^+), a \rangle)^2 + 2\langle \theta - \bar{\theta}(Q_t^+), a \rangle \langle \bar{\theta}(Q_t^+) - \theta_0, a \rangle + (\langle \bar{\theta}(Q_t^+) - \theta_0, a \rangle)^2$$

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

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

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

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

I.2 Generalization of the AM-GM inequality

Dealing with the generalized information ratio requires bounding the cubic root of products. While one could use Hölder's inequality to deal directly with products, we find it more flexible to use a variational form of this inequality. In all that follows, we let p>1 be a real number and q be such 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 Fenchel-Young Inequality which can be seen as an extension of the AM-GM inequality.

Lemma 26. Let $x, y \ge 0$, then

$$xy \le \frac{x^p}{p} + \frac{y^q}{q}. (38)$$

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

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

$$f: \mathbb{R}^+ \longrightarrow \mathbb{R}^+$$
$$x \longmapsto \frac{x^p}{p}$$

is exactly $f^*(y) = \frac{1}{q}y^q$ (for non-negative y). Then the Lemma is a direct consequence of the Fenchel Young inequality and of its equality case.

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

Lemma 27. Let $x, y \ge 0, \lambda > 0$, then

$$\sqrt[p]{xy} \le \frac{x}{\lambda} + c_p^*(\lambda y)^{\frac{1}{p-1}} \tag{39}$$

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

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

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

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

$$\inf_{\lambda>0} \frac{\alpha}{\lambda} + \beta \lambda^{\frac{1}{p-1}} = c_p \alpha^{\frac{1}{p}} \beta^{\frac{p-1}{p}},\tag{40}$$

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

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

Remark An alternative is to pick λ to make both terms equal resulting in the same result but with 2 as a leading constant. Now

$$c_p = p^{\frac{1}{p}} \frac{p}{p-1}$$

$$= \exp\left(\frac{1}{p}\log p + \frac{p-1}{p}\log \frac{p}{p-1}\right)$$

$$\leq \frac{1}{p} \cdot p + \frac{p-1}{p} \cdot \frac{p}{p-1}$$

$$= 2.$$

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

J Experimental details

Here, we describe our implementation of the SOIDS algorithm in more detail, as well as the hyperparameters of all the methods used in our experiments. To run the SOIDS algorithm, one must minimise $\overline{\operatorname{IR}}_t^{(2)}(\pi)$ w.r.t. π in each round t. This is not straightforward, because $\overline{\operatorname{IR}}_t^{(2)}(\pi)$ contains expectations w.r.t. the optimistic posterior Q_t^+ . When we use the Spike-and-Slab prior in Appendix B.2, we are not aware of any efficient method that can be used to maximise $\overline{\operatorname{IR}}_t^{(2)}(\pi)$. Instead, we draw (approximate) samples $\theta^{(1)},\ldots,\theta^{(M)}$ from Q_t^+ to produce the estimates $\widetilde{\Delta}_t(\pi)$ and $\overline{\operatorname{IG}}_t(\pi)$ for the surrogate regret and the surrogate information respectively, where

$$\widetilde{\Delta}_t(\pi) = \sum_{a \in \mathcal{A}} \pi(a) \frac{1}{M} \sum_{i=1}^M \Delta(a, \theta^{(i)}), \qquad \widetilde{\mathrm{IG}}_t(\pi) = \frac{1}{2} \sum_{a \in \mathcal{A}} \pi(a) \frac{1}{M} \sum_{i=1}^M \left(\langle \theta^{(i)} - \overline{\theta}_M, a \rangle \right)^2.$$

Here, $\bar{\theta}_M$ is the sample mean $\frac{1}{M} \sum_{i=1}^M \theta^{(i)}$. We then maximimse the approximate surrogate information ratio $\widetilde{\mathrm{IR}}_t^{(2)}(\pi)$, where

$$\widetilde{\mathrm{IR}}_t^{(2)}(\pi) = \frac{(\widetilde{\Delta}_t(\pi))^2}{\widetilde{\mathrm{IG}}_t(\pi)}.$$

To draw the samples $\theta^{(1)},\ldots,\theta^{(M)}$, we use the empirical Bayesian sparse sampling procedure proposed by Hao et al. [2021], which is designed to draw samples from the Bayesian posterior. To sample from the optimistic posterior, we incorporate the optimistic adjustment into the likelihood. This method replaces the theoretically sound spike-and-slab prior with a relaxation in which the "spikes" are Laplace distributions with small variance, and the "slabs" are Gaussian distributions with large variance. In particular, the density of this prior is

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

Here, $\psi_1(\theta)$ is the density function of a univariate Gaussian distribution, with mean 0 and variance ρ_1 , and ψ_0 is the density function of a univariate Laplace distribution, with mean 0 and scale parameter ρ_0 . $p(\gamma)$ is a product of Bernoulli distributions with mean β . In our experiments, we always use $\rho_1=10$, $\rho_0=0.1$ and $\beta=0.1$. Also, we set the learning rates to $\eta=1/2$ and $\lambda_t=\min(\frac{1}{2},\frac{1}{10}\max(\sqrt{\frac{s\log(edt/s)}{dt}},(\frac{\log(edt/s)}{t})^{2/3}))$.

Implementing the OTCS baseline exactly would require us to compute the means of the distributions played by an exponentially weighted average forecaster with a sparsity prior. These distributions are the same as the optimistic posterior, except $\lambda_t=0$ (i.e. there is no optimistic adjustment). In our implementation of the OTCS baseline, we draw samples using the same empirical Bayesian sparse sampling procedure, and then replace the exact means with the sample means. We use the same

choices for the parameters η , ρ_1 , ρ_0 and β . We set the radii of the confidence sets to the values given in Theorem 4.7 of Clerico et al. [2025]

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