

---

# Optimistic Information Directed Sampling

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We study the problem of online learning in contextual bandit problems where  
2 the loss function is assumed to belong to a known parametric function class.  
3 We propose a new analytic framework for this setting that bridges the Bayesian  
4 theory of information-directed sampling due to [Russo and Van Roy \[2018\]](#) and  
5 the worst-case theory of [Foster, Kakade, Qian, and Rakhlin \[2021\]](#) based on the  
6 decision-estimation coefficient. Drawing from both lines of work, we propose  
7 an algorithmic template called Optimistic Information-Directed Sampling and  
8 show that it can achieve instance-dependent regret guarantees similar to the ones  
9 achievable by the classic Bayesian IDS method, but with the major advantage  
10 of not requiring any Bayesian assumptions. The key technical innovation of our  
11 analysis is introducing an optimistic surrogate model for the regret and using it to  
12 define a frequentist version of the Information Ratio of [Russo and Van Roy \[2018\]](#),  
13 and a less conservative version of the Decision Estimation Coefficient of [Foster  
14 et al. \[2021\]](#).

## 15 1 Introduction

16 We present a framework for the analysis of a family of sequential decision-making algorithms known  
17 as Information-Directed Sampling (**IDS**). First proposed by [Russo and Van Roy \[2018\]](#), **IDS** is a  
18 Bayesian algorithm that selects its policies by optimizing a measure called the *information-ratio*,  
19 which measures the tradeoff between instantaneous regret and information gain about the problem  
20 instance at hand. In a Bayesian setup, both components of the information ratio are explicit functions  
21 of the posterior distribution over models, and can thus be explicitly calculated. As shown by [Russo  
22 and Van Roy \[2018\]](#), the resulting algorithm can guarantee massive statistical gains over more  
23 common approaches like Thompson sampling [[Thompson, 1933](#)] or optimistic exploration methods  
24 [[Lai and Robbins, 1985](#)], and in particular can take advantage of the structure of the problem instance  
25 much more effectively. Realizing the same gains in a non-Bayesian setup (which we will sometimes  
26 call *frequentist*, for lack of a better word) is hard for multiple reasons, the most severe obstacle  
27 being that the true model is entirely unknown and Bayesian posteriors cannot be used to quantify  
28 the uncertainty about the model in a meaningful way. As such, defining appropriate notions of  
29 information gain and information ratio is not straightforward. This is the problem we address in this  
30 paper.

31 Our main contribution is constructing a version of information-directed sampling that is imple-  
32 mentable without Bayesian assumptions, and yields frequentist versions of the same problem-  
33 dependent guarantees as the ones achieved by the original **IDS** method in a Bayesian setup. The  
34 key element in our approach is the introduction of a *surrogate model* that allows for a meaningful  
35 definition of the information ratio that is amenable to a frequentist analysis. This surrogate model is  
36 the function of an optimistically adjusted posterior distribution inspired by the “feel-good Thompson  
37 sampling” algorithm of [Zhang \[2022\]](#), and is used to estimate the components of the information

38 ratio: the regret and the information gain. With these components, it becomes possible to define an  
39 information ratio that is an explicit function of the optimistic posterior, which can then be optimized  
40 to yield a decision-making rule that we call “optimistic information-directed sampling” (**OIDS**).

41 For the sake of concreteness, we focus on the problem of contextual bandits and show that **OIDS**  
42 can not only recover worst-case optimal regret bounds in this case, but also satisfies problem-  
43 dependent guarantees that are commonly referred to as *first-order bounds* [Cesa-Bianchi et al., 2005,  
44 Agarwal et al., 2017, Allen-Zhu et al., 2018, Foster and Krishnamurthy, 2021]. Besides these general  
45 guarantees, we also provide some illustrative examples that show that **OIDS** can reproduce the  
46 expedited learning behavior of **IDS** on easy problems, but without requiring Bayesian assumptions.

47 Our methodology also draws inspiration from the analytic framework of Foster, Kakade, Qian, and  
48 Rakhlin [2021], developed for a very general range of sequential decision-making problems. Their  
49 analysis revolves around the notion of the *decision-estimation coefficient* (DEC), which quantifies  
50 the tradeoffs that need to be made between achieving low regret and gaining information about the  
51 true model in a way that is similar to the information ratio of Russo and Van Roy [2016]. The main  
52 contribution of Foster et al. [2021] is showing that the minimax regret in any sequential decision-  
53 making problem can be lower bounded in terms of the DEC, and they also show that nearly matching  
54 upper bounds can be achieved via a simple algorithm they call *estimation to decisions* (**E2D**). Unlike  
55 the information ratio, the DEC does not make use of a Bayesian posterior to quantify uncertainty,  
56 but is rather defined as a worst-case notion, and as such provides frequentist guarantees that hold  
57 uniformly for all problem instances. However, the worst-case nature of the DEC can also be seen as  
58 an inherent limitation of their framework. In particular, the **E2D** algorithm is also based on the same  
59 conservative notion of regret-information tradeoff, and thus all known guarantees for this algorithm  
60 (and its variants such as the ones proposed by Chen et al., 2022, Foster et al., 2023a,b, Kirschner  
61 et al., 2023) fail to take advantage of problem structures that may facilitate fast learning.

62 Our own framework unifies the advantages of the two threads of literature described above: unlike  
63 **E2D**, it is able to achieve instance-dependent guarantees and learn faster in problems with more  
64 structure, and, unlike standard **IDS**, it can do so without relying on Bayesian assumptions. Our  
65 analysis draws on elements of both lines of work, and also on the techniques introduced by Zhang  
66 [2022], as mentioned above.

67 We are not the first to attempt the generalization of **IDS** beyond the Bayesian setting. Kirschner and  
68 Krause [2018] proposed a frequentist alternative to the information ratio for the special case of loss  
69 functions that are linear in some unknown parameter, and constructed an appropriate version of **IDS**  
70 that is able to take advantage of certain problem structures and obtain guarantees that improve upon  
71 the minimax rates. Their approach has inspired a line of work aiming to prove tighter and tighter  
72 problem-dependent bounds for a range of sequential decision-making problems, but so far all of these  
73 results remained limited to linearly structured losses and observations [Kirschner et al., 2020, 2021,  
74 Hao et al., 2022]. In contrast, our notion of information ratio does not require any specific problem  
75 structure like linearity, and arguably constitutes a more universal generalization of **IDS** beyond the  
76 Bayesian setting.

77 **Notation.** The squared Hellinger distance between two probability distributions  $P$  and  $P'$  (with  
78 a common dominating measure  $Q$ ) is defined as  $\mathcal{D}_H^2(P, P') = \frac{1}{2} \int (\sqrt{\frac{dP}{dQ}} - \sqrt{\frac{dP'}{dQ}})^2 dQ$ , and the  
79 relative entropy (or Kullback–Leibler divergence) as  $\mathcal{D}_{\text{KL}}(P \| P') = \int \log \frac{dP}{dP'} dP$ .

## 80 2 Preliminaries

81 We study contextual bandit problems with finite action spaces and parametric loss functions. The  
82 sequential interaction scheme between the *learner* and the *environment* consists of the following  
83 steps being repeated for a sequence of rounds  $t = 1, 2, \dots, T$ :

- 84 • The environment picks a context  $X_t \in \mathcal{X}$ , possibly using randomization and taking into  
85 account the history of actions, losses and contexts,
- 86 • the learner observes  $X_t$  and picks an action  $A_t \in \mathcal{A}$ , possibly using randomization and  
87 taking into account the history of actions, losses and contexts,

- the learner incurs a loss  $L_t$ , drawn independently of the past from a fixed distribution that depends on  $X_t, A_t$ .

We denote the sigma-algebra generated by the interaction history between the learner and the environment up to the end of round  $t$  as  $\mathcal{F}_t = \sigma(X_1, A_1, L_1, \dots, X_t, A_t, L_t)$ , and the probabilities and expectations conditioned on the history as  $\mathbb{P}_t[\cdot] = \mathbb{P}[\cdot | \mathcal{F}_{t-1}, X_t]$  and  $\mathbb{E}_t[\cdot] = \mathbb{E}[\cdot | \mathcal{F}_{t-1}, X_t]$ .

We will suppose that the action space is finite with cardinality  $|\mathcal{A}| = K$ , and that the loss function belongs to a known parametric class, but is otherwise unknown to the learner. Specifically, we assume that there is a known parameter space  $\Theta$  that parametrizes a class of loss functions  $\ell : \Theta \times \mathcal{X} \times \mathcal{A} \rightarrow \mathbf{R}$ , and a true parameter  $\theta_0 \in \Theta$  such that  $\mathbb{E}_t[L_t | X_t, A_t] = \ell(\theta_0, X_t, A_t)$ . We will refer to this condition as *realizability*. The distribution of random losses under parameter  $\theta$  generated in response to taking action  $a$  in context  $x$  will be denoted by  $p(\theta, x, a)$ , and we will write  $p(\cdot | \theta, x, a)$  to designate the corresponding density with respect to a reference measure (usually the counting measure or Lebesgue measure). Unless stated otherwise, we will assume that the loss distribution is fully supported on the interval  $[0, 1]$  for all parameters  $\theta$ . Furthermore, we will often abbreviate  $\ell(\theta, X_t, a)$  as  $\ell_t(\theta, a)$  and  $p(\theta, X_t, a)$  as  $p_t(\theta, a)$  to lighten our notation. Our formulation will make central use of *policies* which prescribe randomized behavior rules for the learning agent. Precisely, a policy  $\pi : \mathcal{X} \rightarrow \Delta_{\mathcal{A}}$  maps each context  $x$  to a distribution over actions denoted as  $\pi(\cdot | x)$ . Since we will mostly work with action distributions conditioned on the fixed contexts  $X_t$ , we will mostly represent policies as distributions over actions, and use the same notation  $\pi \in \Delta_{\mathcal{A}}$  for this purpose. We will focus on learning algorithms that, in each round  $t$ , select a randomized policy  $\pi_t \in \Delta_{\mathcal{A}}$  based on the interaction history  $\mathcal{F}_{t-1}$  and  $X_t$ . We also define the *optimal loss* in round  $t$  under model parameter  $\theta$  as  $\ell_t^*(\theta) = \min_a \ell_t(\theta, a)$ . The agent aims to make its decisions in a way in that minimizes the expected sum of losses, and in particular aims to incur nearly as little loss as the true optimal policy. The extent to which the learner succeeds in achieving this goal is measured by the (total expected) *regret* defined as

$$R_T(\theta_0) = \mathbb{E} \left[ \sum_{t=1}^T (\ell_t(\theta_0, A_t) - \ell_t^*(\theta_0)) \right]. \quad (1)$$

The expectation is over all sources of randomness: the agent’s randomization over actions, the adversary’s randomization over contexts and the randomness of the realization of the losses. We also define *instantaneous regret* of an action  $a$  under parameter  $\theta$  for each  $t$  as

$$r_t(a; \theta) = \ell_t(\theta, a) - \ell_t^*(\theta),$$

and the instantaneous regret of policy  $\pi$  as  $r_t(\pi; \theta) = \sum_a \pi(a) r_t(a; \theta)$ . With this notation, the regret of the online learning algorithm can be written as  $R_T(\theta_0) = \mathbb{E} \left[ \sum_{t=1}^T r_t(\pi_t; \theta_0) \right]$ .

### 3 Two competing theories of sequential decision making

Our work connects two well-established analytic frameworks for sequential decision making: the Bayesian framework of [Russo and Van Roy \[2018\]](#) and the worst-case framework of [Foster, Kakade, Qian, and Rakhlin \[2021\]](#). We review the two in some detail below, highlighting some of their merits and limitations that we address in this paper.

#### 3.1 The information ratio and Bayesian information-directed sampling

The influential work of [Russo and Van Roy \[2016, 2018\]](#) set forth an analytic framework based on a Bayesian learning paradigm where the true model parameter  $\theta_0$  is supposed to be sampled from a known prior distribution  $Q_0 \in \Delta_{\Theta}$ , and the performance of the learner is measured on expectation with respect to this random choice of instance. We refer to the expected regret under this prior as the *Bayesian regret*. Their work has established that the Bayesian regret of any algorithm can be upper bounded in terms of a quantity called the *Information Ratio* (IR). For the sake of exposition, we will follow the setup and notation of [Neu et al. \[2022\]](#), who study the Bayesian version of our contextual bandit setting, and define the information ratio of policy  $\pi$  in the  $t$ -th round of interaction as

$$\rho_t(\pi) = \frac{(\mathbb{E}_{\theta_0 \sim Q_0} [r_t(\pi; \theta_0)])^2}{\text{IG}_t(\pi)}. \quad (2)$$

132 In the above expression, both the numerator and the denominator are functions of the *posterior*  
 133 *distribution*  $Q_t$  of the parameter  $\theta_0$ , computed based on all information available to the learner up to  
 134 the beginning of round  $t$ . Specifically, the numerator is the squared expected regret in round  $t$ , where  
 135 the expectation is taken under the posterior distribution  $Q_t$ , and the denominator is an appropriately  
 136 defined measure of *information gain* that serves to quantify the amount of new information revealed  
 137 about  $\theta_0$  after having observed the latest loss  $L_t$ . The information gain is formally defined as

$$\text{IG}_t(\pi) = \sum_a \pi(a) \int \mathcal{D}_{\text{KL}}(p_t(\theta, a) \parallel \bar{p}_t(a)) dQ_t(\theta), \quad (3)$$

138 where  $\bar{p}_t(a) = \int p_t(\theta, a) dQ_t(\theta)$  is the posterior predictive distribution of the loss  $L_t$  given that  
 139 action  $a$  is played in context  $X_t$ . In other words, the information gain is the *mutual information*  
 140 between the posterior-sample parameter  $\theta_t \sim Q_t$  and a randomly sampled loss  $\tilde{L}_t \sim p_t(\theta_t, a)$ .

141 Given the above definitions, [Russo and Van Roy \[2016, 2018\]](#) show that the Bayesian regret of any  
 142 algorithm can be upper bounded as follows:

$$\mathbb{E}_{\theta_0 \sim Q_0} [R_T(\theta_0)] \leq \sqrt{\mathbb{E} \left[ \sum_{t=1}^T \rho_t(\pi_t) \right] \cdot \mathbb{E} \left[ \sum_{t=1}^T \text{IG}(\pi_t) \right]}. \quad (4)$$

143 The second sum above can be upper bounded by the entropy of  $\theta_0$  under the prior distribution,  
 144 regardless of what algorithm is used to select the sequence of policies. This suggests that one can  
 145 achieve low regret by picking the sequence of policies in a way that minimizes the information ratio:  
 146  $\pi_t = \arg \min_{\pi} \rho_t(\pi)$ . This algorithm is called *information-directed sampling (IDS)*, and has been  
 147 shown to achieve regret guarantees that often improve significantly over worst-case bounds achieved  
 148 by more traditional methods based on posterior sampling or optimistic exploration methods. In  
 149 particular, for the contextual bandit setting we study in this paper, the works of [Neu et al. \[2022\]](#)  
 150 and [Min and Russo \[2023\]](#) have shown that the information ratio of **IDS** is bounded by the number  
 151 of actions  $K$ . When the parameter space is finite with cardinality  $N$ , this result implies that the  
 152 algorithm achieves the minimax optimal regret bound of  $\mathcal{O}(\sqrt{KT \log N})$  for this Bayesian setting.

153 Despite their appealing properties, **IDS**-style methods have however remained largely limited to  
 154 the Bayesian setting, as there appears to be no universal way of defining an algorithmically useful  
 155 information ratio without Bayesian assumptions. In particular, the instantaneous regret  $r_t(\pi; \theta_0)$   
 156 cannot be computed without knowledge of  $\theta_0$ , and there is no reason to believe that the information  
 157 gain defined in terms of a Bayesian posterior would meaningfully measure the reduction in uncertainty  
 158 about  $\theta_0$  in this more general setting.

### 159 3.2 The decision-estimation coefficient and the estimations-to-decisions algorithm

160 The fundamental work of [Foster et al. \[2021\]](#) provides a general theory of sequential decision making,  
 161 providing a range of upper and lower bounds depending on a quantity they call the *decision-estimation*  
 162 *coefficient* (DEC). With a little deviation from their notation and terminology, the DEC associated  
 163 with a policy  $\pi$ , a model class  $\Theta$  and a “reference model”  $\hat{p}_t : \mathcal{A} \rightarrow \Delta_{\mathbb{R}}$  is defined as

$$\text{DEC}_{\gamma,t}(\pi; \Theta, \hat{p}) = \sup_{\theta \in \Theta} \sum_a \pi(a) (\ell(\theta, X_t, a) - \ell(\theta, X_t, \pi_{\theta}) - \gamma \mathcal{D}_H^2(p_t(\theta, a), \hat{p}_t(a))), \quad (5)$$

164 where  $\gamma > 0$  is a trade-off parameter. With this notation, [Foster et al. \[2021\]](#) define the decision-  
 165 estimation coefficient associated with the model class  $\Theta$  as

$$\text{DEC}_{\gamma}(\Theta) = \sup_t \sup_{\hat{p} \in \Delta_{\mathbb{R}}} \inf_{\pi \in \Delta_{\mathcal{A}}} \text{DEC}_{\gamma,t}(\pi; \Theta, \hat{p}).$$

166 Besides the remarkable feat of showing that the minimax regret can be lower bounded in terms  
 167 of the above quantity, they also show that nearly matching upper bounds can be achieved via a  
 168 simple algorithm they call *estimation to decisions (E2D)*. In each round  $t$ , **E2D** takes as input a  
 169 reference model  $\hat{p}_t$  and outputs the policy achieving the minimum in the definition of the DEC:  
 170  $\pi_t = \arg \min_{\pi} \text{DEC}_{\gamma,t}(\pi; \Theta, \hat{p}_t)$ . They show that the regret of this method can be upper bounded in  
 171 terms of the DEC as follows:

$$R_T(\theta_0) \leq \text{DEC}_{\gamma}(\Theta) \cdot T + \gamma \sum_{t=1}^T \mathcal{D}_H^2(p_t(\theta_0, a), \hat{p}_t(a)).$$

172 This shows that the regret of **E2D** can be upper bounded as the sum of the DEC of the model class  $\Theta$   
 173 and the total *estimation error* associated with the sequence of predictions  $\hat{p}_t$  (measured in terms of  
 174 Hellinger distance). For the contextual bandit setting with finite parameter class of size  $N$ , they show  
 175 that the total estimation error can be upper bounded by  $\gamma \log N$  (under an appropriate choice of the  
 176 predictions  $\hat{p}_t$ ), and that the DEC is upper bounded by  $K/\gamma$ , which once again recovers the minimax  
 177 optimal rate of order  $\mathcal{O}(\sqrt{KT \log N})$  when  $\gamma$  is tuned correctly.

178 A significant problem with the approach outlined above is that the DEC is an inherently worst-case  
 179 measure of complexity due to the supremum taken over  $\theta$  in its definition (5). Since the **E2D** algorithm  
 180 itself is based on this possibly loose bound on the regret-to-information gap, this looseness may not  
 181 only affect the bound but also the actual performance of the algorithm. Intuitively, one may hope to  
 182 be able to do better by replacing the supremum over model parameters by only considering models  
 183 that are still “statistically plausible” in an appropriate sense. In what follows, we provide an algorithm  
 184 that realizes this potential.

## 185 4 Optimistic information-directed sampling

186 Our approach solves the issues outlined in the previous sections with both the Bayesian information  
 187 ratio and the decision estimation coefficient. In particular, our method will extend Bayesian **IDS**  
 188 by being able to provide non-Bayesian performance guarantees, and will be able to address the  
 189 over-conservative nature of the DEC and provide strong instance-dependent guarantees.

190 Following Zhang [2022], we start by defining the *optimistic posterior*  $Q_t^+ \in \Delta_\Theta$  via the following  
 191 recursive update rule (starting from an arbitrary prior  $Q_1^+(\theta) \in \Delta_\Theta$ ):

$$\frac{dQ_{t+1}^+}{dQ_t^+}(\theta) \propto (p_t(L_t|\theta, A_t))^\eta \cdot \exp(-\lambda \cdot \ell_t^*(\theta)). \quad (6)$$

192 Here,  $\eta$  and  $\lambda$  are positive constants that will be specified later. For now, we will only say that  $\eta$   
 193 should be thought of as a “large” constant of order 1, and  $\lambda$  as a “small” parameter of order  $1/\sqrt{T}$   
 194 in the worst case. To proceed, we define the *optimistic posterior predictive distribution* of the loss for  
 195 each  $t$  and  $a$  as the mixture  $\bar{p}_t(a) = \int p_t(\theta, a) dQ_t^+(\theta)$ , and the *surrogate loss function* and *surrogate*  
 196 *optimal loss function* respectively as

$$\bar{\ell}_t(a) = \int \ell_t(\theta, a) dQ_t^+(\theta) \quad \text{and} \quad \bar{\ell}_t^* = \int \ell_t^*(\theta) dQ_t^+(\theta). \quad (7)$$

197 In words, these quantities are averages with respect to a mixture model over all contextual bandit  
 198 instances with mixture weights given by the optimistic posterior  $Q_t^+$ . Notably, they are *improper*  
 199 estimators of the true likelihood, loss, and optimal loss functions respectively, as there may be  
 200 no single  $\theta \in \Theta$  that corresponds to these exact functions (unless one assumes certain convexity  
 201 properties of the relevant objects). With these notations, we define the *surrogate regret* of policy  $\pi$  in  
 202 round  $t$  as  $\bar{r}_t(\pi) = \bar{\ell}_t(\pi) - \bar{\ell}_t^*$ . As we will see in the analysis, the optimistic posterior plays a key  
 203 role in ensuring that the surrogate regret does not overestimate the true regret by too much on average,  
 204 which makes it a sensible target for minimization.

205 It remains to define our notion of information gain that we will call *surrogate information gain*.  
 206 Formally, this quantity is defined for each policy  $\pi$  as follows:

$$\bar{\text{IG}}_t(\pi) = \sum_{a \in \mathcal{A}} \pi(a) \int \mathcal{D}_H^2(p_t(\theta, a), \bar{p}_t(a)) dQ_t^+(\theta). \quad (8)$$

207 Notably, this definition matches the original notion of information gain used by Russo and Van Roy  
 208 [2016, 2018], up to the differences that the divergence being used is the squared Hellinger divergence  
 209 instead of Shannon’s relative entropy, and that the expectation is taken over the optimistic posterior  
 210 instead of the plain Bayesian posterior. We will sometimes write  $\bar{r}_t(\pi; Q_t^+)$  and  $\bar{\text{IG}}_t(\pi; Q_t^+)$  to  
 211 emphasize that these are functions of the optimistic posterior  $Q_t^+$ . With the above definitions, we  
 212 are now ready to introduce the central quantity of our algorithmic framework and our analysis: the  
 213 *surrogate information ratio* defined for each policy  $\pi$  as

$$\bar{\text{IR}}_t(\pi) = \frac{(\bar{r}_t(\pi))^2}{\bar{\text{IG}}_t(\pi)} = \frac{(\sum_{a \in \mathcal{A}} \pi(a) \int (\ell_t(\theta, a) - \bar{\ell}_t^*(\theta)) dQ_t^+(\theta))^2}{\sum_{a \in \mathcal{A}} \pi(a) \int \mathcal{D}_H^2(\bar{p}_t(a), p_t(\theta, a)) dQ_t^+(\theta)}. \quad (9)$$

214 Importantly, computing the surrogate information ratio does not require knowledge of  $\theta_0$ : both its  
 215 denominator and numerator can be expressed in terms of the optimistic posterior  $Q_t^+$ . To emphasize  
 216 this fact, we will sometimes write  $\overline{\text{IR}}_t(\pi; Q_t^+)$  for  $\overline{\text{IR}}_t(\pi)$ .

217 We will also define the “offset” counterpart of the surrogate information ratio that is more closely  
 218 related to the decision-estimation coefficient of Foster et al. [2021]. Following the terminology  
 219 introduced in Section 3.2, we introduce the *averaged decision-estimation coefficient* (ADEC) of  
 220 policy  $\pi$  for each  $\mu > 0$  as

$$\begin{aligned} \overline{\text{DEC}}_{\mu,t}(\pi) &= \bar{r}_t(\pi) - \mu \cdot \overline{\text{IG}}_t(\pi) \\ &= \sum_a \pi(a) \int (\ell_t(\theta, \pi) - \ell_t^*(\theta) - \mu \mathcal{D}_H^2(\ell_t(\theta, \pi), \bar{\ell}_t(\pi))) dQ_t^+(\theta). \end{aligned} \quad (10)$$

221 Once again, we also define the notation  $\overline{\text{DEC}}_{\mu,t}(\pi; Q_t^+) = \overline{\text{DEC}}_{\mu,t}(\pi)$  to emphasize the dependence  
 222 of the ADEC on the posterior distribution  $Q_t^+$ . This definition departs from the classic DEC in that,  
 223 instead of taking a supremum over model parameters, it is defined via an expectation with respect to  
 224 the optimistic posterior, thus preventing overly conservative choices of  $\theta$ . It should be clear from  
 225 this definition that the ADEC is always smaller than its original counterpart defined by Foster et al.  
 226 [2021], as long the latter uses the optimistic posterior predictive distribution as its reference model:  
 227  $\overline{\text{DEC}}_{\mu,t}(\pi; Q_t^+) \leq \text{DEC}_{\mu,t}(\pi; \bar{p}_t, \Theta)$ .

228 The surrogate information ratio and the ADEC are related to each other by the inequality

$$\overline{\text{DEC}}_{\mu,t}(\pi) \leq \frac{\overline{\text{IR}}_t(\pi)}{4\mu} \quad (11)$$

229 that holds for all  $\mu > 0$ . Conversely, it can be seen that

$$\overline{\text{IR}}_t(\pi) = \inf \left\{ C > 0 : \overline{\text{DEC}}_{\mu,t}(\pi) \leq \frac{C}{4\mu} \quad (\forall \mu > 0) \right\}. \quad (12)$$

230 These are both direct consequences of the inequality of arithmetic and geometric means. That is,  
 231 whenever the ADEC behaves as  $C_t/\mu$  for all  $\mu$ , the surrogate information ratio succinctly summarizes  
 232 its behavior at all levels  $\mu$ . We will dedicate special attention to this case below, but we also note that  
 233 there are several important cases where the ADEC behaves differently, and the information ratio is a  
 234 less appropriate notion of complexity. We defer further discussion of this to Section 7.

235 With the above notions, we are now ready to define the algorithmic framework we study in this  
 236 paper, with two separate versions depending on whether we consider the surrogate information ratio  
 237 or the average DEC as the basis of decision making. Both versions are referred to as *optimistic*  
 238 *information-directed sampling* (optimistic **IDS** or **OIDS**). Following the terminology of Hao and  
 239 Lattimore [2022], we call the first variant which selects its policies as  $\pi_t = \arg \min_{\pi} \overline{\text{IR}}(\pi; Q_t^+)$   
 240 *vanilla optimistic information-directed sampling* (**VOIDS**), and the second variant that selects  $\pi_t =$   
 241  $\arg \min_{\pi} \overline{\text{DEC}}_{\mu}(\pi; Q_t^+)$  *regularized optimistic information-directed sampling* (**ROIDS**). We provide  
 242 the pseudocode for these methods for quick reference in Appendix A.

## 243 5 Main results

244 We now present our main results regarding the two varieties of our optimistic **IDS** algorithm. We  
 245 first show a general worst-case regret bound stated in terms of the time horizon  $T$  and the information  
 246 ratio. More importantly, we also show instance-dependent guarantees on the performance of **OIDS**  
 247 that replace the scaling with  $T$  in the upper bounds by the total loss of the best policy after  $T$  steps.  
 248 For simplicity of exposition and easy comparison with existing results, we will present our main  
 249 results assuming that the parameter space  $\Theta$  is finite with cardinality  $N$ , and that the losses are almost  
 250 surely bounded in the interval  $[0, 1]$ . We extend these results to compact metric parameter spaces in  
 251 Section 5.3, and provide an extension to subgaussian losses in Section 5.4. Besides these general  
 252 results, we also present several examples where **OIDS** can achieve very low regret by exploiting  
 253 various flavors of problem structure, in Appendix B.

### 254 5.1 Worst-case bounds

255 We start by stating a general worst-case regret bound that relates the regret of any algorithm to its  
 256 surrogate information ratio. This result is the non-Bayesian counterpart of the bounds stated in Russo

257 and Van Roy [2018], Hao and Lattimore [2022] and Neu et al. [2022] in that it basically says that any  
 258 algorithm with bounded information ratio will enjoy bounded regret.

259 **Theorem 1.** Assume  $|\Theta| = N < \infty$  and let  $\lambda > 0$  be arbitrary. Then, for any choice of prior  
 260  $Q_1 \in \Delta_\Theta$ , the regret of any algorithm satisfies the following bound:

$$\begin{aligned} \mathbb{E}[R_T(\theta_0)] &\leq \frac{\log \frac{1}{Q_1(\theta_0)}}{\lambda} + \lambda T \cdot \left( \frac{\sum_{t=1}^T \mathbb{E}[\overline{\text{DEC}}_{1/10\lambda,t}(\pi_t; Q_t^+)]}{\lambda T} + \frac{21}{4} \right) \\ &\leq \frac{\log \frac{1}{Q_1(\theta_0)}}{\lambda} + \lambda T \cdot \left( 10 \cdot \frac{\sum_{t=1}^T \mathbb{E}[\overline{\text{IR}}_t(\pi_t; Q_t^+)]}{T} + \frac{21}{4} \right). \end{aligned} \quad (13)$$

261 We provide a proof sketch, with pointers to the full technical proof details, in Section 6.1. As is  
 262 common in the information directed sampling literature, we will turn this guarantee into a more  
 263 concrete bound on the regret of **OIDS** by exhibiting a “forerunner” algorithm that is able to control  
 264 the surrogate information ratio and is relatively easier to analyze. Indeed, this will certify a regret  
 265 bound for **OIDS**, since the latter precisely minimizes the surrogate information ratio at every round,  
 266 and as such is guaranteed to achieve the same or a better bound. In particular, we use the *feel-good*  
 267 *Thompson sampling* (**FGTS**) algorithm of Zhang [2022] as our forerunner, which samples a parameter  
 268  $\theta_t$  from the optimistic posterior and then plays the policy  $\pi_t = \arg \max_\pi \sum_a \pi(a) \ell_t(\theta_t, a)$ .

269 **Lemma 1.** The surrogate information ratio and averaged decision-to-estimation-coefficient of  
 270 **VOIDS** and **ROIDS** satisfy for any  $\mu \geq 0$

$$4\mu \overline{\text{DEC}}_{\mu,t}(\mathbf{ROIDS}) \leq 4\mu \overline{\text{DEC}}_{\mu,t}(\mathbf{VOIDS}) \leq \overline{\text{IR}}_t(\mathbf{VOIDS}) \leq \overline{\text{IR}}_t(\mathbf{FGTS}) \leq 8K. \quad (14)$$

271 We note that the above result is more of a property of the posterior sampling policy than **FGTS** itself,  
 272 as the bound holds for any distribution that is handed to **OIDS**. This result is not especially new:  
 273 similar statements have been proven in a variety of papers including Russo and Van Roy [2016, 2018],  
 274 Zhang [2022], Foster et al. [2021], Neu et al. [2022]. We provide a proof in Appendix E.4.1. Putting  
 275 the two previous results together, we get the following upper bound on the regret of **OIDS**:

276 **Corollary 1.** Assume  $|\Theta| = N < \infty$ , and let  $\lambda = \sqrt{\frac{\log N}{(80K + \frac{21}{4})T}}$ . Then, the regret of **ROIDS** with  
 277 input parameter  $\mu = \frac{1}{10\lambda}$  and **VOIDS** both satisfy

$$\mathbb{E}[R_T] \leq \sqrt{(320K + 21) T \log N}. \quad (15)$$

278 In particular, this recovers the minimax optimal rate of  $\mathcal{O}(\sqrt{KT \log N})$  for this problem.

## 279 5.2 First-order bounds

280 We now present a more interesting result that replaces the dependence on  $T$  in the previous bound by  
 281 the cumulative loss of the best policy—constituting an instance-dependent guarantee that is often  
 282 called *first-order regret bound*. In particular, in the important class of “noiseless” problems where the  
 283 optimal loss is zero, the result implies that **OIDS** achieves constant regret.

284 **Theorem 2.** Assume  $|\Theta| = N < \infty$ , let  $L^*$  be such that  $\mathbb{E}[\sum_{t=1}^T \ell_t^*(\theta_0)] \leq L^*$ , and let  $\lambda =$   
 285  $\sqrt{\frac{5 \log N}{(500K + 108)L^*}} \wedge \frac{1}{250K + 54}$ . Then the regret of **ROIDS** with input parameter  $\mu = \frac{1}{10\lambda}$  and **VOIDS**  
 286 both satisfy

$$\mathbb{E}[R_T] \leq \sqrt{(2500K + 540) \log N L^*} + (1250K + 270) \log N. \quad (16)$$

287 We provide a proof in Appendix D.1.

## 288 5.3 Infinite parameter spaces

289 We extend the result of Theorem 1 to work for infinite parameter spaces. For simplicity, we focus on  
 290 the case in which  $\Theta$  is a bounded subset of a finite-dimensional vector space.

291 **Theorem 3.** Assume  $\Theta \subset \mathbb{R}^d$ ,  $\max_{x,y \in \Theta} \|x - y\| = 2R < \infty$ . Assume that for all  $x \in X, a \in \mathcal{A}$ ,  
 292 and  $L \in [0, 1]$ , the log-likelihood of the losses  $p(\cdot, x, a, L)$  is  $C$ -Lipschitz. Assume that a ball of

293 radius  $\frac{1}{C_T}$  containing  $\theta_0$  is included in  $\Theta$  and set  $\lambda = \sqrt{\frac{2d \log(RCT)}{(20K + \frac{21}{4})T}}$  and  $Q_1$  a uniform prior on  $\Theta$ .  
 294 Then the regret of **ROIDS** with input parameter  $\mu = \frac{1}{10\lambda}$  and **VOIDS** both satisfy

$$\mathbb{E}[R_T] \leq \sqrt{(160K + 42)dT \log(CRT)} + 1 = \mathcal{O}(\sqrt{dKT \log(CRT)}). \quad (17)$$

295 We provide a proof in Appendix D.2.

## 296 5.4 Subgaussian losses

297 We also extend the basic result of Theorem 1 to work for a more general family of losses. In particular,  
 298 we drop the assumption that the likelihood model is well-specified and allow the losses to be sub-  
 299 Gaussian. As the following result shows, we can still recover our regret bound of  $\mathcal{O}(\sqrt{KT \log N})$   
 300 with some minor tweaks of the algorithm and the analysis. The resulting method is called **OIDS-SG**,  
 301 and is presented in Appendix D.3 in full detail, along with the proof of the theorem below.

302 **Theorem 4.** Assume that the losses are  $v$ -sub-Gaussian, that  $|\Theta| = N < \infty$  and set  
 303  $\lambda = \sqrt{\frac{\log N}{(\frac{1}{4} + 20(v \wedge 1)(1+K))T}}$ . Then the regret of **ROIDS-SG** with input parameter  $\mu = \frac{1}{80\lambda(v \wedge 1)}$   
 304 and **VOIDS-SG** both satisfy

$$\mathbb{E}[R_T] \leq \sqrt{(1 + 80(v \vee 1)(1+K))T \log N} = \mathcal{O}(\sqrt{KT \log N}). \quad (18)$$

## 305 6 Analysis

306 This section provides an outline of the proofs of our main results. We first give a high-level overview  
 307 of the key ideas that are shared in all proofs, and then fill in provide further technical details that are  
 308 required to prove Theorems 1. Theorems 2, 3 and 4 are proved in Appendices D.1, D.2 and D.3.

309 The core of our analysis is the following decomposition of the instantaneous regret in round  $t$ :

$$\begin{aligned} \mathbb{E}[r_t(\pi_t; \theta_0)] &= \mathbb{E}[\bar{r}_t(\pi_t)] + \mathbb{E}[r_t(\pi_t; \theta_0) - \bar{r}_t(\pi_t)] \\ &= \mathbb{E}[\bar{r}_t(\pi_t)] + \mathbb{E}[\mathbb{E}_t[\ell_t(\theta_0, A_t) - \bar{\ell}_t(A_t)]] + \mathbb{E}[\bar{\ell}_t^* - \ell_t^*(\theta_0)] \\ &= \mathbb{E}[\overline{\text{DEC}}_{\mu,t}(\pi_t) + \mu \overline{\text{IG}}_t(\pi_t) + \text{UE}_t + \text{OG}_t]. \end{aligned} \quad (19)$$

310 Here, in the last line we have introduced the notations  $\text{UE}_t = \mathbb{E}_t[\ell_t(\theta_0, A_t) - \bar{\ell}_t(A_t)]$  to denote the  
 311 *underestimation error* of the losses incurred by our own policy  $\pi_t$ , and  $\text{OG}_t = \bar{\ell}_t^* - \ell_t^*(\theta_0)$  as the  
 312 *optimality gap* between the best loss possible in our mixture of models and the optimal loss attainable  
 313 under the true parameter. The first term is small if the mixture model accurately evaluates the losses  
 314 seen during learning (which is generally easy to ensure on average), and the second term is small  
 315 if the model remains optimistic about the best attainable performance (which is facilitated by the  
 316 optimistic adjustment to the posterior updates). An important quantity in the analysis is the (true)  
 317 *information gain* of policy  $\pi$  defined as

$$\text{IG}_t(\pi) = \sum_{a \in \mathcal{A}} \pi(a) \int \mathcal{D}_H^2(p_t(\theta_0, a, \cdot), p_t(\theta, a, \cdot)) \, dQ_t^+(\theta). \quad (20)$$

318 This quantity is closely related to the surrogate information gain that is optimized by our algorithm,  
 319 and plays a key role in bounding the underestimation errors. In particular, the following simple  
 320 lemma establishes a connection between the true and surrogate information gains:

321 **Lemma 2.** For any  $t$  and policy  $\pi$ , the information gain satisfies  $\overline{\text{IG}}_t(\pi) \leq 4\text{IG}_t(\pi)$ .

322 The proof can be found in Appendix E.2.1. Notably, the proof makes critical use of properties of the  
 323 squared Hellinger distance, and is the main reason that the surrogate information gain is defined the  
 324 way it is. In particular, the proof uses the fact that the Hellinger distance is a metric and as such it  
 325 satisfies the triangle inequality—which is the reason that we were not able to go with the otherwise  
 326 more natural choice of relative entropy in our definition of the information gain.



327 **6.1 The proof of Theorem 1**

328 We first use the following worst-case bound on the underestimation error:

329 **Lemma 3.** For any  $t$  and  $\gamma > 0$ , the underestimation error is bounded as  $|\text{UE}_t| \leq \frac{\gamma}{2} + \frac{\text{IG}_t(\pi_t)}{\gamma}$ .

330 The proof is relegated to Appendix E.1.1. Putting this bound together with the previous derivations,  
 331 we get a regret bound that only depends on the averaged Decision-to-Estimation-Coefficient, the  
 332 information gain and the optimality gap:

$$\mathbb{E}[r_t] \leq \mathbb{E} \left[ \overline{\text{DEC}}_{\mu,t}(\pi_t) + \left(4\mu + \frac{1}{\gamma}\right) \text{IG}_t(\pi_t) + \text{OG}_t \right] + \frac{\gamma}{2}. \quad (21)$$

333 Following the terminology of Foster et al. [2023b], we will refer to the sum  $(4\mu + \frac{1}{\gamma})\text{IG}_t(\pi_t) + \text{OG}_t$   
 334 as the *optimistic estimation error*. The following result establishes that the optimistic posterior  
 335 updates can effectively control a quantity that is closely related to this term.

336 **Lemma 4.** Let  $0 < \eta < \frac{1}{2}$ ,  $\lambda > 0$ , and  $\beta = \frac{1}{1-2\eta}$ . Then, the following inequality holds :

$$\mathbb{E} \left[ \sum_{t=1}^T \left( \frac{2\eta}{\lambda} \cdot \text{IG}_t(\pi_t) + \text{OG}_t \right) \right] \leq \frac{\log \frac{1}{Q_1(\theta_0)}}{\lambda} + \frac{\lambda\beta T}{8}. \quad (22)$$

337 See Appendix E.3.1 for the proof. It remains to pick the hyperparameters in a way that the left-  
 338 hand side matches the total optimistic estimation error, which is achieved when setting way that  
 339  $\frac{2\eta}{\lambda} = 4\mu + \frac{1}{\gamma}$ . To make sure that this holds while minimizing the final constant, we choose  $\eta = \frac{1}{4}$ ,  
 340  $\beta = 2$ , and  $\gamma = \frac{1}{\mu} = 10\lambda$ . Plugging these constants into the bound above, and putting the result  
 341 together with the bound of Equation (21) completes the proof of Theorem 1.

342 **7 Conclusion**

343 We have proposed a new analysis framework that bridges the concepts of information ratio and  
 344 decision-estimation coefficient, and unifies the advantages of both frameworks. We conclude by  
 345 discussing some directions of future work. We expand our discussion of related works and open  
 346 problems in Appendix C.

347 A very important open question is whether our notion of averaged DEC can also serve as a lower  
 348 bound on the minimax regret like its original version proposed by Foster et al. [2021]. Since the  
 349 ADEC is a lower bound on the DEC under a special choice of nominal model, we conjecture that  
 350 it can also be used to lower bound the minimax regret in the same “low-probability” fashion as the  
 351 original results of Foster et al. [2021]. On the same note, we remark that it seems unlikely that our  
 352 DEC variant can be reconciled with the “constrained DEC” of Foster et al. [2023a], which has so far  
 353 yielded the tightest lower bounds on the regret within this family of complexity notions. Whether or  
 354 not the averaging idea we advocate for in this paper will turn out to be useful for fully characterizing  
 355 the minimax regret in sequential decision making remains to be seen.

356 It is interesting to observe that the optimistic posterior updates used by our method simplify drastically  
 357 in the special case of “noiseless” problems where  $\ell^*(\theta, X_t) = 0$  holds for all  $\theta$ . This condition holds  
 358 in two of the examples discussed in Appendix B, and more broadly in all problems where the optimal  
 359 policy is guaranteed to achieve zero loss under all candidate parameters  $\theta$ . As a more concrete example,  
 360 we highlight the problem of bandit linear classification with surrogate losses, which satisfies this  
 361 condition if the data is separable with a margin [Kakade et al., 2008, Beygelzimer et al., 2017, 2019].

362 In such noise-free problems, the optimistic posterior update collapses to  $\frac{dQ_{t+1}^+}{dQ_t^+}(\theta) \propto (p_t(L_t|\theta, A_t))^\eta$ ,  
 363 which is closer to the standard Bayesian update up to the important difference that it involves the  
 364 “stepsize” parameter  $\eta$ . Interestingly, such “generalized” or “safe” Bayesian updates have been  
 365 studied extensively in the context of statistical learning under misspecified models—see, e.g., Zhang  
 366 [2006a,b], Grünwald [2012], de Heide et al. [2020]. This connection leads to a multitude of questions  
 367 that we cannot hope to address in this short discussion, so we close with mentioning only one aspect  
 368 that we find to be particularly exciting. Specifically, we wonder if the techniques established in these  
 369 works could be useful for addressing misspecification in the context of sequential decision making  
 370 under uncertainty, where this issue has been notoriously hard to formalize and handle [Du et al., 2019,  
 371 Lattimore et al., 2020, Weisz et al., 2021]. We leave this exciting question open for future research.

372 **References**

- 373 N. Abe and P. M. Long. Associative reinforcement learning using linear probabilistic concepts. In  
374 *International Conference on Machine Learning*, pages 3–11, 1999.
- 375 A. Agarwal, A. Krishnamurthy, J. Langford, and H. Luo. Open problem: First-order regret bounds  
376 for contextual bandits. In *Conference on Learning Theory*, pages 4–7, 2017.
- 377 S. Agrawal and N. Goyal. Further optimal regret bounds for Thompson sampling. In *Artificial  
378 Intelligence and Statistics*, pages 99–107, 2013.
- 379 Z. Allen-Zhu, S. Bubeck, and Y. Li. Make the minority great again: First-order regret bound for  
380 contextual bandits. In *International Conference on Machine Learning*, pages 186–194, 2018.
- 381 A. Beygelzimer, F. Orabona, and C. Zhang. Efficient online bandit multiclass learning with  $\tilde{O}(\sqrt{T})$   
382 regret. In *International Conference on Machine Learning*, pages 488–497, 2017.
- 383 A. Beygelzimer, D. Pál, B. Szörényi, D. Thiruvengatchari, C.-Y. Wei, and C. Zhang. Bandit  
384 multiclass linear classification: Efficient algorithms for the separable case. In *International  
385 Conference on Machine Learning*, pages 624–633, 2019.
- 386 S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities - A Nonasymptotic Theory  
387 of Independence*. Oxford University Press, 2013. ISBN 978-0-19-953525-5. doi: 10.1093/  
388 ACPROF:OSO/9780199535255.001.0001. URL [https://doi.org/10.1093/acprof:  
389 oso/9780199535255.001.0001](https://doi.org/10.1093/acprof:oso/9780199535255.001.0001).
- 390 S. Bubeck and M. Sellke. First-order Bayesian regret analysis of Thompson sampling. In *Algorithmic  
391 Learning Theory*, pages 196–233, 2020.
- 392 N. Cesa-Bianchi, Y. Mansour, and G. Stoltz. Improved second-order bounds for prediction with  
393 expert advice. In *Conference on Learning Theory*, volume 3559 of *Lecture Notes in Computer  
394 Science*, pages 217–232. Springer, 2005.
- 395 F. Chen, S. Mei, and Y. Bai. Unified Algorithms for RL with Decision-Estimation Coefficients:  
396 No-Regret, PAC, and Reward-Free Learning, 2022.
- 397 R. de Heide, A. Kirichenko, P. Grünwald, and N. Mehta. Safe-bayesian generalized linear regression.  
398 In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*,  
399 pages 2623–2633, 2020.
- 400 S. S. Du, S. M. Kakade, R. Wang, and L. F. Yang. Is a good representation sufficient for sample  
401 efficient reinforcement learning? In *International Conference on Learning Representations*, 2019.
- 402 D. Foster and A. Rakhlin. Beyond UCB: Optimal and efficient contextual bandits with regression  
403 oracles. In *International Conference on Machine Learning*, pages 3199–3210, 2020.
- 404 D. J. Foster and A. Krishnamurthy. Efficient first-order contextual bandits: Prediction, allocation, and  
405 triangular discrimination. *Advances in Neural Information Processing Systems*, 34:18907–18919,  
406 2021.
- 407 D. J. Foster, S. M. Kakade, J. Qian, and A. Rakhlin. The Statistical Complexity of Interactive  
408 Decision Making, 2021.
- 409 D. J. Foster, N. Golowich, and Y. Han. Tight guarantees for interactive decision making with the  
410 decision-estimation coefficient. *arXiv preprint arXiv:2301.08215*, 2023a.
- 411 D. J. Foster, N. Golowich, J. Qian, A. Rakhlin, and A. Sekhari. Model-free reinforcement learning  
412 with the decision-estimation coefficient. In *Thirty-seventh Conference on Neural Information  
413 Processing Systems*, 2023b.
- 414 P. Grünwald. The safe Bayesian: learning the learning rate via the mixability gap. In *Algorithmic  
415 Learning Theory*, pages 169–183, 2012.
- 416 B. Hao and T. Lattimore. Regret bounds for information-directed reinforcement learning. *Advances  
417 in Neural Information Processing Systems*, 35:28575–28587, 2022.

- 418 B. Hao, T. Lattimore, and C. Qin. Contextual information-directed sampling. In *International*  
419 *Conference on Machine Learning*, pages 8446–8464, 2022.
- 420 S. M. Kakade, S. Shalev-Shwartz, and A. Tewari. Efficient bandit algorithms for online multiclass  
421 prediction. In *Proceedings of the 25th international conference on Machine learning*, pages  
422 440–447, 2008.
- 423 J. Kirschner and A. Krause. Information directed sampling and bandits with heteroscedastic noise. In  
424 *Conference On Learning Theory*, pages 358–384, 2018.
- 425 J. Kirschner, T. Lattimore, and A. Krause. Information directed sampling for linear partial monitoring.  
426 In *Conference on Learning Theory*, pages 2328–2369, 2020.
- 427 J. Kirschner, T. Lattimore, C. Vernade, and C. Szepesvári. Asymptotically optimal information-  
428 directed sampling. In *Conference on Learning Theory*, pages 2777–2821, 2021.
- 429 J. Kirschner, A. Bakhtiari, K. Chandak, V. Tkachuk, and C. Szepesvári. Regret minimization via  
430 saddle point optimization. In *Thirty-seventh Conference on Neural Information Processing Systems*,  
431 2023.
- 432 T. L. Lai and H. Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied*  
433 *Mathematics*, 6:4–22, 1985.
- 434 T. Lattimore and A. György. Mirror Descent and the Information Ratio. In *Conference on Learning*  
435 *Theory*, volume 134, pages 2965–2992, 2021.
- 436 T. Lattimore, C. Szepesvári, and G. Weisz. Learning with good feature representations in bandits and  
437 in rl with a generative model. In *International Conference on Machine Learning*, pages 5662–5670,  
438 2020.
- 439 S. Min and D. Russo. An information-theoretic analysis of nonstationary bandit learning. In  
440 *International Conference on Machine Learning*, 2023.
- 441 G. Neu, J. Olkhovskaya, M. Papini, and L. Schwartz. Lifting the information ratio: An information-  
442 theoretic analysis of Thompson sampling for contextual bandits. *Advances in Neural Information*  
443 *Processing Systems*, 35:9486–9498, 2022.
- 444 J. Olkhovskaya, J. Mayo, T. van Erven, G. Neu, and C.-Y. Wei. First- and second-order bounds  
445 for adversarial linear contextual bandits. In *Thirty-seventh Conference on Neural Information*  
446 *Processing Systems*, 2023.
- 447 D. Russo and B. Van Roy. An information-theoretic analysis of Thompson sampling. *J. Mach. Learn.*  
448 *Res.*, 17:68:1–68:30, 2016.
- 449 D. Russo and B. Van Roy. Learning to optimize via information-directed sampling. *Oper. Res.*, 66  
450 (1):230–252, 2018.
- 451 S. Shalev-Shwartz. *Online learning: Theory, algorithms, and applications*. Hebrew University, 2007.
- 452 W. Thompson. On the likelihood that one unknown probability exceeds another in view of the  
453 evidence of two samples. *Bulletin of the American Mathematics Society*, 25:285–294, 1933.
- 454 G. Weisz, P. Amortila, and C. Szepesvári. Exponential lower bounds for planning in mdps with  
455 linearly-realizable optimal action-value functions. In *Algorithmic Learning Theory*, pages 1237–  
456 1264, 2021.
- 457 T. Zhang. From  $\epsilon$ -entropy to KL-entropy: Analysis of minimum information complexity  
458 density estimation. *The Annals of Statistics*, 34(5), 2006a.
- 459 T. Zhang. Information-theoretic upper and lower bounds for statistical estimation. *IEEE Transactions*  
460 *on Information Theory*, 52(4):1307–1321, 2006b.
- 461 T. Zhang. Feel-good Thompson sampling for contextual bandits and reinforcement learning. *SIAM*  
462 *Journal on Mathematics of Data Science*, 4(2):834–857, 2022.

463 **A Pseudocode of OIDS**

464 We provide the pseudocode for OIDS in Algorithm 1 below.

---

**Algorithm 1** Optimistic Information Directed Sampling (**OIDS**)

---

**Input:** prior  $Q_1^+$ , parameters  $\eta, \lambda, \mu$ .

**For**  $t = 1, \dots, T$ , **repeat:**

1. Observe context  $X_t$ ,
  - 2a. **VOIDS:** play policy  $\pi_t = \arg \min_{\pi \in \Delta(\mathcal{A})} \overline{\text{IR}}_t(\pi, Q_t^+)$ ,
  - 2b. **ROIDS:** play policy  $\pi_t = \arg \min_{\pi \in \Delta(\mathcal{A})} \overline{\text{DEC}}_t(\pi, Q_t^+, \mu)$ ,
  3. incur loss  $L_t$ ,
  4. update optimistic prior,  $Q_{t+1}^+(\cdot) \propto Q_t^+(\cdot)(p_t(\cdot, A_t, L_t))^\eta \exp(-\lambda \ell_t^*(\cdot))$ .
- 

465 **B Examples**

466 The most appealing property of **IDS** in the Bayesian setting is that it can take advantage of the  
 467 structure of the problem at hand to achieve extremely good performance that is otherwise not  
 468 achievable by methods like Thompson sampling or UCB. Indeed, unlike these methods, **IDS** has  
 469 the ability to pick actions that are not optimal under any statistically plausible model, but can reveal  
 470 useful information about the problem. Russo and Van Roy [2018] demonstrate several examples of  
 471 situations where **IDS** provably achieves massive speedups via such queries. It is not clear that such  
 472 speedups are achievable without Bayesian assumptions, although some evidence was offered by the  
 473 work of Kirschner and Krause [2018] in the case of linear rewards. In this section, we demonstrate  
 474 that our version of **IDS** can fully reproduce the fast learning behavior of Bayesian **IDS** on the original  
 475 examples of Russo and Van Roy [2018], thus suggesting that **OIDS** may inherit many more good  
 476 properties of its Bayesian counterpart than what our main theoretical results show. We also provide  
 477 an additional example on which we demonstrate that **OIDS** can outperform DEC-based methods by  
 478 addressing the over-conservatism encoded in the definition of the DEC.

479 **B.1 Revealing action**

480 We first adapt the “revealing actions” example of the original work of Russo and Van Roy [2018].  
 481 This example features the action set  $\mathcal{A} = \{0, 1, \dots, K\}$ , the set of parameters  $\Theta = \{1, \dots, K\}$ ,  
 482 and the loss function  $\ell(\theta, a) = \mathbb{I}_{\{a>0, a \neq \theta\}} + \mathbb{I}_{\{a=0\}}(1 - \frac{1}{2^a})$ . The losses are deterministic and the  
 483 agent gets loss 0 by picking the action corresponding to the unknown parameter  $\theta_0$ . Action 0 is  
 484 special, it results in a large loss that however encodes the identity of the optimal action. Thus, the  
 485 optimal exploration strategy is to pick this revealing action once, read out the identity of the optimal  
 486 action, and play that action until the end of time. Russo and Van Roy [2018] show that **IDS** follows  
 487 this exact strategy, and here we show that **OIDS** does the same when taking as input a (completely  
 488 noninformative) uniform prior over the parameters.

489 To show this, we will compute for any action the surrogate reward and surrogate information gain  
 490 under the optimistic posterior (which is identical to the uniform prior, given that we are in the first  
 491 round). For  $a \neq 0$ , the surrogate regret is written as

$$\bar{r}_1(a) = \int_{\Theta} \ell(\theta, a) - \ell(\theta) dQ_0(\theta) = \frac{1}{K} \sum_{\theta=1}^K (1 - \mathbb{I}_{\{a=\theta\}}) = 1 - \frac{1}{K},$$

492 while for the revealing action, the surrogate regret is

$$\bar{r}_1(0) = 1 - \frac{1}{K} + \frac{2^{-K}}{K}.$$

493 In particular  $\bar{r}_t(0) > \bar{r}_t(a)$  so the action 0 has the worst expected reward under our model. As for the  
 494 information gain, we an explicit computation of the Hellinger distance for  $a \neq 0$  shows

$$\text{IG}_t(a) = \frac{1}{K} \cdot \left(1 - \sqrt{\frac{1}{K}}\right) + \frac{K-1}{K} \cdot \left(1 - \sqrt{\frac{K-1}{K}}\right) = \mathcal{O}\left(\frac{1}{K}\right).$$

495 Meanwhile, for action 0 we have

$$\text{IG}_t(0) = 1 - \sqrt{\frac{1}{K}} = \Theta(1).$$

## 496 B.2 Sparse linear model

497 Our second example is a linear bandit problem where the action space corresponds to a finite subset  
 498 of the Euclidean unit ball  $\mathcal{A} = \{\frac{x}{\|x\|_1} : x \in \{0, 1\}^d, x \neq 0\}$ , the parameter space consists of the set  
 499 of coordinate vectors  $\Theta = \{\theta' \in \{0, 1\}^d, \|\theta'\|_1 = 1\}$ , and the loss function is  $\ell(\theta, a) = 1 - \langle a, \theta \rangle$ .  
 500 As in the previous example, the losses are again deterministic. This is a linear bandit problem where  
 501 the parameter  $\theta$  is known to be 1-sparse. In particular, the optimal action under the model  $\theta$  consists  
 502 in only selecting action  $a = \theta$  so any Thompson Sampling based algorithm will only select one  
 503 coordinate at a time and will take up to  $d$  steps to determine the true parameter  $\theta_0$ . In contrast, the  
 504 optimal exploration policy will perform binary search on the action space and find the optimal action  
 505 exponentially faster.

506 To investigate the behaviour of **OIDS** on this problem, we will compute the surrogate regret and  
 507 surrogate information gain of an action  $a$ . Since our prior is uniform, we have

$$\bar{r}_1(a) = \bar{\ell}_1(a) = \mathbb{P}[\langle \theta_0, a \rangle > 0] \cdot \frac{1}{\|a\|_1} = \frac{\|a\|_1}{d} \cdot \frac{1}{\|a\|_1} = \frac{1}{d}$$

508 and

$$\begin{aligned} \text{IG}_1(a) &= \frac{\|a\|_1}{d} \cdot \left(1 - \sqrt{\frac{\|a\|_1}{d}}\right) + \frac{d - \|a\|_1}{d} \cdot \left(1 - \sqrt{\frac{d - \|a\|_1}{d}}\right) \\ &= 1 - \left(\frac{\|a\|_1}{d}\right)^{\frac{3}{2}} - \left(1 - \frac{\|a\|_1}{d}\right)^{\frac{3}{2}} \end{aligned}$$

509 Thus, the expected reward of all actions is the same, and the information gain is maximized for  
 510 actions with norm  $\|a\|_1 = \frac{d}{2}$ . **IDS** thus picks an action  $A_1$  uniformly at random and updates the  
 511 posterior as follows. If the observed loss is 1, all parameters with  $\langle \theta, A_1 \rangle > 0$  will be eliminated by  
 512 the posterior update. If the observed loss is smaller than 1, all parameters satisfying  $\langle \theta, A_1 \rangle = 0$  are  
 513 excluded. The posterior is thus set as uniform over all surviving parameters and the process repeats.  
 514 Continuing along the same lines, we can see that both versions of **OIDS** will continue performing  
 515 binary search and identify the true parameter in  $\log_2 d$  time steps.

## 516 B.3 Bandits with a revelatory zero

517 Our final example is a multi-armed bandit problem where the losses keep looking exactly the same  
 518 until a low-probability event happens that reveals the optimal action perfectly. In this setup (vaguely  
 519 inspired by Example 3.3 of Foster et al., 2021),  $\Theta = [K]$ , and the losses are defined as uniformly  
 520 distributed random variables in  $[0, 1]$  for all actions except  $a = \theta$ . For this special action, the loss is  
 521 defined as  $B_t U_t$ , with  $U_t$  uniform on  $[0, 1]$ , and  $B_t$  is Bernoulli with mean  $1 - 2\Delta \in [0, 1]$ . The mean  
 522 loss for this action is  $\frac{1}{2} - \Delta$ . For this model, there is essentially no way for any algorithm to discover  
 523 the optimal action until the first time that a loss of zero is observed. In this case, the (optimistic)  
 524 posterior immediately collapses on  $\theta_0$ . Consequently, **OIDS** keeps drawing uniform random actions  
 525 until the first zero is observed, and plays the optimal action in all remaining rounds. The number of  
 526 time steps spent with uniform exploration are geometrically distributed with mean  $\frac{K}{2\Delta}$ , thus making  
 527 for a total regret of approximately  $\frac{K}{2}$ . Note that in this instance, the optimistic adjustment to the  
 528 posterior is not necessary as the optimal loss of all models are the same, so the performance of the  
 529 algorithm is unaffected by the choice of  $\lambda$  or  $\mu$ .

530 Interestingly, the **E2D** algorithm of Foster et al. [2021] cannot take advantage of the structure of this  
 531 problem so effectively. When using the posterior predictive distribution  $\bar{p}_t$  as the nominal model, the  
 532 Hellinger distance will approximately behave as  $\mathcal{D}_H^2(p(\theta, a), \bar{p}_t(a)) \approx \mathbb{1}_{\{\theta \neq \theta_0\}}$  after observing the

533 first zero. Thus, the worst-case DEC associated with policy  $\pi$  is written as

$$\begin{aligned} \text{DEC}_\gamma(\pi; \bar{p}_t, \Theta) &= \sup_{\theta} \{ \ell(\theta, \pi) - \ell(\theta, a_\theta) - \gamma \mathbb{I}_{\{\theta \neq \theta_0\}} \} = \sup_{\theta} \left\{ \Delta \sum_{a \neq \theta} \pi(a) - \gamma \mathbb{I}_{\{\theta \neq \theta_0\}} \right\} \\ &= \sup_{\theta} \{ \Delta(1 - \pi(\theta)) - \gamma \mathbb{I}_{\{\theta \neq \theta_0\}} \}. \end{aligned}$$

534 When  $\gamma \geq \Delta$ , the expression in the supremum can be positive for certain policies  $\pi$  and parameters  
 535  $\theta \neq \theta_0$ , and thus the  $\theta$  player will prefer picking  $\theta \neq \theta_0$  for some policies. More precisely, the DEC  
 536 for any policy will be given as

$$\text{DEC}(\pi; \Theta, \hat{p}_t) = \max \left\{ \Delta(1 - \min_{a \neq \theta_0} \pi(a)) - \gamma, \Delta(1 - \pi(\theta_0)) \right\}.$$

537 In the extreme case  $\gamma = 0$ , the policy achieving maximum value is approximately uniform, and it  
 538 approximates the optimal policy  $\pi^*$  gradually as  $\gamma$  increases. When  $\gamma$  is large enough, the alternative  
 539  $\theta \neq \theta_0$  stops being attractive to the max player and **E2D** starts outputting  $\pi^*$ . This happens at  
 540 the threshold  $\gamma > \Delta$  at the latest. This observation matches the discussion of [Foster et al. \[2021,](#)  
 541 [Example 3.3\]](#) and [Foster et al. \[2023a, p. 8\]](#), who demonstrate the same threshold behavior of the  
 542 DEC and point out that this leads to tight lower bounds, without discussing the potential shortcomings  
 543 of **E2D** that prevents it from obtaining tight upper bounds. It is easy to see that **E2D** fails because of  
 544 the over-conservative definition of the DEC: while there is sufficient evidence to reject all alternative  
 545 parameters, **E2D** still computes its optimization objective by taking a supremum over *all* model  
 546 parameters  $\theta$ , including ones that have already been ruled out by the observations. This clearly  
 547 demonstrates the advantage of the surrogate model used by **OIDS**, which computes its objective with  
 548 the help of the optimistic posterior distribution that allows faster elimination of unlikely parameters.

## 549 C Further discussion

550 In this section, we expand our discussion of related works and open questions.

551 **General bounded losses.** At the surface level, it may seem that our results only apply to well-  
 552 specified models where the likelihood model correctly captures the distribution of the random losses.  
 553 This is of course a very restrictive assumption. However, it is easy to see that our framework can tackle  
 554 arbitrary bounded losses via a standard binarization trick [[Agrawal and Goyal, 2013](#)]: supposing  
 555 that the losses are bounded in  $[0, 1]$ , they can be randomly rounded to  $\{0, 1\}$  to apply **OIDS** with a  
 556 Bernoulli likelihood. It is easy to see that the regret bounds for these post-processed losses continue  
 557 to hold for the original losses as well. We presume that our approach can be generalized beyond such  
 558 sub-Bernoulli and sub-Gaussian losses to more general sub-exponential-family losses, but we leave  
 559 the investigation of this generalization open for future work.

560 **Beyond contextual bandits.** For the sake of simplicity, we have presented our results within the  
 561 relatively modest framework of contextual bandits. That said, it is clear that our framework can be  
 562 generalized to the much broader setting of “decision making with structured observations” studied  
 563 by [Foster et al. \[2021\]](#), and that it can be used to prove regret bounds of the form of [Theorem 1](#)  
 564 straightforwardly in said setting. However, so far we could only prove quantitative improvements  
 565 over the DEC for contextual bandits, and thus we decided not to let down the reader by introducing a  
 566 very general setting and then only providing interesting results in a narrow special case. Nevertheless,  
 567 our results demonstrate that our framework can achieve strictly superior upper bounds on the regret in  
 568 a highly nontrivial setting that has been studied extensively (see, e.g., [Agarwal et al., 2017](#), [Allen-Zhu](#)  
 569 [et al., 2018](#), [Foster and Krishnamurthy, 2021](#), [Bubeck and Sellke, 2020](#), [Olkhovskaya et al., 2023](#)).

570 **Multiplicative or additive tradeoff?** All of our results are stated in terms of both the surrogate  
 571 information ratio, which measures the regret-to-information tradeoff multiplicatively, and the averaged  
 572 DEC, which does so in an additive fashion. Based on these results, it is not immediately clear which  
 573 of the two notions is more useful. Equations (11) and (12) suggest that the ADEC is always smaller  
 574 than the information ratio, which may suggest that it may yield better guarantees. To a certain degree,  
 575 [Russo and Van Roy \[2018\]](#) have already addressed this question: their Proposition 11 shows that

576 measuring the regret-information tradeoff additively results in strictly *worse* regret for a range of  
577 hyperparameter choices. While at the surface, this seems to defy the intuition provided our results, in  
578 reality their additive tradeoff is only vaguely related to the one we consider, and the regularization  
579 range for which the result holds does not seem to be practical in the first place. On the other hand,  
580 Foster et al. [2021] make a more robust argument against the information ratio in comparison with the  
581 DEC, showing that there are some hard problems for which the information ratio is infinite but the  
582 DEC remains finite (see their Section 9.3). Besides the fact that their information ratio is defined in  
583 an unorthodox way via the same conservative supremum as what appears in the definition of the DEC,  
584 this claim seems to miss some important follow-up work on **IDS** that has already addressed this  
585 issue. Specifically, Lattimore and György [2021] have pointed out that the information ratio is only  
586 suitable for problems where the minimax regret is of the order  $\sqrt{T}$  (which one can already notice by  
587 inspecting the general bound of Equation 4), and studying harder games with larger minimax regret  
588 may be done by introducing a generalized notion of information ratio that features a different power  
589 of the regret in the denominator. In the present paper, we decided to stay impartial and state our  
590 results for both flavors of optimistic **IDS**, and we hope that this debate will progress productively in  
591 the future.

592 **Connection with the Bayesian DEC.** The attentive reader may have noticed that a notion closely  
593 related to our averaged DEC has already been mentioned in the original work of Foster et al.  
594 [2021]. Indeed, their Section 4.2 proposes a Bayesian version of the **E2D** algorithm that optimizes  
595  $\overline{\text{DEC}}_{\gamma,t}(\cdot; Q_t)$ , where  $Q_t$  is the exact Bayesian posterior over the model parameters. They show that  
596 the resulting algorithm enjoys essentially the same guarantees on the Bayesian regret as the worst-  
597 case guarantees obtained by the standard **E2D** method. Our approach effectively considers the same  
598 optimization objective, with the important change that the standard Bayesian posterior is replaced  
599 with the optimistic posterior of Zhang [2022]. This not only strengthens the mentioned results of  
600 Foster et al. [2021] by removing the Bayesian assumption necessary for its analysis, but also allows  
601 us to obtain instance-dependent guarantees as well. We believe that the same instance-dependent  
602 improvements (and more) should be directly provable for the Bayesian **E2D** method of Foster et al.  
603 [2021], but we did not pursue this direction as we preferred to focus on pointwise regret guarantees  
604 this time.

## 605 D Proofs of the main results

606 We now give the complete proofs of our main results. We relegate most of the technical content into  
607 Appendix E and only provide the main arguments here for better readability.

### 608 D.1 The proof of Theorem 2

609 We start our analysis from the regret decomposition of Equation (19) and apply Lemma 2 to obtain

$$\mathbb{E}[r_t] \leq \mathbb{E}[\overline{\text{DEC}}_{\mu,t}(\pi_t) + 4\mu\text{IG}_t(\pi_t) + \text{UE}_t + \text{OG}_t].$$

610 As before, we can control the ADEC of **OIDS** by producing a suitable forerunner. In particular, we  
611 use the *inverse-gap weighting* **IGW** algorithm of Foster and Krishnamurthy [2021]

612 **Lemma 5.** *The surrogate information ratio and averaged decision-to-estimation-coefficient of*  
613 ***VOIDS** and **ROIDS** satisfy for any  $\mu \geq 0$*

$$4\mu\overline{\text{DEC}}_{\mu,t}(\mathbf{ROIDS}) \leq 4\mu\overline{\text{DEC}}_{\mu,t}(\mathbf{VOIDS}) \leq \overline{\text{IR}}_t(\mathbf{VOIDS}) \leq \overline{\text{IR}}_t(\mathbf{IGW}) \leq 40K \min_{a \in \mathcal{A}} \bar{\ell}_t(a). \quad (23)$$

614 See Appendix E.4.2 for a definition of the (**IGW**) algorithm and the proof. The term on the right-hand  
615 side can be further bounded as

$$\overline{\text{DEC}}_{\mu,t}(\pi_t) \leq \frac{10K}{\mu} \min_a \bar{\ell}_t(a) \leq \frac{10K}{\mu} (\mathbb{E}_t[\bar{\ell}_t(A_t)]) = \frac{10K}{\mu} (\mathbb{E}_t[\ell_t(\theta_0, A_t)] - \text{UE}_t)$$

616 The final tool is a refined version of Lemma 3 that controls the underestimation error in terms of the  
617 information gain and the current estimate of the loss.

618 **Lemma 6.** *For any  $t$  and  $\gamma > 0$ , the underestimation error is bounded as*

$$\text{UE}_t \leq \frac{\text{IG}_t(\pi_t)}{\gamma} + 2\gamma\mathbb{E}_t[\ell_t(\theta_0, A_t)]. \quad (24)$$

619 See Appendix E.1.2 for the proof. Putting this together with the previous regret decomposition, as  
 620 long as  $\frac{10K}{\mu} \leq 1$ , we get:

$$\mathbb{E}[r_t] \leq \mathbb{E} \left[ \left( 4\mu + \frac{1}{\gamma} \cdot \left( 1 - \frac{10K}{\mu} \right) \right) \mathbf{IG}_t(\pi_t) + \mathbf{OG}_t + \left( 2\gamma \left( 1 - \frac{10K}{\mu} \right) + \frac{10K}{\mu} \right) \ell_t(\theta_0, A_t) \right], \quad (25)$$

621 As before, we will regard the term  $(4\mu + \frac{1}{\gamma} \cdot (1 - \frac{10K}{\mu}))\mathbf{IG}_t + \mathbf{OG}_t$  as the optimistic estimation error,  
 622 and adapt Lemma 4 to provide a refined bound on this quantity:

623 **Lemma 7.** *Let  $0 < \eta < \frac{1}{2}$ ,  $\lambda > 0$ , and  $\beta = \frac{1}{1-2\eta}$ . Then, the optimistic estimation error satisfies*

$$\sum_{t=1}^T \left( \frac{2\eta}{\lambda} \cdot \mathbf{IG}_t(\pi_t) + \left( 1 - \frac{\lambda\beta}{2} \right) \mathbf{OG}_t \right) \leq \frac{\log N}{\lambda} + \frac{\lambda\beta}{2} \sum_{t=1}^T \ell_t^*(\theta_0). \quad (26)$$

624 See Appendix E.3.2 for the proof. The claim of the theorem is then proved by tuning the hyperparam-  
 625 eters in a way that the quantity bounded in the previous Lemma matches the optimistic estimation  
 626 error.

627 Under the condition  $\frac{10K}{\mu} \leq 1$ , the following holds

$$\begin{aligned} \mathbb{E}[r_t] &\leq \mathbb{E} \left[ \left( 4\mu + \frac{1}{\gamma} \cdot \left( 1 - \frac{10K}{\mu} \right) \right) \mathbf{IG}_t(\pi_t) + \mathbf{OG}_t + \left( 2\gamma \left( 1 - \frac{10K}{\mu} \right) + \frac{10K}{\mu} \right) \ell_t(\theta_0, A_t) \right] \\ &\leq \mathbb{E} \left[ \left( 4\mu + \frac{1}{\gamma} \right) \mathbf{IG}_t(\pi_t) + \mathbf{OG}_t + \left( 2\gamma + \frac{10K}{\mu} \right) \ell_t(\theta_0, A_t) \right], \end{aligned}$$

628 where in the last line we also used that  $\mathbf{IG}_t$  and  $\ell_t(\theta_0, A_t)$  are nonnegative to upper bound  $1 - \frac{10K}{\mu} \leq 1$ .

629 In order to apply Lemma 7, we would like to manipulate the above expression so that the coefficients of  
 630  $\mathbf{IG}_t$  and  $\mathbf{OG}_t$  match. To this end, we use the condition that  $\frac{\lambda\beta}{2} \leq \frac{1}{5}$ , which ensures that  $1 \leq \frac{1}{1-\frac{\lambda\beta}{2}} \leq \frac{5}{4}$

631 and thus we can continue the above bound as

$$\mathbb{E}[r_t] \leq \mathbb{E} \left[ \frac{5}{4} \cdot \left( \left( 4\mu + \frac{1}{\gamma} \right) \mathbf{IG}_t(\pi_t) + \left( 1 - \frac{\lambda\beta}{2} \right) \mathbf{OG}_t + \left( 2\gamma + \frac{10K}{\mu} \right) \ell_t(\theta_0, A_t) \right) \right].$$

632 To apply Lemma 7, we choose  $\eta = \frac{1}{4}$ ,  $\beta = 2$ ,  $\gamma = \frac{1}{\mu} = 10\lambda$ , and sum over all rounds to obtain

$$\begin{aligned} \mathbb{E}[R_T] &\leq \mathbb{E} \left[ \frac{5}{4} \cdot \frac{\log N}{\lambda} + \frac{5\lambda}{4} \sum_{t=1}^T \ell_t^*(\theta_0) + (125K + 25)\lambda \sum_{t=1}^T \ell_t(\theta_0, A_t) \right] \\ &\leq \mathbb{E} \left[ \frac{5}{4} \cdot \frac{\log N}{\lambda} + (125K + 27)\lambda \sum_{t=1}^T \ell_t(\theta_0, A_t) \right], \end{aligned}$$

633 where we upper-bounded the optimal loss  $\frac{5\lambda}{4} \ell_t^*(\theta_0)$  by  $2\lambda \ell_t(\theta_0, A_t)$  in the last step. Introducing  
 634 the notation  $\widehat{L}_T = \sum_{t=1}^T \ell_t(\theta^*, A_t)$  and  $L_t^* = \sum_{t=1}^T \ell_t^*(\theta_0)$ , the two sides of the equation can be  
 635 rewritten as

$$R_T = \widehat{L}_T - L_t^* \leq \mathbb{E} \left[ \frac{5}{4} \cdot \frac{\log N}{\lambda} + (125K + 27)\lambda \widehat{L}_T \right].$$

636 Hence, after some reordering we arrive at

$$\mathbb{E}[R_T] \cdot (1 - (125K + 27)\lambda) \leq \mathbb{E} \left[ \frac{5}{4} \cdot \frac{\log N}{\lambda} + (125K + 27)\lambda L_T^* \right].$$

637 If  $\lambda < \frac{1}{2(125K+30)}$ , we can divide both sides of the inequality by  $(1 - (125K + 27)\lambda)$  to obtain

$$\mathbb{E}[R_T] \leq \mathbb{E} \left[ \frac{5}{2} \cdot \frac{\log N}{\lambda} + (250K + 54)\lambda L^* \right],$$

638 where  $L^*$  is an upper bound on  $\mathbb{E}[L_T^*]$ . Finally, we plug the value  $\lambda = \sqrt{\frac{5 \log N}{(500K+108)L^*}} \wedge \frac{1}{250K+54}$   
 639 to get the regret bound of Theorem 2.



640 **D.2 The proof of Theorem 3**

641 The only difference with the finite parameter space analysis is in the control of the optimistic  
 642 estimation error. In particular, we only need to adapt our analysis of the optimistic posterior and  
 643 Lemma 4 to get the regret bound claimed in Theorem 3. We do this with the following lemma.

644 **Lemma 8.** *Let  $0 < \eta < \frac{1}{2}$ ,  $\lambda > 0$ , and  $\beta = \frac{1}{1-2\eta}$ , assume the hypothesis of Theorem 3 hold. Then,*  
 645 *the following inequality holds :*

$$\mathbb{E} \left[ \sum_{t=1}^T \left( \frac{2\eta}{\lambda} \cdot \text{IG}_t(\pi_t) + \text{OG}_t \right) \right] \leq \frac{d \log \frac{R}{\epsilon}}{\lambda} + \frac{\lambda \beta T}{8} + \left( \frac{\eta}{\lambda} + 1 \right) \cdot CT\epsilon. \quad (27)$$

646 The proof is found in Appendix E.3.4. We can now put this together with the regret decomposition  
 647 of Equation (21). As in the proof of Theorem 1, we need to pick the hyperparameters such that the  
 648 optimistic estimation error matches the left hand side of the previous lemma. The same choice of  
 649 hyperparameters  $\eta = \frac{1}{4}$ ,  $\beta = 2$ , and  $\gamma = \frac{1}{\mu} = 10\lambda$  combined with Lemma 1 gives us the following  
 650 bound

$$\mathbb{E} [R_T] \leq \lambda T (20K + \frac{1}{4} + 5) + \frac{d \log \frac{R}{\epsilon}}{\lambda} + \left( \frac{1}{4\lambda} + 1 \right) \cdot CT\epsilon. \quad (28)$$

651 Picking  $\epsilon = 1/(CT)$  gives us

$$\mathbb{E} [R_T] \leq \frac{2d \log RCT}{\lambda} + \lambda T \left( 20K + \frac{21}{4} \right) + 1, \quad (29)$$

652 where we used  $\frac{1}{4} \leq d \log RCT$ . Finally picking  $\lambda = \sqrt{\frac{2d \log(RCT)}{T(20K + \frac{21}{4})}}$  recovers the claim of Theorem 3.

653 **D.3 The proof of Theorem 4**

654 One of the appeals of our approach is that with minor tweaking, we can extend the previous guarantees  
 655 so subgaussian losses. To do that, we consider the following family of likelihoods:

$$p(c|\theta, x, a) \propto \exp \left( -\frac{(c - \ell(\theta, x, a))^2}{2} \right).$$

656 We also readjust our definition of information gain for this setting by replacing the squared Hellinger  
 657 distance by the square loss. In particular, the *Gaussian surrogate information gain* is defined as

$$\overline{\text{IG}}_t^{\mathcal{G}}(\pi) = \sum_{a \in \mathcal{A}} \pi(a) \int (\ell_t(\theta, a) - \bar{\ell}_t(a))^2 dQ_t^+(\theta)$$

658 and the (*true*) *Gaussian information gain* as

$$\text{IG}_t^{\mathcal{G}}(\pi) = \sum_{a \in \mathcal{A}} \pi(a) \int (\ell_t(\theta, a) - \ell_t(\theta_0, a))^2 dQ_t^+(\theta).$$

659 The surrogate information ratio and averaged DEC are adapted as any policy  $\pi$

$$\overline{\text{IR}}_t^{\mathcal{G}}(\pi) = \frac{\bar{r}_t(\pi)}{\overline{\text{IG}}_t^{\mathcal{G}}(\pi)} \quad \text{and} \quad \overline{\text{DEC}}_{\mu, t}^{\mathcal{G}}(\pi) = \bar{r}_t(\pi) - \mu \cdot \overline{\text{IG}}_t^{\mathcal{G}}(\pi). \quad (30)$$

660 Then, we define the corresponding algorithm template (called Optimistic Information Directed  
 661 Sampling for subgaussian losses, **OIDS-SG**) as the method that either picks  $\pi_t$  as the minimizer of  
 662  $\overline{\text{IR}}_t^{\mathcal{G}}$  or  $\overline{\text{DEC}}_T^{\mathcal{G}}$ . The two varieties are referred to as **VOIDS-SG** and **ROIDS-SG**.

663 Replacing the surrogate information gain by its Gaussian counterpart, the regret decomposition of  
 664 Equation (19) is still valid:

$$\mathbb{E} [r_t] = \mathbb{E} \left[ \overline{\text{DEC}}_t^{\mathcal{G}}(\pi_t, \mu) + \mu \overline{\text{IG}}_t^{\mathcal{G}}(\pi_t) + \text{UE}_t + \text{OG}_t \right].$$

665 The surrogate and true information gains are related to each other by the following lemma:

666 **Lemma 9.** For any  $t$  and policy  $\pi$ , the information gain for Gaussians satisfies  $\overline{\text{IG}}_t^{\mathcal{G}}(\pi) \leq 4\text{IG}_t^{\mathcal{G}}(\pi)$ .

667 See Appendix E.2.2 for the proof. We also relate the underestimation error to the information gain  
668 through the following lemma

669 **Lemma 10.** For any  $t$  and  $\gamma > 0$ , the underestimation error is bounded as

$$|\text{UE}_t| \leq \frac{\gamma}{4} + \frac{\text{IG}_t^{\mathcal{G}}(\pi_t)}{\gamma}.$$

670 The proof is presented in Appendix E.1.3. Putting these together, we get a regret bound that only  
671 depends on the average DEC, the information gain and optimality gap:

$$\mathbb{E}[r_t] \leq \mathbb{E}\left[\overline{\text{DEC}}_{\mu,t}^{\mathcal{G}}(\pi_t) + \left(4\mu + \frac{1}{\gamma}\right)\text{IG}_t^{\mathcal{G}}(\pi_t) + \text{OG}_t + \frac{\gamma}{4}\right]. \quad (31)$$

672 We again refer to the sum  $\left(4\mu + \frac{1}{\gamma}\right)\text{IG}_t^{\mathcal{G}}(\pi_t)$  as the optimistic estimation error and will control it  
673 through an analysis of the optimistic posterior adapted to the sub-Gaussianity of the losses. This is  
674 done in the following lemma, whose proof we relegate to Appendix E.3.3.

675 **Lemma 11.** Assume that the losses are  $v$  sub-Gaussian and that for all  $\theta \in \Theta, x \in X, a \in A$ ,  
676  $\ell(\theta, x, a) \in [0, 1]$ , then setting  $\eta = \frac{1+\sqrt{1-1\wedge v}}{2v}$  the following inequality holds :

$$\mathbb{E}\left[\sum_{t=1}^T \frac{1}{16\lambda(v \vee 1)} \cdot \text{IG}_t^{\mathcal{G}}(\pi_t) + \text{OG}_t\right] \leq \frac{\log N}{\lambda} + \frac{\lambda T}{4}. \quad (32)$$

677 Now we pick  $\mu = \frac{1}{\gamma} = \frac{1}{80\lambda(v \vee 1)}$  and apply the previous lemma to obtain the bound

$$\mathbb{E}[R_T] \leq \mathbb{E}\left[\sum_{t=1}^T \overline{\text{DEC}}_{\frac{1}{80\lambda(v \vee 1)},t}^{\mathcal{G}}(\pi_t)\right] + \frac{\log N}{\lambda} + \lambda T \left(\frac{1}{4} + 20(v \vee 1)\right). \quad (33)$$

678 It remains to bound the ADEC. We do this by exhibiting a “forerunner” algorithm that is able to  
679 control the *Surrogate Information Ratio*. In particular, we use again the feel-good Thompson sampling  
680 (**FGTS**) algorithm of Zhang [2022] for this purpose.

681 **Lemma 12.** The surrogate information and averaged decision-to-estimation-coefficient of **OIDS** and  
682 **VOIDS** satisfy the following bound for any  $\mu > 0$ :

$$4\mu\overline{\text{DEC}}_{\mu,t}^{\mathcal{G}}(\mathbf{ROIDS-SG}) \leq 4\mu\overline{\text{DEC}}_{\mu,t}^{\mathcal{G}}(\mathbf{VOIDS-SG}) \leq \overline{\text{IR}}_t^{\mathcal{G}}(\mathbf{VOIDS-SG}) \leq \overline{\text{IR}}_t^{\mathcal{G}}(\mathbf{FGTS}) = K \quad (34)$$

683 Putting everything together, we obtain the bound

$$\mathbb{E}[R_T] \leq \frac{\log N}{\lambda} + \lambda T \left(\frac{1}{4} + 20(v \vee 1)(1 + K)\right), \quad (35)$$

684 from which the bound claimed in Theorem 4 follows by picking the optimal choice of  $\lambda$ .

## 685 E Technical proofs

686 This section presents the more technical parts of the analysis, along with detailed proofs. The content  
687 is organized into four main parts: Appendix E.1 presents techniques for bounding the underestimation  
688 error, Appendix E.2 provides techniques for relating the surrogate information gain to the true  
689 information gain, Appendix E.3 presents the analysis of the optimistic posterior updates to control  
690 the optimistic estimation error, and Appendix E.4 provides bounds on the surrogate information ratio  
691 and the ADEC. All subsections include a variety of results, stated respectively for the worst-case  
692 bounds, first-order bounds, and subgaussian losses.

693 **E.1 Analysis of the Underestimation error**

694 **E.1.1 Worst case analysis: The proof of Lemma 3**

695 We define the total variation distance between two distributions  $P, Q$  sharing a common dominating  
696 measure  $\lambda$  as

$$\text{TV}(P, Q) = \frac{1}{2} \int |p(x) - q(x)| d\lambda(x),$$

697 where  $p, q$  are their densities with respect to  $\lambda$ . The total variation distance can be upper bounded by  
698 the Hellinger distance as follows:

$$\begin{aligned} \text{TV}(P, Q) &= \frac{1}{2} \int \left| (\sqrt{p(x)} - \sqrt{q(x)}) \cdot (\sqrt{p(x)} + \sqrt{q(x)}) \right| d\lambda(x) \\ &\leq \frac{1}{2} \sqrt{\int (\sqrt{p(x)} - \sqrt{q(x)})^2 d\lambda(x) \cdot \int (\sqrt{p(x)} + \sqrt{q(x)})^2 d\lambda(x)} \\ &\leq \frac{1}{2} \sqrt{2\mathcal{D}_H^2(P, Q) \cdot 2 \int (p(x) + q(x)) d\lambda(x)} \\ &= \sqrt{2\mathcal{D}_H^2(P, Q)} \\ &\leq \frac{\gamma}{2} + \frac{\mathcal{D}_H^2(P, Q)}{\gamma}. \end{aligned}$$

699 Here, the first two inequalities follow from applying Cauchy–Schwarz, and the last one from the  
700 inequality of arithmetic and geometric means. Thus, we proceed as

$$\begin{aligned} |\text{UE}_t| &= \left| \sum_a \pi_t(a) \int \ell_t(\theta_0, a) - \ell_t(\theta, a) dQ_t^+(\theta) \right| \\ &\leq \sum_a \pi_t(a) \int |\ell_t(\theta_0, a) - \ell_t(\theta, a)| dQ_t^+(\theta) \\ &= \sum_a \pi_t(a) \int \text{TV}(\text{Ber}(\ell_t(\theta_0, a)), \text{Ber}(\ell_t(\theta, a))) dQ_t^+(\theta) \\ &\leq \sum_a \pi_t(a) \int \text{TV}(p_t(\theta_0, a), p_t(\theta, a)) dQ_t^+(\theta) \\ &\leq \frac{\gamma}{2} + \frac{\sum_a \pi_t(a) \int \mathcal{D}_H^2(p_t(\theta_0, a), p_t(\theta, a)) dQ_t^+(\theta)}{\gamma} \\ &= \frac{\gamma}{2} + \frac{IG_t}{\gamma}. \end{aligned}$$

701 The first inequality above uses the boundedness of the losses in  $[0, 1]$ , the second inequality is the  
702 data-processing inequality for the total variation distance (applied on the noisy channel  $X \rightarrow Y$  that  
703 randomly rounds  $X \in [0, 1]$  to  $Y \in \{0, 1\}$ ), and the last one is the inequality we have just proved  
704 above. This concludes the proof.

705 **E.1.2 Instance-dependent analysis: The proof of Lemma 6**

706 This proof requires a more sophisticated technique based on careful specialized handling of the  
707 “underestimated” and “overestimated” actions. The argument is vaguely inspired by the techniques  
708 of [Bubeck and Sellke \[2020\]](#) and [Foster and Krishnamurthy \[2021\]](#). Specifically, for a parameter  $\theta$ ,  
709 we define  $\mathcal{A}_\theta^- = \{a \in \mathcal{A} : \ell_t(\theta, a) < \ell_t(\theta_0, a)\}$  as the set of actions where  $\ell_t(\theta, a)$  underestimates

710  $\ell_t(\theta_0, a)$ . With this notation, we write

$$\begin{aligned}
\text{UE}_t &= \sum_a \pi_t(a) (\ell_t(\theta_0, a) - \bar{\ell}_t(a)) \\
&= \int \sum_a \pi_t(a) (\ell_t(\theta_0, a) - \ell_t(\theta, a)) dQ_t^+(\theta) \\
&\leq \int \sum_{a \in \mathcal{A}_\theta^-} \pi_t(a) (\ell_t(\theta_0, a) - \ell_t(\theta, a)) dQ_t^+(\theta) \\
&= \int \sum_{a \in \mathcal{A}_\theta^-} \pi_t(a) \cdot \frac{\sqrt{\gamma(\ell_t(\theta_0, a) + \ell_t(\theta, a))}}{\sqrt{\gamma(\ell_t(\theta_0, a) + \ell_t(\theta, a))}} (\ell_t(\theta_0, a) - \ell_t(\theta, a)) dQ_t^+(\theta),
\end{aligned}$$

711 where the inequality follows by dropping the negative terms of the sum. Now, the inequality of  
712 arithmetic and geometric means implies that for any  $x, y \geq 0$ ,  $xy \leq \frac{x^2 + y^2}{2}$ . We apply it to  
713  $x = 2\sqrt{\gamma(\ell_t(\theta_0, a) + \ell_t(\theta, a))}$  and  $y = \frac{(\ell_t(\theta_0, a) - \ell_t(\theta, a))}{2\sqrt{\gamma(\ell_t(\theta_0, a) + \ell_t(\theta, a))}}$  to obtain

$$\text{UE}_t \leq \int \left( \gamma \sum_{a \in \mathcal{A}_\theta^-} \pi_t(a) \cdot (\ell_t(\theta_0, a) + \ell_t(\theta, a)) + \frac{1}{4\gamma} \sum_{a \in \mathcal{A}_\theta^-} \pi_t(a) \frac{(\ell_t(\theta_0, a) - \ell_t(\theta, a))^2}{\ell_t(\theta_0, a) + \ell_t(\theta, a)} \right) dQ_t^+(\theta).$$

714 To proceed, we use the inequality  $\frac{(\ell_t(\theta_0, a) - \ell_t(\theta, a))^2}{\ell_t(\theta_0, a) + \ell_t(\theta, a)} \leq 4\mathcal{D}_H^2(p_t(\theta, a), p_t(\theta_0, a))$  that holds for all  $a$   
715 and  $\theta$ , and is proved separately as Lemma 23. Hence,

$$\begin{aligned}
\text{UE}_t &\leq 2\gamma \sum_a \pi_t(a) \ell_t(\theta_0, a) + \frac{1}{\gamma} \int \sum_a \mathcal{D}_H^2(p_t(\theta, a), p_t(\theta_0, a)) dQ_t^+(\theta) \\
&\leq 2\gamma \sum_a \pi_t(a) \ell_t(\theta_0, a) + \frac{\text{IG}_t}{\gamma},
\end{aligned}$$

716 which concludes the proof.

### 717 E.1.3 Subgaussian analysis: The proof of Lemma 10

718 The claim follows from the following calculations:

$$\begin{aligned}
|\text{UE}_t| &= \left| \sum_a \pi_t(a) \int \ell(\theta_0, a) - \bar{\ell}_t(a) dQ_t^+(\theta) \right| \\
&\leq \sum_a \pi_t(a) \int |\ell(\theta_0, a) - \bar{\ell}_t(a)| dQ_t^+(\theta) \\
&\leq \sqrt{\sum_a \pi_t(a) \int (\ell(\theta_0, a) - \bar{\ell}_t(a))^2 dQ_t^+(\theta)} \\
&= \sqrt{\text{IG}_t^{\mathcal{G}}(\pi_t)} \\
&\leq \frac{\gamma}{4} + \frac{\text{IG}_t^{\mathcal{G}}(\pi_t)}{\gamma}.
\end{aligned}$$

719 Here, the second inequality is Cauchy–Schwarz and the last one is the inequality of arithmetic and  
720 geometric means.

721 **E.2 Analysis of the Surrogate Information Gain and the True Information Gain**

722 **E.2.1 Bounded losses: The proof of Lemma 2**

723 The claim is proved as

$$\begin{aligned}
\overline{\text{IG}}_t(\pi) &= \sum_a \pi(a) \int \mathcal{D}_H^2(\bar{\ell}_t(a), \ell_t(\theta, a)) \, dQ_t^+(\theta) \\
&\leq 2 \cdot \sum_a \pi(a) \int \mathcal{D}_H^2(\bar{\ell}_t(a), \ell_t(\theta_0, a)) \, dQ_t^+(\theta) \\
&\quad + 2 \cdot \sum_a \pi(a) \int \mathcal{D}_H^2(\ell_t(\theta_0, a), \ell_t(\theta, a)) \, dQ_t^+(\theta) \\
&\leq 4 \cdot \sum_a \pi(a) \int \mathcal{D}_H^2(\ell_t(\theta_0, a), \ell_t(\theta, a)) \, dQ_t^+(\theta) \\
&= 4\text{IG}_t(\pi),
\end{aligned}$$

724 where the first inequality critically uses that the Hellinger distance is a metric and as such it satisfies  
725 the triangle inequality, and thus  $\mathcal{D}_H^2(P, P') \leq 2\mathcal{D}_H^2(P, Q) + 2\mathcal{D}_H^2(Q, P')$  holds for any  $P, P'$  and  
726  $Q$  by an additional application of Cauchy–Schwarz. The final inequality then uses the convexity of  
727 the Hellinger distance and Jensen’s inequality.

728 **E.2.2 Subgaussian losses: The proof of Lemma 9**

729 The claims follows from writing

$$\begin{aligned}
\overline{\text{IG}}_t^{\mathcal{G}}(\pi) &= \sum_a \pi(a) \int (\bar{\ell}_t(a) - \ell(\theta, a))^2 \, dQ_t^+(\theta) \\
&\leq 2 \cdot \sum_a \pi(a) \int (\bar{\ell}_t(a) - \ell(\theta_0, a))^2 \, dQ_t^+(\theta) \\
&\quad + 2 \cdot \sum_a \pi(a) \int (\ell(\theta_0, a) - \ell(\theta, a))^2 \, dQ_t^+(\theta) \\
&\leq 4 \cdot \sum_a \pi(a) \int (\ell(\theta_0, a) - \ell(\theta, a))^2 \, dQ_t^+(\theta) \\
&= 4\text{IG}_t^{\mathcal{G}}(\pi),
\end{aligned}$$

730 where the first inequality comes an application of the triangle inequality and Cauchy–Schwarz, and  
731 the second one comes from the convexity of the squared loss and Jensen’s inequality.

732 **E.3 Analysis of the Optimistic Posterior**

733 We start by providing a general statement about the properties of the optimistic posterior updates,  
734 which will then prove useful for bounding the optimistic estimation error.

735 **Lemma 13.** *Consider the optimistic posterior defined recursively by*

$$\frac{dQ_{t+1}^+}{dQ_t^+}(\theta) = \frac{\exp\left(-\eta \log\left(\frac{1}{p_t(L_t|\theta, A_t)}\right) - \lambda \ell_t^*(\theta)\right)}{\int \exp\left(-\eta \log\left(\frac{1}{p_t(L_t|\theta', A_t)}\right) - \lambda \ell_t^*(\theta')\right) \, dQ_t^+(\theta')}, \quad (36)$$

736 where  $Q_1^+ = Q_1$  is some prior distribution on  $\Theta$  and  $p_t(\cdot|\theta, a) \in \Delta_{\mathbb{R}^+}$  is the density the loss  
737 distribution associated with parameter  $\theta$ . For any  $T > 0$ , for any  $\alpha, \beta > 0$  such that  $\frac{1}{\alpha} + \frac{1}{\beta} = 1$ , for  
738 any distribution  $Q^* \in \Delta(\Theta)$ , and for any sequence of actions  $A_1, \dots, A_T$  and losses  $L_1, \dots, L_T$ ,

739 *the following inequality holds:*

$$\begin{aligned}
& -\frac{1}{\lambda\alpha} \sum_{t=1}^T \log \int p_t(\theta, A_t, L_t)^{\eta\alpha} dQ_t^+(\theta) - \frac{1}{\lambda\beta} \sum_{t=1}^T \log \int \exp(-\lambda\beta\ell_t^*(\theta)) dQ_t^+(\theta) \\
& \leq \int \left( \frac{1}{\lambda\alpha} \cdot \sum_{t=1}^T \log \frac{1}{p_t(\theta, A_t, L_t)^{\eta\alpha}} + \sum_{t=1}^T \ell_t^*(\theta) \right) dQ^*(\theta) + \frac{1}{\lambda} \cdot \mathcal{D}_{\text{KL}}(Q^* \| Q_1).
\end{aligned} \tag{37}$$

740 *Proof.* We study the potential function  $\Phi$  defined for all  $c \in \mathbb{R}^\Theta$  as

$$\Phi(c) = \frac{1}{\lambda} \log \int_{\Theta} \exp(-\lambda c(\theta)) dQ_1(\theta).$$

741 We define  $c_t(\theta) = \frac{\eta}{\lambda} \log \frac{1}{p_t(\theta, A_t, L_t)} + \ell_t^*(\theta)$  and evaluate  $\Phi \left( \sum_{t=1}^T c_t \right)$ :

$$\Phi \left( \sum_{t=1}^T c_t \right) = \frac{1}{\lambda} \log \int_{\Theta} \exp \left( -\lambda \sum_{t=1}^T c_t(\theta) \right) dQ_1(\theta) \geq - \int_{\Theta} \sum_{t=1}^T c_t(\theta) dQ^*(\theta) - \frac{\mathcal{D}_{\text{KL}}(Q^* \| Q_1)}{\lambda}.$$

742 where the inequality is the Donsker-Varadhan variational formula [cf. Section 4.9 in [Boucheron et al., 2013](#)]. We also have

$$\begin{aligned}
\Phi \left( \sum_{t=1}^T c_t \right) &= \sum_{t=1}^T \left( \Phi \left( \sum_{k=1}^t c_k \right) - \Phi \left( \sum_{k=1}^{t-1} c_k \right) \right) \\
&= \sum_{t=1}^T \frac{1}{\lambda} \log \frac{\int_{\Theta} \exp \left( -\lambda \sum_{k=1}^t c_k(\theta) \right) dQ_1(\theta)}{\int_{\Theta} \exp \left( -\lambda \sum_{k=1}^{t-1} c_k(\theta) \right) dQ_1(\theta)} \\
&= \sum_{t=1}^T \frac{1}{\lambda} \log \int_{\Theta} \frac{\exp \left( -\lambda \sum_{k=1}^{t-1} c_k(\theta) \right)}{\int_{\Theta} \exp \left( -\lambda \sum_{k=1}^{t-1} c_k(\theta) \right) dQ_1(\theta)} \cdot \exp(-\lambda c_t(\theta)) dQ_1(\theta) \\
&= \sum_{t=1}^T \frac{1}{\lambda} \log \int_{\Theta} \exp(-\lambda c_t(\theta)) dQ_t^+(\theta) \\
&= \sum_{t=1}^T \frac{1}{\lambda} \log \int_{\Theta} p_t(\theta, A_t, L_t)^{\eta} \cdot \exp(-\lambda \ell_t^*(\theta)) dQ_t^+(\theta),
\end{aligned}$$

744 where the fourth equality is by definition of  $Q_t^+$  and  $c_t$ .

745 We can now apply Hölder's inequality with  $\alpha, \beta > 0$  such that  $\frac{1}{\alpha} + \frac{1}{\beta} = 1$ , obtaining

$$\Phi \left( \sum_{t=1}^T c_t \right) \leq \frac{1}{\lambda} \cdot \sum_{t=1}^T \left( \frac{1}{\alpha} \log \int_{\Theta} p_t(\theta, A_t, L_t)^{\eta\alpha} dQ_t^+(\theta) + \frac{1}{\beta} \log \int_{\Theta} \exp(-\lambda\beta\ell_t^*(\theta)) dQ_t^+(\theta) \right).$$

746 Plugging both bounds together, we get the claim of the lemma.  $\square$

747 The following statement will be useful for turning the above guarantee into a bound on the information  
748 gain:

749 **Lemma 14.** *For any  $t \geq 1$  and any policy  $\pi_t \in \Delta(\mathcal{A})$ , the following inequality holds:*

$$\mathbb{E} [\text{IG}_t(\pi_t)] \leq \mathbb{E} \left[ -\log \int_{\Theta} \left( \frac{p_t(\theta, A_t, L_t)}{p_t(\theta_0, A_t, L_t)} \right)^{\frac{1}{2}} dQ_t^+(\theta) \right]. \tag{38}$$

750 *Proof.* Let  $\tau$  be the dominating measure used to define the densities  $p_t(\cdot | \theta, a)$ . We write:

$$\mathbb{E} [\text{IG}_t(\pi_t)] = \mathbb{E} \left[ \int_{\Theta} \sum_a \pi_t(a) \mathcal{D}_H^2(p_t(\theta_0, a), p_t(\theta, a)) dQ_t^+(\theta) \right]$$

$$\begin{aligned}
&= \mathbb{E} \left[ \int_{\Theta} \sum_a \pi_t(a) \left( 1 - \int_{\mathbb{R}} (p_t(c|\theta, A_t) p_t(c|\theta_0, A_t))^{\frac{1}{2}} d\tau(c) \right) dQ_t^+(\theta) \right] \\
&= \mathbb{E} \left[ \int_{\Theta} \mathbb{E}_t \left[ \int_{\mathbb{R}} \left( 1 - \left( \frac{p_t(c|\theta, A_t)}{p_t(c|\theta_0, A_t)} \right)^{\frac{1}{2}} \right) p_t(c|\theta_0, A_t) d\tau(c) \right] dQ_t^+(\theta) \right] \\
&= \mathbb{E} \left[ \int_{\Theta} \mathbb{E}_t \left[ \int_{\mathbb{R}} \left( 1 - \left( \frac{p_t(L_t|\theta, A_t)}{p_t(L_t|\theta_0, A_t)} \right)^{\frac{1}{2}} \right) p_t(L_t|\theta_0, A_t) \right] dQ_t^+(\theta) \right] \\
&\leq \mathbb{E} \left[ \mathbb{E}_t \left[ -\log \int \left( \frac{p_t(L_t|\theta, A_t)}{p_t(L_t|\theta_0, A_t)} \right)^{\frac{1}{2}} dQ_t^+(\theta) \right] \right] \\
&= \mathbb{E} \left[ -\log \int_{\Theta} \left( \frac{p_t(\theta, A_t, L_t)}{p_t(\theta_0, A_t, L_t)} \right)^{\frac{1}{2}} dQ_t^+(\theta) \right].
\end{aligned}$$

751 Here, we used the tower rule of expectation several times, and also the elementary inequality  
752  $\log(x) \leq x - 1$  that holds for all  $x$ . This concludes the proof.  $\square$

### 753 E.3.1 Worst case analysis: The proof of Lemma 4

754 **Lemma 15.** For any  $t \geq 1$ ,  $\beta, \lambda > 0$ , as long as  $\ell_t^*(\theta) \in [0, 1]$  for all values of  $\theta$ , the following  
755 inequality holds

$$756 \frac{1}{\lambda\beta} \log \int_{\Theta} \exp(-\lambda\beta\ell_t^*(\theta)) dQ_t^+(\theta) \leq -\bar{\ell}_t^* + \frac{\lambda\beta}{8}. \quad (39)$$

756 *Proof.* This is a direct consequence of Hoeffding's lemma for bounded random variables, see for  
757 example Section 2.3 of [Boucheron et al. \[2013\]](#).  $\square$

758 The proof of Lemma 4 then follows directly by applying Lemma 13 with  $\eta, \alpha$  such that  $\eta\alpha = \frac{1}{2}$   
759 (which means  $\beta = 1/(1 - 2\eta)$ ) and with  $Q^*$  a dirac distribution in  $\theta_0$ , and combining the result with  
760 Lemmas 14 and 15 above.

### 761 E.3.2 Instance dependent analysis and proof of Lemma 7

762 **Lemma 16.** For any  $t \geq 1$ ,  $\beta, \lambda > 0$ , as long as  $\ell_t^*(\theta) \in [0, 1]$  for all values of  $\theta$ , the following  
763 inequality holds

$$764 \frac{1}{\lambda\beta} \log \int_{\Theta} \exp(-\lambda\beta\ell_t^*(\theta)) dQ_t^+(\theta) \leq -\bar{\ell}_t^* \left( 1 - \frac{\lambda\beta}{2} \right). \quad (40)$$

764 *Proof.* We use the two elementary inequalities  $\log(x) \leq x - 1$  that holds for all  $x \in \mathbb{R}$  and  
765  $e^{-x} \leq 1 - x + \frac{x^2}{2}$  that holds for all  $x \geq 0$  to show

$$\begin{aligned}
\frac{1}{\lambda\beta} \log \int_{\Theta} \exp(-\lambda\beta\ell_t^*(\theta)) dQ_t^+(\theta) &\leq \frac{1}{\lambda\beta} \left( \int_{\Theta} 1 - \lambda\beta\ell_t^*(\theta) + \left( \frac{\lambda\beta}{2} \ell_t^*(\theta) \right)^2 dQ_t^+(\theta) - 1 \right) \\
&\leq \frac{1}{\lambda\beta} \left( \int_{\Theta} -\lambda\beta\ell_t^*(\theta) + \left( \frac{\lambda\beta}{2} \right)^2 \ell_t^*(\theta) dQ_t^+(\theta) \right) \\
&= -\bar{\ell}_t^* \left( 1 - \frac{\lambda\beta}{2} \right),
\end{aligned}$$

766 where we used the fact that for all  $\theta \in \Theta$ , we have  $\ell_t^*(\theta) \in [0, 1]$  and thus  $\ell_t^*(\theta)^2 \leq \ell_t^*(\theta)$ .  $\square$

767 We use again Lemma 13 with  $\eta, \alpha$  such that  $\eta\alpha = 1/2$  and with  $Q^*$  a dirac distribution in  $\theta_0$ . Then  
768 we apply Lemma 16 and Lemma 14 to conclude the proof of Lemma 7.

769 **E.3.3 Subgaussian analysis: The proof of Lemma 11**

770 **Lemma 17.** Assume that the losses are  $v$  sub-Gaussian and that for all  $\theta \in \Theta, x \in \mathcal{X}, a \in$   
 771  $\mathcal{A}, \ell(\theta, x, a) \in [0, 1]$ . For any  $t \geq 1, \eta, \alpha \geq 0$  such that  $\delta = \frac{\eta\alpha}{2} \left(1 - \frac{\eta\alpha v}{2}\right) \geq 0$  and any policy  
 772  $\pi_t \in \Delta(\mathcal{A})$ , the following inequality holds

$$\delta(1 - 2\delta) \cdot \mathbb{E} [\text{IG}_t^{\mathcal{G}}(\pi_t)] \leq \mathbb{E} \left[ -\log \int_{\Theta} \left( \frac{p_t(\theta, A_t, L_t)}{p_t(\theta_0, A_t, L_t)} \right)^{\eta\alpha} dQ_t^+(\theta) \right]. \quad (41)$$

773 *Proof.* We remind the reader that  $\mathcal{F}_t = \theta(X_1, A_1, L_1, \dots, X_{t-1}, A_{t-1}, L_{t-1})$  is the  $\sigma$ -algebra  
 774 generated by the interaction history between the learner and the environment up to the end of round  $t$ .  
 775 By the tower rule of expectation, we have

$$\begin{aligned} & \mathbb{E} \left[ -\log \int_{\Theta} \left( \frac{p_t(\theta, A_t, L_t)}{p_t(\theta_0, A_t, L_t)} \right)^{\eta\alpha} dQ_t^+(\theta) \right] \\ &= \mathbb{E} \left[ \mathbb{E} \left[ -\log \int_{\Theta} \left( \frac{p_t(\theta, A_t, L_t)}{p_t(\theta_0, A_t, L_t)} \right)^{\eta\alpha} dQ_t^+(\theta) \middle| \mathcal{F}_{t-1}, X_t, A_t \right] \right] \\ &\leq \mathbb{E} \left[ -\log \mathbb{E} \left[ \int_{\Theta} \left( \frac{p_t(\theta, A_t, L_t)}{p_t(\theta_0, A_t, L_t)} \right)^{\eta\alpha} dQ_t^+(\theta) \middle| \mathcal{F}_{t-1}, X_t, A_t \right] \right] \\ &= \mathbb{E} \left[ -\log \int_{\Theta} \int_{\mathbb{R}} \left( \frac{p_t(\theta, A_t, L)}{p_t(\theta_0, A_t, L)} \right)^{\eta\alpha} d\mathcal{P}_{L_t|X_t, A_t}(L) dQ_t^+(\theta) \right]. \end{aligned} \quad (42)$$

776 Where the first inequality comes from Jensen's Inequality applied to the logarithm and  $\mathcal{P}_{L_t|X_t, A_t}$   
 777 is the conditional law of  $L_t$  given  $X_t$  and  $A_t$ . We fix  $\theta \in \Theta$ , drop the subscripts for simplicity and  
 778 define  $\ell = \ell_t(A_t), \ell_0 = \ell_t(\theta_0, A_t)$  and  $\mathcal{P}_t = \mathcal{P}_{L_t|X_t, A_t}$ . Using the definition of the likelihood  $p_t$ , we  
 779 get

$$\begin{aligned} & \int \left( \frac{p_t(\theta, A_t, L)}{p_t(\theta_0, A_t, L)} \right)^{\eta\alpha} d\mathcal{P}_t(L) \\ &= \int \exp \left( -\eta\alpha \left( \frac{(L - \ell_t(\theta, A_t))^2}{2} + \frac{(L - \ell(\theta_0, A_t))^2}{2} \right) \right) d\mathcal{P}_t(L) \\ &= \int \exp \left( \frac{\eta\alpha}{2} (2L - \ell - \ell_0) \cdot (\ell - \ell_0) \right) d\mathcal{P}_t(L) \\ &= \exp \left( -\frac{\eta\alpha}{2} (\ell + \ell_0) \cdot (\ell - \ell_0) \right) \cdot \int \exp(\eta\alpha L(\ell - \ell_0)) d\mathcal{P}_t(L) \\ &= \exp \left( \frac{\eta\alpha}{2} (\ell_0^2 - \ell^2) \right) \cdot \int \exp(\eta\alpha L(\ell - \ell_0)) d\mathcal{P}_t(L) \\ &\leq \exp \left( \frac{\eta\alpha}{2} (\ell_0^2 - \ell^2) \right) \cdot \exp(\eta\alpha \ell_0 \cdot (\ell - \ell_0)) \exp \left( \frac{\eta^2 \alpha^2 v}{2} (\ell - \ell_0)^2 \right) \\ &\leq \exp \left( -(\ell - \ell_0)^2 \cdot \frac{\eta\alpha}{2} \left( 1 - \frac{\eta\alpha v}{2} \right) \right). \end{aligned}$$

780 Now we define  $\delta = \frac{\eta\alpha}{2} \left( 1 - \frac{\eta\alpha v}{2} \right)$  we have :

$$\begin{aligned} & \int \left( \frac{p_t(\theta, A_t, L)}{p_t(\theta_0, A_t, L)} \right)^{\eta\alpha} d\mathcal{P}_t(L) \\ &\leq \exp \left( -(\ell - \ell_0)^2 \cdot \delta \right) \\ &\leq 1 - \delta(\ell - \ell_0)^2 + \frac{\delta^2}{2} (\ell - \ell_0)^4 \\ &\leq 1 - \delta(\ell - \ell_0)^2 + 4\delta^2 (\ell - \ell_0)^2 \\ &\leq 1 - \delta(1 - 2\delta)(\ell - \ell_0)^2. \end{aligned}$$

781 Where we use that  $|\ell - \ell_0| \leq 2$ . Finally using that for any  $x > 0, \log x \leq x - 1$  and equation 42, we  
 782 get the claim of the Lemma.  $\square$



783 It remains to pick the best values for  $\eta$ ,  $\alpha$  and  $\beta$  and apply Lemma 13 with  $Q^*$  a dirac distribution  
784 in  $\theta_0$  and Lemma 15. To finish the proof of Lemma 17, we combine the previous Lemma (17) with  
785 Lemma 15 and Lemma 13. We want the quantity  $\delta(1 - 2\delta)$  to be as big as possible, this happens  
786 when  $\delta = \frac{1}{4}$ . This is only possible if  $v \leq 1$  and  $\frac{\eta\alpha}{2} = \frac{1+\sqrt{1-v}}{2v}$ . If  $v > 1$ , our best choice of  $\frac{\eta\alpha}{2}$  is  
787  $\frac{1}{2v}$  and in that case  $\delta(1 - 2\delta) = \frac{1}{4v} \left(1 - \frac{1}{2v}\right) \geq \frac{1}{8v}$ . Finally, uniting both cases, we set  $\alpha = \beta = 2$ ,  
788  $\eta = \frac{1+\sqrt{1-v\wedge 1}}{2v}$  and we have that  $\delta(1 - 2\delta) \geq \frac{1}{8(1\vee v)}$ .

### 789 E.3.4 Metric Parameter Analysis : the proof of Lemma 8

790 We start by a technical lemma on the Lipschitzness of the losses and the optimal losses.

791 **Lemma 18.** *For any  $x, \theta, a, \ell_t(\cdot, x, a)$  and  $\ell_t^*(\cdot, x)$  are C-Lipschitz.*

792 *Proof.* Let  $\tau$  be the measure against which the densities  $p(\cdot|\theta, x, a)$  are defined. Without loss of  
793 generality, we can assume that  $\int_{[0,1]} d\tau(c) = 1$ . Letting  $\theta_1, \theta_2 \in \Theta$ , we have

$$\begin{aligned} |\ell(\theta_1, x, a) - \ell(\theta_2, x, a)| &= \left| \int_{[0,1]} c(p(c|\theta_1, x, a) - p(c|\theta_2, x, a)) d\tau(c) \right| \\ &\leq \int_{[0,1]} |(p(c|\theta_1, x, a) - p(c|\theta_2, x, a))| d\tau(c) \\ &= \int_{[0,1]} |\exp(\log(p(c|\theta_1, x, a))) - \exp(\log(p(c|\theta_2, x, a)))| d\tau(c) \\ &\leq \int_{[0,1]} C \|\theta_1 - \theta_2\| d\tau(c) \\ &= C \|\theta_1 - \theta_2\|, \end{aligned}$$

794 where the second inequality comes from the C-Lipschitzness of the composition of the exponential  
795 that is 1-Lipschitz on the negative numbers and the log likelihood that is C-Lipschitz. This proves  
796 the C-Lipschitzness of  $\ell_t(\cdot, x, a)$ . Now it easily follows that  $\ell^*(\cdot, x)$  is also C-Lipschitz, being an  
797 infimum of a family of C-Lipschitz functions.  $\square$

798 Now we introduce two further lemmas related to Lemma 13 when  $Q^*$  is chosen as a uniform  
799 distribution on a ball of radius  $\epsilon$ .

800 **Lemma 19.** *Fix  $\theta_0 \in \Theta$ , and  $\epsilon > 0$ , and assume that a ball including  $\theta_0$  with radius  $\epsilon$  is contained*  
801 *in  $\Theta$ . Letting  $Q^*$  be the uniform distribution on such a ball, we have*

$$802 \mathcal{D}_{KL}(Q^* \| Q_1) = d \log \left( \frac{R}{\epsilon} \right). \quad (43)$$

802 *Proof.* Since both  $Q^*$  and  $Q_1$  are uniform, the ratio of their density is equal to the ratio of the volume  
803 of  $\Theta$  and the volume of a ball of radius  $\epsilon$ . Since  $\Theta$  is included in a ball of radius  $R$ , this ratio is  
804 bounded by  $\left(\frac{R}{\epsilon}\right)^d$ . Finally

$$805 \mathcal{D}_{KL}(Q^* \| Q_1) = \int_{\Theta} \frac{dQ^*}{dQ_1}(\theta) \log \left( \frac{dQ^*}{dQ_1}(\theta) \right) dQ_1(\theta) \leq \log \left( \frac{R}{\epsilon} \right)^d \int_{\Theta} dQ^*(\theta) = d \log \left( \frac{R}{\epsilon} \right). \quad \square$$

806 **Lemma 20.** *Under the same conditions as Lemma 19, we have*

$$\left| \int \left( \frac{1}{\lambda\alpha} \cdot \sum_{t=1}^T \log \frac{p_t(\theta_0, A_t, L_t)^{\eta\alpha}}{p_t(\theta, A_t, L_t)^{\eta\alpha}} + \sum_{t=1}^T (\ell_t^*(\theta) - \ell_t^*(\theta_0)) \right) dQ^*(\theta) \right| \leq \left( \frac{\eta}{\lambda} + 1 \right) \cdot CT\epsilon. \quad (44)$$

807 *Proof.* This is a direct consequence of the Lipschitzness of the log-likelihood and Lemma 18.  $\square$

808 Putting Lemma 13 together with this choice of  $Q^*$  and with Lemma 14 and Lemma 15, we finish the  
809 proof of Lemma 8

810 **E.4 Upper bounds on the averaged DEC and the Surrogate Information ratio**

811 Here we provide the technical tools to bound the surrogate information ratio and the averaged DEC  
812 for some appropriately chosen forerunner algorithms.

813 **E.4.1 Worst-case analysis: The proof of Lemmas 1 and 12**

814 Here we study the performance of Thompson sampling as the forerunner algorithm, which will  
815 certify a bound on the surrogate information ratio of **OIDS**. The Thompson sampling policy  $\pi_t$  works  
816 by sampling  $\theta_t$  according to the posterior  $Q_t^+$  and then playing the action  $A_t \in \arg \min_a \ell_t(\theta_t, a)$ .  
817 To facilitate the derivations below, we define  $a_t^* : \Theta \rightarrow \mathcal{A}$  the greedy action selector by  $a_t^*(\theta) =$   
818  $\arg \min_a \ell_t(\theta, a)$  (with ties broken arbitrarily). By definition of the policy, sampling according to  
819  $\pi_t$  is the same as sampling according to  $dQ_t^+$  and then applying the greedy action selector. More  
820 formally, for any measurable function  $f$ , we have

$$\int_{\Theta} f(a_t^*(\theta)) dQ_t^+(\theta) = \sum_a \pi_t(a) f(a).$$

821 Moreover, we have that  $\bar{\ell}_t^* = \int_{\Theta} \ell_t^*(\theta) dQ_t^+(\theta) = \int_{\Theta} \ell_t(\theta, a_t^*(\theta)) dQ_t^+(\theta)$ . Putting these observations  
822 together, we can write the surrogate regret as

$$\bar{r}_t(\pi_t) = \sum_a \pi_t(a) (\bar{\ell}_t(a) - \bar{\ell}_t^*) = \int_{\Theta} \bar{\ell}_t(a_t^*(\theta)) - \ell_t(\theta, a_t^*(\theta)). \quad (45)$$

823 Observe that the regret is the difference of the expectation of the same function under the joint  
824 distribution of  $\theta_t$  and  $A_t$  and their product distribution, and thus measures the extent to which the  
825 two are ‘‘coupled’’. We will analyze this quantity by a decoupling argument inspired by [Zhang \[2022\]](#)  
826 and [Neu et al. \[2022\]](#).

827 For setting up the decoupling analysis, we first need some technical lemmas. We start by a corollary  
828 of the Fenchel–Young inequality for strongly convex functions that will come handy.

829 **Lemma 21.** *Let  $I$  be an interval on the real line and let  $\mathcal{D} : I^2 \rightarrow \mathbb{R}$  be a convex function satisfying  
830 the following conditions:*

- 831 • For any  $y \in I$ , the function  $x \rightarrow \mathcal{D}(x, y)$  is proper, closed and  $C$ -strongly convex.
- 832 • For any  $x \in I$ ,  $\mathcal{D}(x, x) = 0$ .

833 Then for any  $x, y \in I$  and any  $\mu \in \mathbb{R}$  we have

$$(x - y)u \leq \mathcal{D}(x, y) + \frac{u^2}{2C}. \quad (46)$$

834 *Proof.* Let  $y \in I$ . We compute the Legendre–Fenchel conjugate of  $x \rightarrow \mathcal{D}(x, y)$ , defined for any  
835  $u \in \mathbb{R}$  as

$$\mathcal{D}^*(u, y) = \sup_{x \in I} \{xu - \mathcal{D}(x, y)\}.$$

836 Since  $y$  is a minimum of  $x \rightarrow \mathcal{D}(x, y)$  and  $\mathcal{D}(y, y) = 0$ , we have that  $\mathcal{D}^*(0, y) = 0$ . Moreover using  
837 Lemma 15 of [Shalev-Shwartz \[2007\]](#), we directly have that  $\mathcal{D}^*$  is  $\frac{1}{C}$  smooth in its first coordinate and  
838 that  $\frac{\partial \mathcal{D}^*}{\partial u}(0, y) = y$ , so that for any  $u \in \mathbb{R}$  we have

$$\mathcal{D}^*(u, y) \leq \mathcal{D}^*(0, y) + u \frac{\partial \mathcal{D}^*}{\partial u}(0, y) + \frac{u^2}{2C} \leq yu + \frac{u^2}{2C}.$$

839 Then, by the Fenchel–Young inequality, this implies the following for any  $x \in I$  and any  $u \in \mathbb{R}$ :

$$x \cdot \mu \leq \mathcal{D}(x, y) + \mathcal{D}^*(u, y) \leq y \cdot u + \frac{u^2}{2C}.$$

840 This proves the statement. □

841 We use this inequality to prove the following general decoupling lemma that can handle arbitrary  
842 joint distributions of random variables.

843 **Lemma 22.** Let  $\mathcal{D} : [0, 1]^2 \rightarrow \mathbb{R}$  be  $C$ -strongly convex and satisfy the same hypothesis as for the  
844 previous lemma. Let  $Q \in \Delta(\Theta)$ ,  $f : \Theta \times \mathcal{A} \rightarrow [0, 1]$  and  $a^* : \Theta \rightarrow \mathcal{A}$ . Assume  $f$  and  $a^*$  are  
845 measurable. We define  $\pi \in \Delta(\mathcal{A})$  by  $\pi(a) = \int_{\Theta} \mathbb{I}_{\{a^*(\theta)=a\}} dQ(\theta)$  and  $\bar{f}(a) = \int_{\Theta} f(\theta, a) dQ(\theta)$ .  
846 Then for any  $\mu > 0$  the following holds

$$\int_{\Theta} \bar{f}(a^*(\theta)) - f(\theta, a^*(\theta)) dQ(\theta) \leq \mu \int_{\Theta} \sum_a \pi(a) \mathcal{D}(\bar{f}(a), f(\theta, a)) dQ(\theta) + \frac{K}{2\mu C} \quad (47)$$

847 *Proof.* We start by writing

$$\begin{aligned} \int_{\Theta} \bar{f}(a^*(\theta)) - f(\theta, a^*(\theta)) &= \int_{\Theta} \sum_a \frac{\mu\pi(a)}{\mu\pi(a)} \mathbb{I}_{\{a^*(\theta)=a\}} (\bar{f}(a) - f(\theta, a)) dQ(\theta) \\ &= \int_{\Theta} \sum_a \mu\pi(a) \left( \frac{\mathbb{I}_{\{a^*(\theta)=a\}}}{\mu\pi(a)} (\bar{f}(a) - f(\theta, a)) \right) dQ(\theta) \\ &\leq \int_{\Theta} \sum_a \mu\pi(a) \left( \mathcal{D}(\bar{f}(a), f(\theta, a)) + \frac{\mathbb{I}_{\{a^*(\theta)=a\}}}{2C\mu^2\pi(a)^2} \right) dQ(\theta), \end{aligned}$$

848 where we used Lemma 21 with  $u = \frac{\mathbb{I}_{\{a^*(\theta)=a\}}}{\mu\pi(a)}$  in the last line. Finally, we have

$$\begin{aligned} \int_{\Theta} \bar{f}(a^*(\theta)) - f(\theta, a^*(\theta)) &\leq \mu \int_{\Theta} \sum_a \pi(a) \mathcal{D}(\bar{f}(a), f(\theta, a)) dQ(\theta) + \frac{1}{2\mu C} \sum_a \int_{\Theta} \frac{\mathbb{I}_{\{a^*(\theta)=a\}}}{\pi(a)} dQ(\theta) \\ &\leq \mu \int_{\Theta} \sum_a \pi(a) \mathcal{D}(\bar{f}(a), f(\theta, a)) dQ(\theta) + \frac{K}{2\mu C}, \end{aligned}$$

849 where we used  $\pi(a) = \int_{\Theta} \mathbb{I}_{\{a^*(\theta)=a\}} dQ(\theta)$  in the last line. □

850 To prove Lemma 1, we use the above result with  $Q = Q_t^+$ ,  $f = \ell_t$  and  $a^* = a_{j_t}$  and  $\mathcal{D}$  chosen as  
851 the squared Hellinger distance  $\mathcal{D}_H^2$ , which is  $\frac{1}{4}$ -strongly convex in its first argument by Lemma 24  
852 provided in Appendix E.5. Thus, applying Lemma 22 we get for any  $\mu > 0$  that

$$\bar{r}_t(\pi_t) \leq \mu \int_{\Theta} \sum_a \pi_t(a) \mathcal{D}_H^2(\bar{\ell}_t(a), \ell_t(\theta, a)) dQ_t^+(\theta) + \frac{2K}{\mu}.$$

853 This concludes the proof of Lemma 1.

854 Lemma 12 is proved by choosing  $\mathcal{D}(x, y) = (x - y)^2$  that is 2-strongly convex in its first argument,  
855 which yields the advertised result as

$$\bar{r}_t(\pi_t) \leq \mu \int_{\Theta} \sum_a \pi_t(a) (\bar{\ell}_t(a) - \ell_t(\theta, a))^2 dQ_t^+(\theta) + \frac{K}{4\mu}.$$

#### 856 E.4.2 Instance-dependent analysis: The proof of Lemma 5

857 This analysis uses the so-called *inverse-gap weighting* algorithm of Abe and Long [1999] as  
858 forerunner—see also the works of Foster and Rakhlin [2020] and Foster and Krishnamurthy [2021]  
859 that reignited interest in this method. Our analysis below is especially inspired by the latter work.

860 We define the inverse gap weighting policy with scale parameter  $\gamma$  and with respect to a nominal loss  
861 function  $f : \mathcal{A} \rightarrow \mathbb{R}^+$  as

$$\pi_{\gamma, f}^{(\text{IGW})}(a) = \begin{cases} \frac{f(b)}{Kf(b) + \gamma(f(b) - f(a))} & \text{if } a \neq b \\ 1 - \sum_{a \neq b} \pi_{\gamma, f}^{(\text{IGW})}(a) & \text{if } a = b \end{cases}$$

862 where  $b \in \arg \min_a f(a)$  is fixed (with ties broken arbitrarily). We fix  $\theta$  and apply Lemma 4 of  
863 Foster and Krishnamurthy [2021] with nominal loss  $\bar{\ell}_t : \mathcal{A} \rightarrow \mathbb{R}$  and true loss  $\ell_t(\theta) : \mathcal{A} \rightarrow \mathbb{R}$  to get

$$\bar{\ell}_t(b) - \ell_t(\theta, a_t^*(\theta)) \leq \frac{K}{4\gamma} \bar{\ell}_t(b) + 2\gamma \cdot \pi_{\gamma, \bar{\ell}_t}^{(\text{IGW})}(a_t^*(\theta)) \frac{(\bar{\ell}_t(a_t^*(\theta)) - \ell_t(\theta, a_t^*(\theta)))^2}{\bar{\ell}_t(a_t^*(\theta)) + \ell_t(\theta, a_t^*(\theta))}$$

$$\leq \frac{K}{4\gamma} \bar{\ell}_t(b) + 2\gamma \cdot \sum_a \pi_{\gamma, \bar{\ell}_t}^{(\text{IGW})}(a) \frac{(\bar{\ell}_t(a) - \ell_t(\theta, a))^2}{\bar{\ell}_t(a) + \ell_t(\theta, a)},$$

864 where  $b \in \arg \min_a \bar{\ell}_t(a)$  and  $a_t^*(\theta) \in \arg \min_a \ell_t(\theta, a)$ . To proceed, we use that for any  $p, q \in$   
 865  $[0, 1]$ , we have  $\frac{(p-q)^2}{p+q} \leq 4 \cdot \mathcal{D}_H^2(\text{Ber}(p), \text{Ber}(q))$  (cf. Lemma 23 in Appendix E.5). We combine this  
 866 with the data processing inequality for  $f$ -divergences to obtain

$$\begin{aligned} \bar{\ell}_t(b) - \ell_t(\theta, a_t^*(\theta)) &\leq \frac{K}{4\gamma} \bar{\ell}_t(b) + 8\gamma \cdot \sum_a \pi_{\gamma, \bar{\ell}_t}^{(\text{IGW})}(a) \mathcal{D}_H^2(\text{Ber}(\bar{\ell}_t(a)), \text{Ber}(\ell_t(\theta, a))) \\ &\leq \frac{K}{4\gamma} \bar{\ell}_t(b) + 8\gamma \cdot \sum_a \pi_{\gamma, \bar{\ell}_t}^{(\text{IGW})}(a) \mathcal{D}_H^2(\bar{p}_t(a, \cdot), p_t(\theta, a, \cdot)). \end{aligned} \quad (48)$$

867 On the other hand, we can rewrite the surrogate regret of the inverse gap weighting policy as

$$\begin{aligned} \bar{r}_t(\pi_{\gamma, \bar{\ell}_t}^{(\text{IGW})}) &= \int \sum_a \pi_{\gamma, \bar{\ell}_t}^{(\text{IGW})}(a) (\bar{\ell}_t(a) - \ell_t^*(\theta)) dQ_t^+(\theta) \\ &= \int \sum_a \pi_{\gamma, \bar{\ell}_t}^{(\text{IGW})}(a) (\bar{\ell}_t(a) - \ell_t(\theta, a_t^*(\theta))) dQ_t^+(\theta) \\ &= \int \sum_{a \neq b} \pi_{\gamma, \bar{\ell}_t}^{(\text{IGW})}(a) (\bar{\ell}_t(a) - \bar{\ell}_t(b)) dQ_t^+(\theta) + \int \sum_a \pi(a) (\bar{\ell}_t(b) - \ell_t(\theta, a_t^*(\theta))) dQ_t^+(\theta). \end{aligned}$$

868 The second term in the above decomposition can be bounded using Equation (48). As for the first  
 869 term, we can exploit the definition of the policy to write

$$\sum_{a \neq b} \pi_{\gamma, \bar{\ell}_t}^{(\text{IGW})}(a) (\bar{\ell}_t(a) - \bar{\ell}_t(b)) = \sum_{a \neq b} \frac{\bar{\ell}_t(b) (\bar{\ell}_t(a) - \bar{\ell}_t(b))}{K \bar{\ell}_t(b) + \gamma (\bar{\ell}_t(a) - \bar{\ell}_t(b))} \leq \frac{K \bar{\ell}_t(b)}{\gamma}.$$

870 Putting these bounds together gives

$$\begin{aligned} \bar{r}_t(\pi_{\gamma, \bar{\ell}_t}^{(\text{IGW})}) &\leq \frac{K \bar{\ell}_t(b)}{\gamma} + \frac{K \bar{\ell}_t(b)}{4\gamma} + 8\gamma \cdot \int \sum_a \mathcal{D}_H^2(\bar{p}_t(a, \cdot), p_t(\theta, a, \cdot)) dQ_t^+(\theta) \\ &\leq \frac{5K \bar{\ell}_t(b)}{4\gamma} + 8\gamma \cdot \bar{\text{IG}}_t, \end{aligned}$$

871 Optimizing over  $\gamma$ , we get the claim of Lemma 5.

## 872 E.5 Auxiliary results

873 **Lemma 23** (Proposition 3 Foster and Krishnamurthy, 2021). For any  $p, q \in [0, 1]$ , we have

$$\frac{(p-q)^2}{p+q} \leq 4 \mathcal{D}_H^2(\text{Ber}(p), \text{Ber}(q)).$$

874 *Proof.* The statement follows from the simple calculation

$$\mathcal{D}_H^2(p, q) \geq \frac{1}{2} (\sqrt{p} - \sqrt{q})^2 = \frac{1}{2} \left( \frac{(\sqrt{p} - \sqrt{q})(\sqrt{p} + \sqrt{q})}{\sqrt{p} + \sqrt{q}} \right)^2 = \frac{1}{2} \frac{(p-q)^2}{(\sqrt{p} + \sqrt{q})^2} \geq \frac{1}{4} \frac{(p-q)^2}{p+q},$$

875 where the last step uses the elementary inequality  $(x+y)^2 \leq 2(x^2+y^2)$  that holds for any  $x, y$ .  $\square$

876 **Lemma 24.** For any fixed  $q \in [0, 1]$ , the function  $p \mapsto \mathcal{D}_H^2(\text{Ber}(p), \text{Ber}(q))$  is  $\frac{1}{4}$ -strongly convex.

877 *Proof.* The proof is based on showing that the second derivative of the function of interest is uniformly  
 878 lower bounded by a positive constant. This follows from calculating the first derivative as

$$\frac{\partial \mathcal{D}_H^2(p, q)}{\partial p} = \frac{1}{2} \left( -\sqrt{\frac{q}{p}} + \sqrt{\frac{1-q}{1-p}} \right),$$

879 and then lower-bounding the second derivative as

$$\frac{\partial^2 \mathcal{D}_H^2(p, q)}{\partial^2 p} = \frac{1}{4} \left( \sqrt{\frac{q}{p^3}} + \sqrt{\frac{1-q}{(1-p)^3}} \right) \geq \frac{1}{4} (\sqrt{q} + \sqrt{1-q}) \geq \frac{1}{4}.$$

880 This inequality is tight when  $q = 0$  or  $q = 1$ .

□