
On the Complexity of Adversarial Decision Making

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A central problem in online learning and decision making—from bandits to rein-
2 forcement learning—is to understand what modeling assumptions lead to sample-
3 efficient learning guarantees. With a focus on stochastic environments, a recent
4 line of research provides general structural conditions under which sample-efficient
5 learning is possible, but robust learning guarantees for agnostic or adversarial
6 settings have remained elusive. We consider a general *adversarial decision making*
7 framework that encompasses (structured) bandit problems with adversarial rewards
8 and reinforcement learning problems with adversarial dynamics. Our main result
9 is to show—via new upper and lower bounds—that the Decision-Estimation Co-
10 efficient, a complexity measure introduced by Foster et al. [18] in the stochastic
11 counterpart to our setting, is both necessary and sufficient for low regret in the
12 adversarial setting. However, compared to the stochastic setting, one must apply
13 the Decision-Estimation Coefficient to the *convex hull* of the class of models (or,
14 hypotheses) under consideration. This establishes that the price of accommodating
15 adversarial rewards or dynamics is governed by the behavior of the model class
16 under convexification, and recovers a number of existing results—both positive and
17 negative. En route to obtaining these guarantees, we provide new structural results
18 that connect the Decision-Estimation Coefficient to variants of other well-known
19 complexity measures, including the Information Ratio of Russo and Van Roy [52]
20 and the Exploration-by-Optimization objective of Lattimore and György [34].

21 1 Introduction

22 We consider the problem of robust data-driven decision making in bandits, reinforcement learn-
23 ing, and beyond. The last decade has seen development of data-driven decision algorithms with
24 strong empirical performance in domains including robotics [28, 40], dialogue systems [38], and
25 personalization [2, 57]. Reliably deploying data-driven decision making methods in safety-critical
26 systems requires principled algorithms with provable robustness in the face of dynamic or even
27 adversarial environments. Furthermore, for such algorithms to be applicable, they must effectively
28 take advantage of problem structure as modeled by the practitioner. In high-dimensional problems,
29 this means efficiently generalizing across states and actions while delicately exploring new decisions.

30 For decision making in static, stochastic environments, recent years have seen extensive investigation
31 into optimal sample complexity and algorithm design principles, and the foundations are beginning to
32 take shape. In particular, with an emphasis on reinforcement learning, a burgeoning body of research
33 identifies specific modeling assumptions under which sample-efficient interactive decision making
34 is possible [13, 60, 22, 44, 6, 29, 15, 39, 14, 63], as well as general structural conditions that aim
35 to unify these settings [50, 21, 56, 59, 16, 23, 18]. For dynamic or adversarial settings, however,
36 comparatively little is known outside of (i) positive results for special cases such as adversarial bandit
37 problems [5, 4, 20, 11, 1, 8, 27, 17, 9, 31], and (ii) a handful of negative results suggesting that online
38 reinforcement learning in agnostic or adversarial settings can actually be statistically intractable

39 [53, 41]. These developments raise the following questions: (a) what are the underlying phenomena
 40 that determine the statistical complexity of decision making in adversarial settings? (b) what are the
 41 corresponding algorithmic design principles that attain optimal sample complexity?

42 **Contributions.** We consider an adversarial variant of the *Decision Making with Structured*
 43 *Observations* (DMSO) framework introduced in Foster et al. [18], where learner or decision-maker
 44 interacts with a sequence of *models* (reward distributions in the case of bandits, or MDPs in the case
 45 of reinforcement learning) chosen by an adaptive adversary, and aims to minimize regret against the
 46 best decision in hindsight. The models are assumed to belong to a known model class which reflects
 47 the learner’s prior knowledge about the problem. The main question we investigate is: *How does the*
 48 *structure of the model class determine the minimax regret for adversarial decision making?* We show:

- 49 1. For *any* model class, one can obtain high-probability regret bounds based on a *convexified* version
 50 of the *Decision-Estimation Coefficient* (DEC) complexity measure introduced in Foster et al. [18].
- 51 2. For any algorithm with reasonable tail behavior, the optimal regret for adversarial decision making
 52 is lower bounded by (a suitably localized version of) the convexified DEC.

53 In the process, we draw new connections to several existing complexity measures.

54 1.1 Problem Setting

55 We adopt an adversarial variant of the DMSO framework of Foster et al. [18]. The protocol consists
 56 of T rounds, where at each round $t = 1, \dots, T$:

- 57 1. The learner selects a *decision* $\pi^{(t)} \in \Pi$, where Π is the *decision space*.
- 58 2. Nature selects a *model* $M^{(t)} \in \mathcal{M}$, where \mathcal{M} is a *model class*.
- 59 3. The learner receives a reward $r^{(t)} \in \mathcal{R} \subseteq \mathbb{R}$ and observation $o^{(t)} \in \mathcal{O}$ sampled via $(r^{(t)}, o^{(t)}) \sim$
 60 $M^{(t)}(\pi^{(t)})$, where \mathcal{O} is the *observation space*. We abbreviate $z^{(t)} := (r^{(t)}, o^{(t)})$ and $\mathcal{Z} := \mathcal{R} \times \mathcal{O}$.

61 Here, each model $M = M(\cdot, \cdot | \cdot) \in \mathcal{M}$ is a conditional distribution $M : \Pi \rightarrow \Delta(\mathcal{R} \times \mathcal{O})$ that
 62 maps the learner’s decision to a distribution over rewards and outcomes. This setting subsumes
 63 (adversarial) bandit problems, where models consist of reward functions/distributions, as well as
 64 adversarial reinforcement learning, where models correspond to Markov decision processes (MDPs).
 65 In both cases, the model class \mathcal{M} encodes prior knowledge about the decision making problem such
 66 as structure of rewards or dynamics (e.g., linearity or convexity), and might be parameterized by linear
 67 models, neural networks, or other rich function approximators depending on the problem domain.

68 For a model $M \in \mathcal{M}$, $\mathbb{E}^{M, \pi}[\cdot]$ denotes expectation under the process $(r, o) \sim M(\pi)$. We define
 69 $f^M(\pi) := \mathbb{E}^{M, \pi}[r]$ as the mean reward function and $\pi_M := \arg \max_{\pi \in \Pi} f^M(\pi)$ as the decision with
 70 greatest reward for M . We let $\mathcal{F}_{\mathcal{M}} = \{f^M \mid M \in \mathcal{M}\}$ denote the induced class of reward functions.
 71 We measure performance via *regret* to the best fixed decision in hindsight:

$$\text{Reg}_{\text{DM}} := \sup_{\pi^* \in \Pi} \sum_{t=1}^T \mathbb{E}_{\pi^{(t)} \sim p^{(t)}} \left[f^{M^{(t)}}(\pi^*) - f^{M^{(t)}}(\pi^{(t)}) \right]. \quad (1)$$

72 This formulation generalizes Foster et al. [18], who considered the *stochastic* setting where $M^{(t)} =$
 73 M^* is fixed across all rounds. Examples include:

- 74 • **Adversarial bandits.** With no observations ($\mathcal{O} = \{\emptyset\}$), the adversarial DMSO framework is
 75 equivalent to the *adversarial bandit* problem with structured rewards. In this context, $\pi^{(t)}$ is
 76 typically referred to as an *action* or *arm* and Π is referred to as the *action space*. The most basic
 77 example here is the adversarial finite-armed bandit problem with A actions [5, 4, 20], where
 78 $\Pi = \{1, \dots, A\}$ and $\mathcal{F}_{\mathcal{M}} = \mathbb{R}^A$. Other well-studied examples include adversarial linear bandits
 79 [11, 1, 8], bandit convex optimization [27, 17, 9, 31], and nonparametric bandits [27, 7, 43].¹
- 80 • **Reinforcement learning.** The adversarial DMSO framework encompasses finite-horizon, episodic
 81 online reinforcement learning, with each round t corresponding to a single episode: $\pi^{(t)}$ is a
 82 *policy* (a mapping from state to actions) to play in the episode, $r^{(t)}$ is the cumulative reward in the

¹Typically, these examples are formulated with deterministic rewards, which we encompass by restricting
 models in \mathcal{M} to be deterministic. Our formulation is more general and allows for, semi-stochastic adversaries.

83 episode, and the observation $o^{(t)}$ is the episode’s trajectory (sequence of observed states, actions,
 84 and rewards). Online reinforcement learning in the stochastic setting where $M^{(t)} = M^*$ is fixed
 85 has received extensive attention [21, 56, 22, 59, 16, 23, 18], but the adversarial setting we study has
 86 received less investigation. Examples include the setting in which the adversary chooses a sequence
 87 of tabular MDPs, which is known to be intractable [41], and the easier setting in which there is a
 88 fixed (known) MDP but rewards are adversarial [45, 64, 46, 24]. See Appendix D for more details.

89 We refer to Appendix B for additional measure-theoretic details and background, and to Foster et al.
 90 [18] for further examples and detailed discussion.²

91 Understanding sample complexity for the DMSO setting at this level of generality is a challenging
 92 problem. Even if one restricts only to bandit-type problems (with no observations), any complexity
 93 measure must capture the role of structural assumptions such as convexity or smoothness in determin-
 94 ing the optimal rates. To go beyond bandit problems and handle the general setting, one must accom-
 95 modate problems with rich, structured feedback such as reinforcement learning, where observations
 96 (as well as subtle features of the noise distribution) can reveal information about the underlying model.

97 1.2 Overview of Results

98 For a model class \mathcal{M} , reference model $\bar{M} \in \mathcal{M}$, and scale parameter $\gamma > 0$, the Decision-Estimation
 99 Coefficient [18] is defined via

$$\text{dec}_\gamma(\mathcal{M}, \bar{M}) = \inf_{p \in \Delta(\Pi)} \sup_{M \in \mathcal{M}} \mathbb{E}_{\pi \sim p} [f^M(\pi_M) - f^M(\pi) - \gamma \cdot D_{\mathbb{H}}^2(M(\pi), \bar{M}(\pi))], \quad (2)$$

100 where we recall that for probability measures \mathbb{P} and \mathbb{Q} with a common dominating measure ν ,
 101 (squared) Hellinger distance is given by $D_{\mathbb{H}}^2(\mathbb{P}, \mathbb{Q}) = \int (\sqrt{d\mathbb{P}/d\nu} - \sqrt{d\mathbb{Q}/d\nu})^2$. We define
 102 $\text{dec}_\gamma(\mathcal{M}) = \sup_{\bar{M} \in \mathcal{M}} \text{dec}_\gamma(\mathcal{M}, \bar{M})$, and let $\text{co}(\mathcal{M})$ denote the convex hull of \mathcal{M} , which can
 103 equivalently be viewed as the set of all mixtures of models in \mathcal{M} . Our main results show that the
 104 *convexified Decision-Estimation Coefficient* $\text{dec}_\gamma(\text{co}(\mathcal{M}))$ leads to upper and lower bounds on the
 105 optimal regret for adversarial decision making.

106 **Theorem (informal).** *For any model class \mathcal{M} , Algorithm 1 ensures that with high probability,*

$$\mathbf{Reg}_{\text{DM}} \lesssim \text{dec}_\gamma(\text{co}(\mathcal{M})) \cdot T, \quad (3)$$

107 where γ satisfies the balance $\text{dec}_\gamma(\text{co}(\mathcal{M})) \propto \frac{\gamma}{T} \log |\Pi|$. Moreover, for any algorithm with “rea-
 108 sonable” tail behavior (Section 2.2), regret must scale with a localized version of the same quantity.

109 As a consequence, there exists an algorithm for which $\mathbb{E}[\mathbf{Reg}_{\text{DM}}] \leq \tilde{o}(T)$ if and only if
 110 $\text{dec}_\gamma(\text{co}(\mathcal{M})) \propto \gamma^{-\rho}$ for some $\rho > 0$.

111 For the stochastic version of our setting, Foster et al. [18] give upper and lower bounds that scale with
 112 $\text{dec}_\gamma(\mathcal{M})$ (under appropriate technical assumptions; cf. Section 2.3). Hence, our results show that in
 113 general, the gap between optimal regret for stochastic and adversarial settings (or, “price of adversarial
 114 outcomes”) is governed by the behavior of the DEC under convexification. For example, multi-armed
 115 bandits, linear bandits, and convex bandits are convex model classes (where $\text{co}(\mathcal{M}) = \mathcal{M}$), which
 116 gives a post-hoc explanation for why these models are tractable in the adversarial setting. Finite
 117 state/action Markov decision processes are not a convex class, and have $\text{dec}_\gamma(\text{co}(\mathcal{M}))$ exponentially
 118 large compared to $\text{dec}_\gamma(\mathcal{M})$; in this case, our results recover lower bounds of Liu et al. [41].

119 Beyond these results, we prove that the convexified Decision-Estimation Coefficient is equivalent to:

- 120 1. a “parameterized” variant of the generalized information ratio of Lattimore and György [34].
- 121 2. a novel high-probability variant of the *Exploration-by-Optimization* of Lattimore and György [34].

122 Overall, while our results heavily draw on the work of Foster et al. [18] and Lattimore and György [34],
 123 we believe they play a valuable role in bridging these lines of research and formalizing connections.

124 **Our techniques.** On the lower bound side, we strengthen the lower bound from Foster et al. [18]
 125 with an improved change-of-measure argument (leading to improved results even in the stochastic

²We mention in passing that the upper bounds in this paper encompass the more general setting where rewards are not observed by the learner (i.e., $z^{(t)}$ does not contain the reward), thus subsuming the partial monitoring problem. Our lower bounds, however, require that rewards are observed. See Appendix A.

126 setting), and combine this with the simple idea of choosing a static mixture model as the adversary.
127 On the upper bound side, we extend the powerful Exploration-by-Optimization machinery of
128 Lattimore and György [34] to the DMSO setting, and give a novel high-probability variant of the
129 technique. We show that the performance of this method is controlled by a complexity measure
130 whose value is equivalent to the convexified DEC, as well as parameterized variant of the information
131 ratio (we present results in terms of the former to draw comparison to the stochastic setting).

132 **Organization.** Section 2 presents our main results, including upper and lower bounds on regret and
133 a characterization of learnability. In Section 3, we provide new structural results connecting the DEC
134 to Exploration-by-Optimization and the information ratio. We close with future directions (Section 4).
135 Additional comparison to related work is deferred to Appendix A. The appendix also contains proofs
136 and additional results, including examples (Appendix D) and further structural results (Appendix E).

137 2 Main Results

138 We now present our main results. First, using a new high-probability variant of the Exploration-by-
139 Optimization technique [37, 34], we provide an upper bound on regret via the (convexified) Decision-
140 Estimation Coefficient (Section 2.1). Next, we present a lower bound that scales with a localized
141 version of the same quantity (Section 2.2), and use these results to give a characterization for learn-
142 ability (Section 2.3). Finally, we discuss the gap between stochastic and adversarial decision making.

143 For the sake of keeping presentation as simple as possible, we make the following assumption.

144 **Assumption 2.1.** *The decision space Π has $|\Pi| < \infty$, and we have $\mathcal{R} = [0, 1]$.*

145 This assumption only serves to facilitate the use of the minimax theorem, and we expect that our results
146 can be generalized substantially (e.g., with covering numbers as in Section 3.4 of Foster et al. [18]).

147 2.1 Upper Bound

148 In this section we give regret bounds for adversarial decision making based on the (convexified)
149 Decision-Estimation Coefficient. A-priori, it is not obvious why the DEC should bear any relevance
150 to the adversarial setting: The algorithms and regret bounds based on the DEC that Foster et al. [18]
151 introduce for the stochastic setting heavily rely on the ability to estimate a static underlying model,
152 yet in the adversarial setting the learner may only interact with each model a single time. This renders
153 any sort of global estimation (e.g., for dynamics of an MDP) impossible. In spite of this difficulty, we
154 show that regret bounds can be achieved by building on the powerful *Exploration-by-Optimization*
155 technique of Lattimore and Szepesvári [37], Lattimore and György [34], which provides an elegant
156 approach to estimating rewards while exploiting the structure of the model class under consideration.

157 Exploration-by-Optimization—which was introduced in Lattimore and Szepesvári [37] and substan-
158 tially expanded in Lattimore and György [34]—can be thought of as a generalization of the classical
159 EXP3 algorithm [5], which we recall applies the exponential weights method for full-information
160 online learning to a sequence of unbiased estimators for the rewards (formed via importance weight-
161 ing). The naive reward estimator used by EXP3 is unsuitable for general model classes because it
162 does not exploit the structure of the decision space. Consequently, the regret scales linearly with $|\Pi|$
163 rather than with, e.g., dimension, as one might hope for linear bandits. The idea behind Exploration-
164 by-Optimization is to solve an optimization problem at each round to find a reward estimator and
165 modified sampling distribution that better exploit the structure of the model class \mathcal{M} for improved re-
166 gret. Lattimore and György [34] showed that for a general partial monitoring setting (cf. Appendix A),
167 the expected regret of this method—for exponential weights and a more general family of algorithms
168 based on Bregman divergences—is bounded by a generalization of the information ratio [51, 52].

169 Our development builds on that of Lattimore and György [34], but we pursue *high-probability* guaran-
170 tees rather than in-expectation guarantees.³ While high-probability guarantees are useful in their own
171 right, our motivation for studying such guarantees comes from the lower bound in the sequel (Sec-
172 tion 2.2), which shows that the convexified Decision-Estimation Coefficient lower bounds the regret
173 for algorithms with “reasonable” tail behavior. To develop high-probability regret bounds and com-
174 plement this lower bound, we use a novel high-probability variant of the Exploration-by-Optimization
175 objective and a specialized analysis which goes beyond the Bregman divergence framework.

³In general, in-expectation regret bounds do not imply high-probability bounds. For example, in adversarial bandits, the EXP3 algorithm can experience linear regret with constant probability [36].

Algorithm 1 High-Probability Exploration-by-Optimization (ExO⁺)

- 1: **parameters:** Learning rate $\eta > 0$.
- 2: **for** $t = 1, 2, \dots, T$ **do**
- 3: Define $q^{(t)} \in \Delta(\Pi)$ via exponential weights update:

$$q^{(t)}(\pi) = \exp\left(\eta \sum_{i=1}^{t-1} \widehat{f}^{(i)}(\pi)\right) / \sum_{\pi' \in \Pi} \exp\left(\eta \sum_{i=1}^{t-1} \widehat{f}^{(i)}(\pi')\right). \quad (4)$$

- 4: Solve *high-probability exploration-by-optimization* objective: // See Eq. (7)

$$(p^{(t)}, g^{(t)}) \leftarrow \arg \min_{p \in \Delta(\Pi), g \in \mathcal{G}} \sup_{M \in \mathcal{M}, \pi^* \in \Pi} \Gamma_{q^{(t)}, \eta}(p, g; \pi^*, M). \quad (5)$$

- 5: Sample decision $\pi^{(t)} \sim p^{(t)}$ and observe $z^{(t)} = (r^{(t)}, o^{(t)})$.
- 6: Form reward estimator:

$$\widehat{f}^{(t)}(\pi) = \frac{g^{(t)}(\pi; \pi^{(t)}, z^{(t)})}{p^{(t)}(\pi^{(t)})}. \quad (6)$$

Our algorithm, ExO⁺, is displayed in [Algorithm 1](#). At each round t , the algorithm proceeds by forming a *reference distribution* $q^{(t)} \in \Delta(\Pi)$ by applying the standard exponential weights update (with learning rate $\eta > 0$) to a sequence of reward estimators $\widehat{f}^{(1)}, \dots, \widehat{f}^{(t-1)}$ from previous rounds ([Line 3](#)). Next, for the main step of the algorithm ([Line 4](#)), we obtain a sampling distribution $p^{(t)} \in \Delta(\Pi)$ and an *estimation function* $g^{(t)} \in \mathcal{G} := (\Pi \times \Pi \times \mathcal{X} \rightarrow \mathbb{R})$ by solving a minimax optimization problem based on a new objective we term *high-probability exploration-by-optimization*:

$$\Gamma_{q, \eta}(p, g; \pi^*, M) := \mathbb{E}_{\pi \sim p}[f^M(\pi^*) - f^M(\pi)] + \eta^{-1} \cdot \mathbb{E}_{\pi \sim p, z \sim M(\pi)} \mathbb{E}_{\pi' \sim q} \left[\exp\left(\frac{\eta}{p(\pi)}(g(\pi'; \pi, z) - g(\pi^*; \pi, z))\right) - 1 \right]. \quad (7)$$

Finally ([Lines 5 and 6](#)), the algorithm samples $\pi^{(t)} \sim p^{(t)}$, observes $z^{(t)} = (r^{(t)}, o^{(t)})$, and then forms an importance-weighted reward estimator via $\widehat{f}^{(t)}(\pi) := g^{(t)}(\pi; \pi^{(t)}, z^{(t)})/p^{(t)}(\pi^{(t)})$.

The interpretation of the high-probability Exploration-by-Optimization objective (7) is as follows: For a given round t , the model $M \in \mathcal{M}$ and decision $\pi^* \in \Pi$ should be thought of as a proxy for the true model and optimal decision, respectively. By solving the minimax problem in (5), the min-player aims to—in the face of an unknown/worst-case model—find a sampling distribution that minimizes instantaneous regret, yet ensures good tail behavior for the importance-weighted estimator $g(\cdot; \pi, z)/p(\pi)$. Here, tail behavior is captured by the MGF-like term in (7), which penalizes the learner for over-estimating rewards under the reference distribution q or under-estimating rewards under π^* .

We show that this approach leads to a bound on regret that scales with the convexified DEC.

Theorem 2.1 (Main upper bound). *For any choice of $\eta > 0$, [Algorithm 1](#) ensures that for all $\delta > 0$, with probability at least $1 - \delta$,*

$$\mathbf{Reg}_{\text{DM}} \leq \text{dec}_{(8\eta)^{-1}}(\text{co}(\mathcal{M})) \cdot T + 2\eta^{-1} \cdot \log(|\Pi|/\delta). \quad (8)$$

In particular, for any $\delta > 0$, with appropriate η , the algorithm has that with probability at least $1 - \delta$,

$$\mathbf{Reg}_{\text{DM}} \leq O(1) \cdot \inf_{\gamma > 0} \{\text{dec}_{\gamma}(\text{co}(\mathcal{M})) \cdot T + \gamma \cdot \log(|\Pi|/\delta)\}. \quad (9)$$

This should be compared to the upper bound for the stochastic setting in Foster et al. [18] (e.g., [Theorem 3.3](#)), which takes a similar form, but scales with the weaker quantity $\sup_{\bar{M} \in \text{co}(\mathcal{M})} \text{dec}_{\gamma}(\mathcal{M}, \bar{M})$.⁴ See also [Appendix A](#) for a comparison to Lattimore and Szepesvári [37], Lattimore and György [34].

Equivalence of Exploration-by-Optimization and Decision-Estimation Coefficient. We now discuss a deeper connection between Exploration-by-Optimization and the DEC. Define the minimax value of the high-probability Exploration-by-Optimization objective via

$$\text{exo}_{\eta}(\mathcal{M}, q) := \inf_{p \in \Delta(\Pi), g \in \mathcal{G}} \sup_{M \in \mathcal{M}, \pi^* \in \Pi} \Gamma_{q, \eta}(p, g; \pi^*, M), \quad (10)$$

⁴If a proper estimator is available, Foster et al. [18] (Thm. 4.1) gives tighter bounds scaling with $\text{dec}_{\gamma}(\mathcal{M})$.

201 and let $\text{exo}_\eta(\mathcal{M}) := \sup_{q \in \Delta(\Pi)} \text{exo}_\eta(\mathcal{M}, q)$. This quantity can be interpreted as a complexity
 202 measure for \mathcal{M} whose value reflects the difficulty of exploration. The following structural result
 203 (Corollary 3.1 in Section 3), which is critical to the proof of Theorem 2.1, shows that this complexity
 204 measure is equivalent to the convexified Decision-Estimation Coefficient:

$$\text{dec}_{(4\eta)^{-1}}(\text{co}(\mathcal{M})) \leq \text{exo}_\eta(\mathcal{M}) \leq \text{dec}_{(8\eta)^{-1}}(\text{co}(\mathcal{M})), \quad \forall \eta > 0. \quad (11)$$

205 As we show, the regret of Algorithm 1 is controlled by the value of $\text{exo}_\eta(\mathcal{M})$, and thus Theorem 2.1
 206 follows. This result builds on, but goes beyond the Bregman divergence-based framework in Lattimore
 207 and György [34], and exploits a somewhat obscure connection between Hellinger distance and the
 208 moment generating function (MGF) for the logarithmic loss. In particular, we use a technical lemma
 209 (proven in Appendix C), which shows that up to constants, the value of Hellinger distance between
 210 two probability distributions can be expressed as variational problem based on the associated MGFs.

211 **Lemma 2.1.** *Let \mathbb{P} and \mathbb{Q} be probability distributions over a measurable space $(\mathcal{X}, \mathcal{F})$. Then*

$$\frac{1}{2} D_{\text{H}}^2(\mathbb{P}, \mathbb{Q}) \leq \sup_{g: \mathcal{X} \rightarrow \mathbb{R}} \{1 - \mathbb{E}_{\mathbb{P}}[e^g] \cdot \mathbb{E}_{\mathbb{Q}}[e^{-g}]\} \leq D_{\text{H}}^2(\mathbb{P}, \mathbb{Q}). \quad (12)$$

212 The lower inequality in Lemma 2.1 is proven using a trick similar to one used by Zhang [62] to prove
 213 high-probability bounds for maximum likelihood estimation based on Hellinger distance. In the
 214 process of proving (11), we also establish equivalence of the Exploration-by-Optimization objective
 215 and a *parameterized* version of the information ratio, which is of independent interest (cf. Section 3).
 216

217 **Further remarks.** The main focus of this work is sample complexity, and the runtime and memory
 218 requirements of Algorithm 1—which are linear in $|\Pi|$ —are not practical for large decision spaces.
 219 Improving the computational efficiency is an interesting question for future work. We mention in
 220 passing that Theorem 2.1 answers a question raised by Foster et al. [18] of obtaining in the frequentist
 221 setting a regret bound matching the Bayesian regret bound in their Theorem 3.6.

222 2.2 Lower Bound

223 We now complement the regret bound in the prequel with a lower bound based on the convexified
 224 DEC. Our most general result shows that for any algorithm, either the expected regret or its (one-sided)
 225 second moment must scale with a localized version of the convexified DEC.

226 To state the result, we define the *localized model class* around a model \bar{M} via

$$\mathcal{M}_\varepsilon(\bar{M}) = \{M \in \mathcal{M} : f^{\bar{M}}(\pi_{\bar{M}}) \geq f^M(\pi_M) - \varepsilon\},$$

227 and let $\text{dec}_{\gamma, \varepsilon}(\mathcal{M}) := \sup_{\bar{M} \in \mathcal{M}} \text{dec}_\gamma(\mathcal{M}_\varepsilon(\bar{M}), \bar{M})$ be the *localized Decision-Estimation Coefficient*.

228 Let $(x)_+ := \max\{x, 0\}$ and $V(\mathcal{M}) := \sup_{M, M' \in \mathcal{M}} \sup_{\pi \in \Pi} \sup_{A \in \mathcal{A} \otimes \mathcal{O}} \left\{ \frac{M(A|\pi)}{M'(A|\pi)} \right\} \vee e$.⁵

229 **Theorem 2.2** (Main lower bound). *Let $C(T) := c \cdot \log(T \wedge V(\mathcal{M}))$ for a sufficiently large numerical
 230 constant $c > 0$. Set $\varepsilon_\gamma := \frac{\gamma}{4C(T)T}$. For any algorithm, there exists an oblivious adversary for which*

$$\mathbb{E}[\mathbf{Reg}_{\text{DM}}] + \sqrt{\mathbb{E}(\mathbf{Reg}_{\text{DM}})_+^2} \geq \Omega(1) \cdot \sup_{\gamma > \sqrt{2C(T)T}} \text{dec}_{\gamma, \varepsilon_\gamma}(\text{co}(\mathcal{M})) \cdot T - O(T^{1/2}). \quad (13)$$

231 Theorem 2.2 implies that for any algorithm with “reasonable” tail behavior beyond what is granted by
 232 control of the first moment (such as Algorithm 1), the regret in Theorem 2.1 cannot be substantially
 233 improved. In more detail, consider the notion of a *sub-Chebychev* algorithm.

234 **Definition 2.1** (Sub-Chebychev Algorithm). *We say that a regret minimization algorithm is sub-
 235 Chebychev with parameter R if for all $t > 0$,*

$$\mathbb{P}((\mathbf{Reg}_{\text{DM}})_+ \geq t) \leq R^2/t^2. \quad (14)$$

236 For sub-Chebychev algorithms, both the mean and (root) second moment of regret are bounded by
 237 the parameter R (cf. Appendix F.4), which has the following consequence.

⁵Recall (Appendix B) that $M(\cdot, \cdot | \pi)$ is the conditional distribution given π ; finiteness of $V(\mathcal{M})$ is not necessary, but removes a $\log(T)$ factor from Theorem 2.2.

238 **Corollary 2.1.** Any regret minimization algorithm with sub-Chebychev parameter $R > 0$ must have

$$R \geq \tilde{\Omega}(1) \cdot \sup_{\gamma > \sqrt{2C(T)T}} \text{dec}_{\gamma, \varepsilon_\gamma}(\text{co}(\mathcal{M})) \cdot T - O(T^{1/2}). \quad (15)$$

239 To interpret this result, suppose for simplicity that $\text{dec}_\gamma(\text{co}(\mathcal{M}))$ and $\text{dec}_{\gamma, \varepsilon_\gamma}(\text{co}(\mathcal{M}))$ are continuous
 240 with respect to $\gamma > 0$, and that $\text{dec}_{\gamma, \varepsilon_\gamma}(\text{co}(\mathcal{M})) \gtrsim \gamma^{-1}$, which is satisfied for all non-trivial classes.⁶
 241 In this case, one can show (cf. [Proposition F.2](#) for a proof) that by setting $\delta = 1/T^2$, [Theorem 2.1](#)
 242 implies that [Algorithm 1](#) is sub-Chebychev with parameter

$$R = \tilde{O}\left(\inf_{\gamma > 0} \{\text{dec}_\gamma(\text{co}(\mathcal{M})) \cdot T + \gamma \cdot \log(|\Pi|)\}\right) = \tilde{O}(\text{dec}_{\gamma_u}(\text{co}(\mathcal{M})) \cdot T), \quad (16)$$

243 where γ_u satisfies the balance $\text{dec}_{\gamma_u}(\text{co}(\mathcal{M})) \propto \frac{\gamma_u}{T} \log|\Pi|$. On the other hand, the lower bound in
 244 [\(15\)](#) can be shown to scale with

$$R \geq \tilde{\Omega}\left(\text{dec}_{\gamma_\ell, \varepsilon_{\gamma_\ell}}(\text{co}(\mathcal{M})) \cdot T\right), \quad (17)$$

245 where γ_ℓ satisfies the balance $\text{dec}_{\gamma_\ell, \varepsilon_{\gamma_\ell}}(\text{co}(\mathcal{M})) \propto \frac{\gamma_\ell}{T}$. We conclude that the upper bound from
 246 [Theorem 2.1](#) cannot be improved beyond (i) localization and (ii) dependence on $\log|\Pi|$.

247 As an example, we show in [Appendix D.3](#) that for the multi-armed bandit problem with $\Pi =$
 248 $\{1, \dots, A\}$, the upper bound in [\(16\)](#) yields $R = O(\sqrt{AT \log A})$, while the lower bound in [\(17\)](#) yields
 249 $R = \Omega(\sqrt{AT})$. See [Appendix D](#) for additional examples which further illustrate the scaling above.

250 The dependence on $\log|\Pi|$ cannot be removed from the upper bound or made to appear in the lower
 251 bound in general (cf. Section 3.5 of Foster et al. [18]). As shown in Foster et al. [18], localization is
 252 inconsequential for essentially all model classes commonly studied in the literature, and the same is
 253 true for the examples we consider here ([Appendix D](#)), where [Theorem 2.2](#) leads to the correct rate up
 254 to small polynomial factors. However, improving the upper bound to achieve localization (which
 255 Foster et al. [18] show is possible in the stochastic setting) is an interesting future direction.

256 See [Appendix A](#) for further discussion and for comparison to a related lower bound in Lattimore [33].

257 **Why convexity?** At this point, a natural question is *why* the convex hull $\text{co}(\mathcal{M})$ plays a fundamental
 258 role in the adversarial setting. For the lower bound, the intuition is simple: Given a model class \mathcal{M} ,
 259 the adversary can pick any mixture distribution $\mu \in \Delta(\mathcal{M})$, then choose the sequence of models
 260 $M^{(1)}, \dots, M^{(T)}$ by sampling $M^{(t)} \sim \mu$ independently at each round. This is equivalent to playing a
 261 static mixture model $M^* = \mathbb{E}_{M \sim \mu}[M] \in \text{co}(\mathcal{M})$, which is what allows us to prove a lower bound
 262 based on the DEC for the set $\text{co}(\mathcal{M})$ of all such models. In view of the fact that the lower bound is
 263 obtained through this static, stochastic adversary, we believe the more surprising result here is that
 264 good behavior of the convexified DEC is also *sufficient* for low regret.

265 2.3 Learnability and Comparison to Stochastic Setting

266 Building on the upper and lower bounds in the prequel, we give a characterization for *learnability*
 267 (i.e., when non-trivial regret is possible) in the adversarial setting. This extends the learnability result
 268 for the stochastic setting in Foster et al. [18], and follows a long tradition of such characterizations in
 269 learning theory [58, 3, 54, 49, 12]. To state the result, we define the minimax regret as

$$\mathfrak{R}(\mathcal{M}, T) = \inf_{\mathbf{p}^{(1)}, \dots, \mathbf{p}^{(T)}} \sup_{M^{(1)}, \dots, M^{(T)}} \mathbb{E}[\mathbf{Reg}_{\text{DM}}],$$

270 where $\mathbf{p}^{(t)} : (\Pi \times \mathcal{Z})^{t-1} \rightarrow \Delta(\Pi)$ and $M^{(t)} : (\Pi \times \mathcal{Z})^{t-1} \rightarrow \mathcal{M}$ are policies for the learner and
 271 adversary, respectively. Our characterization is as follows.

272 **Theorem 2.3.** *Suppose there exists $M_0 \in \mathcal{M}$ such that f^{M_0} is a constant function, and that $|\Pi| < \infty$.*

273 1. *If there exists $\rho > 0$ s.t. $\lim_{\gamma \rightarrow \infty} \text{dec}_\gamma(\text{co}(\mathcal{M})) \cdot \gamma^\rho = 0$, then $\lim_{T \rightarrow \infty} \frac{\mathfrak{R}(\mathcal{M}, T)}{T^p} = 0$ for $p < 1$.*

274 2. *If $\lim_{\gamma \rightarrow \infty} \text{dec}_\gamma(\text{co}(\mathcal{M})) \cdot \gamma^\rho > 0$ for all $\rho > 0$, then $\lim_{T \rightarrow \infty} \frac{\mathfrak{R}(\mathcal{M}, T)}{T^p} = \infty$ for all $p < 1$.*

⁶Note that the dominant term $\text{dec}_{\gamma, \varepsilon_\gamma}(\text{co}(\mathcal{M})) \cdot T$ in [\(13\)](#) scales with \sqrt{T} any “non-trivial” class that embeds the two-armed bandit problem, so that the $-O(T^{1/2})$ term can be discarded.

275 The same conclusion holds when $\Pi = \Pi_T$ grows with T , but has $\log|\Pi_T| = O(T^q)$ for any $q < 1$.⁷

276 **Theorem 2.3** shows that polynomial decay of the convexified DEC is necessary and sufficient for low
 277 regret. We emphasize that this result is complementary to [Theorem 2.2](#), and does not require local-
 278 ization or any assumption on the tail behavior of the algorithm. This is a consequence of the coarse,
 279 asymptotic nature of the result, which allows us to perform rescaling tricks to remove these conditions.

280 **Comparison to stochastic setting.** Having shown that the convexified Decision-Estimation
 281 Coefficient plays a fundamental role in determining the optimal regret for the adversarial DMSO
 282 setting, now is a good time to make comparisons to the stochastic setting. There, Foster et al. [18]
 283 obtain upper bounds on regret that have the same form as (9), but scale with the weaker quantity
 284 $\max_{\bar{M} \in \text{co}(\mathcal{M})} \text{dec}_\gamma(\mathcal{M}, \bar{M})$.⁸ For classes that are not convex, but where “proper” estimators are
 285 available (e.g., tabular MDPs), the upper bounds in Foster et al. [18] can further be improved to scale
 286 with $\text{dec}_\gamma(\mathcal{M})$. Hence, our results show that in general, the price of adversarial outcomes can be as
 287 large as $\text{dec}_\gamma(\text{co}(\mathcal{M}))/\text{dec}_\gamma(\mathcal{M})$. Examples (see [Appendix D](#) for details and more) include:

- 288 • For tabular MDPs with horizon H , S states, and A actions, Foster et al. [18] show that $\text{dec}_\gamma(\mathcal{M}) =$
 289 $\text{poly}(H, S, A)/\gamma$, and use this to obtain regret $\sqrt{\text{poly}(H, S, A) \cdot T}$. Tabular MDPs are *not* a
 290 convex class, and $\text{co}(\mathcal{M})$ is equivalent to the class of so-called *latent MDPs*, which are known to
 291 be intractable [30, 41]. Indeed, we show ([Appendix D](#)) that $\text{dec}_\gamma(\text{co}(\mathcal{M})) \geq \Omega(A^{\min\{S, H\}})$. This
 292 example highlights that in general, the gap between stochastic and adversarial can be quite large.
- 293 • For many common bandit problems, one has $\text{co}(\mathcal{M}) = \mathcal{M}$, leading to polynomial bounds on
 294 regret in the adversarial setting. For example the multi-armed bandit problem with A actions has
 295 $\text{dec}_\gamma(\text{co}(\mathcal{M})) \leq O(A/\gamma)$, leading to $\sqrt{AT \log A}$ regret from [Theorem 2.1](#), and the linear bandit
 296 problem in d dimensions has $\text{dec}_\gamma(\text{co}(\mathcal{M})) \leq O(d/\gamma)$, leading to regret $\sqrt{dT \log |\Pi|}$.

297 3 Connections Between Complexity Measures

298 The Decision-Estimation Coefficient bears a resemblance to the notion of *generalized information*
 299 *ratio* introduced by Lattimore and György [34], Lattimore [32] which extends the original information
 300 ratio of Russo and Van Roy [51, 52]. In what follows, we establish deeper connections between these
 301 complexity measures. All of the results in this section are proven in [Appendix E](#).

302 Let us recall the definition of the generalized information ratio from Lattimore [32], which we state
 303 here for a general divergence-like function $D(\cdot \| \cdot) \rightarrow \mathbb{R}^+$ (typically, KL divergence or another
 304 Bregman divergence). For a distribution $\mu \in \Delta(\mathcal{M} \times \Pi)$ and decision distribution $p \in \Delta(\Pi)$, define
 305 $\mu_{\text{pr}}(\pi') := \mathbb{P}(\pi^* = \pi')$ and $\mu_{\text{po}}(\pi'; \pi, z) := \mathbb{P}(\pi^* = \pi' \mid (\pi, z))$, where \mathbb{P} is the law of the process
 306 $(M, \pi^*) \sim \mu, \pi \sim p, z \sim M(\pi)$. μ_{pr} should be thought of as the prior over π^* , and μ_{po} as the
 307 posterior having observed (z, π) ; note that the law μ_{po} does not depend on the distribution p . For
 308 parameter $\lambda > 1$, Lattimore [33] defines the generalized information ratio for a class \mathcal{M} via⁹

$$\Psi_\lambda(\mathcal{M}) = \sup_{\mu \in \Delta(\mathcal{M} \times \Pi)} \inf_{p \in \Delta(\Pi)} \left\{ \frac{(\mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)])^\lambda}{\mathbb{E}_{\pi \sim p} \mathbb{E}_{z \mid \pi} [D(\mu_{\text{po}}(\cdot; \pi, z) \parallel \mu_{\text{pr}})]} \right\}. \quad (18)$$

309 Here, we have slightly generalized the original definition in Lattimore [33] by incorporating models in
 310 \mathcal{M} rather than placing an arbitrary prior over observations z directly. We also use a general divergence,
 311 while Lattimore [33] uses KL divergence and Lattimore and György [34] use Bregman divergences.

312 To understand the connection to the Decision-Estimation Coefficient, it will be helpful introduce
 313 another variant of the information ratio that we call the *parameterized information ratio*.

314 **Definition 3.1.** For a divergence $D(\cdot \| \cdot)$, the parameterized information ratio is given by

$$\begin{aligned} & \inf_\gamma^D(\mathcal{M}) && (19) \\ &= \sup_{\mu \in \Delta(\mathcal{M} \times \Pi)} \inf_{p \in \Delta(\Pi)} \mathbb{E}_{\pi \sim p} [\mathbb{E}_{(M, \pi^*) \sim \mu} [f^M(\pi^*) - f^M(\pi)] - \gamma \cdot \mathbb{E}_{\pi \sim p} \mathbb{E}_{z \mid \pi} [D(\mu_{\text{po}}(\cdot; \pi, z) \parallel \mu_{\text{pr}})]] \end{aligned}$$

⁷Allowing Π to grow with T can be used to handle infinite decision spaces using covering arguments.

⁸Theorem 3.1 of Foster et al. [18] attains $\mathbf{Reg}_{\text{DM}} \lesssim \inf_{\gamma > 0} \{ \max_{\bar{M} \in \text{co}(\mathcal{M})} \text{dec}_\gamma(\mathcal{M}, \bar{M}) + \gamma \cdot \log |\mathcal{M}| \}$.

⁹Lattimore and György [34] give a slightly different but essentially equivalent definition; cf. [Appendix E](#).

315 The parameterized information ratio is always bounded by the generalized information ratio in (18); in
 316 particular, we have $\inf_{\gamma}^D(\mathcal{M}) \leq (\Psi_{\lambda}(\mathcal{M})/\gamma)^{\frac{1}{\lambda-1}} \forall \gamma > 0$. All of the regret bounds based on the gen-
 317 eralized information ratio that we are aware of [34, 33] implicitly bound regret by the parameterized
 318 information ratio, and then invoke the inequality above to move to the generalized information ratio.
 319 In general though, it does not appear that these notions are equivalent. Informally, this is because
 320 the notion in (18) is equivalent to requiring that a single distribution p certify a certain bound on the
 321 value in (19) for all values of the parameter γ simultaneously, while the parameterized information
 322 ratio allows the distribution p to vary as a function of $\gamma > 0$ (hence the name); see also Appendix E.

323 Letting $\inf_{\gamma}^H(\mathcal{M})$ denote the parameterized information ratio with $D = D_{\Pi}^2(\cdot, \cdot)$, we show that this
 324 notion is equivalent to the convexified Decision-Estimation Coefficient.

325 **Theorem 3.1.** *For all $\gamma > 0$, $\inf_{\gamma}^H(\mathcal{M}) \leq \text{dec}_{\gamma}(\text{co}(\mathcal{M})) \leq \inf_{\gamma/4}^H(\mathcal{M})$.*

326 This result is a special case of Theorem E.1 in Appendix E, which shows that a similar equivalence
 327 holds for a class of “well-behaved” f -divergences that includes KL divergence (but not necessarily
 328 for general Bregman divergences). The basic idea is to use Bayes’ rule to move from the Decision-
 329 Estimation Coefficient, which considers distance between distributions over *observations*, to the
 330 information ratio, which considers distance between distributions over *decisions*.

331 In light of this characterization, the results in this paper could have equivalently been presented in
 332 terms of the parameterized information ratio. We chose to present them in terms of the Decision-
 333 Estimation Coefficient in order to draw parallels to the stochastic setting, where guarantees that scale
 334 with $\text{dec}_{\gamma}(\mathcal{M})$ (without convexification) are available. It is unclear whether the information ratio
 335 can accurately reflect the complexity for both stochastic and adversarial settings in the same fashion,
 336 because—unlike the DEC—it is invariant under convexification.¹⁰

337 **Proposition 3.1.** *For any divergence-like function $D(\cdot \parallel \cdot) : \Delta(\Pi) \times \Delta(\Pi) \rightarrow \mathbb{R}_+$, we have*

$$\inf_{\gamma}^D(\mathcal{M}) = \inf_{\gamma}^D(\text{co}(\mathcal{M})), \quad \forall \gamma > 0.$$

338 For a final structural result, we show that up to constants, the parameterized information ratio is
 339 equivalent to the high-probability Exploration-by-Optimization objective.

340 **Theorem 3.2.** *For all $\eta > 0$, $\inf_{\eta-1}^H(\mathcal{M}) \leq \text{exo}_{\eta}(\mathcal{M}) \leq \inf_{(8\eta)-1}^H(\mathcal{M})$.*

341 This result is proven through a direct argument, and the equivalence of the DEC and Exploration-
 342 by-Optimization in (11) is proven by combining with Theorem 3.1. Summarizing the equivalence:

343 **Corollary 3.1.** *For all $\eta > 0$,*

$$\text{dec}_{(4\eta)-1}(\text{co}(\mathcal{M})) \leq \inf_{\eta-1}^H(\mathcal{M}) \leq \text{exo}_{\eta}(\mathcal{M}) \leq \inf_{(8\eta)-1}^H(\mathcal{M}) \leq \text{dec}_{(8\eta)-1}(\text{co}(\mathcal{M})).$$

344 Since this equivalence depends of the value of the parameter $\gamma > 0$ in the parameterized information
 345 ratio, it seems unlikely that a similar equivalence can be established using the generalized information
 346 ratio in (18). We note in passing that one can use similar techniques to lower bound the Bregman
 347 divergence-based Exploration-by-Optimization objective in Lattimore and György [34] by the param-
 348 eterized information ratio for the Bregman divergence of interest, complementing their upper bound.

349 4 Discussion

350 We have shown that the convexified Decision-Estimation Coefficient is necessary and sufficient
 351 to achieve low regret for adversarial interactive decision making, establishing that convexity
 352 governs the price of adversarial outcomes. Our results elucidate the relationship between the DEC,
 353 Exploration-by-Optimization, and the information ratio, and we hope they will find broader use.

354 Our results add to a growing body of research which shows that online reinforcement learning with
 355 agnostic or adversarial outcomes can be statistically intractable [53, 41]. A promising future direction
 356 is to extend our techniques to natural semi-adversarial models in which reinforcement learning is
 357 tractable (for example, the so-called *adversarially corrupted* setting [42, 19]). Other interesting ques-
 358 tions include (i) extending our lower bounds beyond the observable-reward setting and to directly han-
 359 dle expected regret, and (ii) developing computationally efficient algorithms for large decision spaces.

¹⁰The variants in Lattimore and György [34], Lattimore [33] are also invariant under convexification.

References

- [1] J. Abernethy, E. Hazan, and A. Rakhlin. Competing in the dark: An efficient algorithm for bandit linear optimization. In *Proc. of the 21st Annual Conference on Learning Theory (COLT)*, 2008.
- [2] A. Agarwal, S. Bird, M. Cozowicz, L. Hoang, J. Langford, S. Lee, J. Li, D. Melamed, G. Oshri, O. Ribas, S. Sen, and A. Slivkins. Making contextual decisions with low technical debt. *arXiv:1606.03966*, 2016.
- [3] N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haussler. Scale-sensitive dimensions, uniform convergence, and learnability. *Journal of the ACM*, 44:615–631, 1997.
- [4] J.-Y. Audibert and S. Bubeck. Minimax policies for adversarial and stochastic bandits. In *COLT*, volume 7, pages 1–122, 2009.
- [5] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.
- [6] A. Ayoub, Z. Jia, C. Szepesvari, M. Wang, and L. Yang. Model-based reinforcement learning with value-targeted regression. In *International Conference on Machine Learning*, pages 463–474. PMLR, 2020.
- [7] S. Bubeck, R. Munos, G. Stoltz, and C. Szepesvári. X-armed bandits. *Journal of Machine Learning Research*, 12(5), 2011.
- [8] S. Bubeck, N. Cesa-Bianchi, and S. M. Kakade. Towards minimax policies for online linear optimization with bandit feedback. In *Conference on Learning Theory*, pages 41–1. JMLR Workshop and Conference Proceedings, 2012.
- [9] S. Bubeck, Y. T. Lee, and R. Eldan. Kernel-based methods for bandit convex optimization. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pages 72–85, 2017.
- [10] N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, New York, NY, USA, 2006. ISBN 0521841089.
- [11] V. Dani, T. P. Hayes, and S. Kakade. The price of bandit information for online optimization. 2007.
- [12] A. Daniely, S. Sabato, S. Ben-David, and S. Shalev-Shwartz. Multiclass learnability and the erm principle. In *Proceedings of the 24th Annual Conference on Learning Theory*, pages 207–232, 2011.
- [13] S. Dean, H. Mania, N. Matni, B. Recht, and S. Tu. On the sample complexity of the linear quadratic regulator. *Foundations of Computational Mathematics*, 20(4):633–679, 2020.
- [14] S. Dong, B. Van Roy, and Z. Zhou. Provably efficient reinforcement learning with aggregated states. *arXiv preprint arXiv:1912.06366*, 2019.
- [15] S. Du, A. Krishnamurthy, N. Jiang, A. Agarwal, M. Dudik, and J. Langford. Provably efficient RL with rich observations via latent state decoding. In *International Conference on Machine Learning*, pages 1665–1674. PMLR, 2019.
- [16] S. S. Du, S. M. Kakade, J. D. Lee, S. Lovett, G. Mahajan, W. Sun, and R. Wang. Bilinear classes: A structural framework for provable generalization in RL. *International Conference on Machine Learning*, 2021.
- [17] A. D. Flaxman, A. T. Kalai, and H. B. McMahan. Online convex optimization in the bandit setting: gradient descent without a gradient. In *Proceedings of the sixteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 385–394, 2005.
- [18] D. J. Foster, S. M. Kakade, J. Qian, and A. Rakhlin. The statistical complexity of interactive decision making. *arXiv preprint arXiv:2112.13487*, 2021.
- [19] A. Gupta, T. Koren, and K. Talwar. Better algorithms for stochastic bandits with adversarial corruptions. In *Conference on Learning Theory*, pages 1562–1578. PMLR, 2019.
- [20] E. Hazan and S. Kale. Better algorithms for benign bandits. *Journal of Machine Learning Research*, 12(4), 2011.

- 410 [21] N. Jiang, A. Krishnamurthy, A. Agarwal, J. Langford, and R. E. Schapire. Contextual decision
411 processes with low Bellman rank are PAC-learnable. In *International Conference on Machine*
412 *Learning*, pages 1704–1713, 2017.
- 413 [22] C. Jin, Z. Yang, Z. Wang, and M. I. Jordan. Provably efficient reinforcement learning with
414 linear function approximation. In *Conference on Learning Theory*, pages 2137–2143, 2020.
- 415 [23] C. Jin, Q. Liu, and S. Miryoosefi. Bellman eluder dimension: New rich classes of RL problems,
416 and sample-efficient algorithms. *Neural Information Processing Systems*, 2021.
- 417 [24] T. Jin and H. Luo. Simultaneously learning stochastic and adversarial episodic mdps with
418 known transition. *Advances in neural information processing systems*, 33:16557–16566, 2020.
- 419 [25] J. Kirschner, T. Lattimore, and A. Krause. Information directed sampling for linear partial
420 monitoring. In *Conference on Learning Theory*, pages 2328–2369. PMLR, 2020.
- 421 [26] J. Kirschner, T. Lattimore, C. Vernade, and C. Szepesvári. Asymptotically optimal information-
422 directed sampling. In *Conference on Learning Theory*, pages 2777–2821. PMLR, 2021.
- 423 [27] R. Kleinberg. Nearly tight bounds for the continuum-armed bandit problem. *Advances in*
424 *Neural Information Processing Systems*, 17:697–704, 2004.
- 425 [28] J. Kober, J. A. Bagnell, and J. Peters. Reinforcement learning in robotics: A survey. *The*
426 *International Journal of Robotics Research*, 32(11):1238–1274, 2013.
- 427 [29] A. Krishnamurthy, A. Agarwal, and J. Langford. PAC reinforcement learning with rich observa-
428 tions. In *Advances in Neural Information Processing Systems*, pages 1840–1848, 2016.
- 429 [30] J. Kwon, Y. Efroni, C. Caramanis, and S. Mannor. RL for latent mdps: Regret guarantees and a
430 lower bound. *Advances in Neural Information Processing Systems*, 34, 2021.
- 431 [31] T. Lattimore. Improved regret for zeroth-order adversarial bandit convex optimisation. *arXiv*
432 *preprint arXiv:2006.00475*, 2020.
- 433 [32] T. Lattimore. Minimax regret for bandit convex optimisation of ridge functions. *arXiv preprint*
434 *arXiv:2106.00444*, 2021.
- 435 [33] T. Lattimore. Minimax regret for partial monitoring: Infinite outcomes and rustichini’s regret.
436 *arXiv preprint arXiv:2202.10997*, 2022.
- 437 [34] T. Lattimore and A. György. Mirror descent and the information ratio. In *Conference on*
438 *Learning Theory*, pages 2965–2992. PMLR, 2021.
- 439 [35] T. Lattimore and C. Szepesvári. An information-theoretic approach to minimax regret in partial
440 monitoring. In *Conference on Learning Theory*, pages 2111–2139. PMLR, 2019.
- 441 [36] T. Lattimore and C. Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- 442 [37] T. Lattimore and C. Szepesvári. Exploration by optimisation in partial monitoring. In *Conference*
443 *on Learning Theory*, pages 2488–2515. PMLR, 2020.
- 444 [38] J. Li, W. Monroe, A. Ritter, D. Jurafsky, M. Galley, and J. Gao. Deep reinforcement learning
445 for dialogue generation. In *EMNLP*, 2016.
- 446 [39] L. Li. *A unifying framework for computational reinforcement learning theory*. Rutgers, The
447 State University of New Jersey—New Brunswick, 2009.
- 448 [40] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra.
449 Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- 450 [41] Q. Liu, Y. Wang, and C. Jin. Learning markov games with adversarial opponents: Efficient
451 algorithms and fundamental limits. *arXiv preprint arXiv:2203.06803*, 2022.
- 452 [42] T. Lykouris, V. Mirrokni, and R. Paes Leme. Stochastic bandits robust to adversarial corruptions.
453 In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages
454 114–122, 2018.
- 455 [43] S. Magureanu, R. Combes, and A. Proutiere. Lipschitz bandits: Regret lower bound and optimal
456 algorithms. In *Conference on Learning Theory*, pages 975–999. PMLR, 2014.
- 457 [44] A. Modi, N. Jiang, A. Tewari, and S. Singh. Sample complexity of reinforcement learning using
458 linearly combined model ensembles. In *International Conference on Artificial Intelligence and*
459 *Statistics*, pages 2010–2020. PMLR, 2020.

- 460 [45] G. Neu, A. György, C. Szepesvári, et al. The online loop-free stochastic shortest-path problem.
461 In *COLT*, volume 2010, pages 231–243. Citeseer, 2010.
- 462 [46] G. Neu, A. György, C. Szepesvári, and A. Antos. Online markov decision processes under
463 bandit feedback. *IEEE Transactions on Automatic Control*, 59:676–691, 2014.
- 464 [47] Y. Polyanskiy. Information theoretic methods in statistics and computer science. 2020. URL
465 https://people.lids.mit.edu/yp/homepage/sdpi_course.html.
- 466 [48] Y. Polyanskiy and Y. Wu. Lecture notes on information theory. 2014.
- 467 [49] A. Rakhlin, K. Sridharan, and A. Tewari. Online learning: Random averages, combinatorial
468 parameters, and learnability. *Advances in Neural Information Processing Systems 23*, pages
469 1984–1992, 2010.
- 470 [50] D. Russo and B. Van Roy. Eluder dimension and the sample complexity of optimistic exploration.
471 In *Advances in Neural Information Processing Systems*, pages 2256–2264, 2013.
- 472 [51] D. Russo and B. Van Roy. Learning to optimize via posterior sampling. *Mathematics of*
473 *Operations Research*, 39(4):1221–1243, 2014.
- 474 [52] D. Russo and B. Van Roy. Learning to optimize via information-directed sampling. *Operations*
475 *Research*, 66(1):230–252, 2018.
- 476 [53] A. Sekhari, C. Dann, M. Mohri, Y. Mansour, and K. Sridharan. Agnostic reinforcement learning
477 with low-rank mdps and rich observations. *Advances in Neural Information Processing Systems*,
478 34, 2021.
- 479 [54] S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridharan. Learnability, stability, and uniform
480 convergence. *Journal of Machine Learning Research (JMLR)*, 2010.
- 481 [55] M. Sion. On general minimax theorems. *Pacific J. Math.*, 8:171–176, 1958.
- 482 [56] W. Sun, N. Jiang, A. Krishnamurthy, A. Agarwal, and J. Langford. Model-based RL in
483 contextual decision processes: PAC bounds and exponential improvements over model-free
484 approaches. In *Conference on learning theory*, pages 2898–2933. PMLR, 2019.
- 485 [57] A. Tewari and S. A. Murphy. From ads to interventions: Contextual bandits in mobile health.
486 In *Mobile Health*, 2017.
- 487 [58] V. N. Vapnik. The nature of statistical learning theory. 1995.
- 488 [59] R. Wang, R. R. Salakhutdinov, and L. Yang. Reinforcement learning with general value function
489 approximation: Provably efficient approach via bounded eluder dimension. *Advances in Neural*
490 *Information Processing Systems*, 33, 2020.
- 491 [60] L. Yang and M. Wang. Sample-optimal parametric Q-learning using linearly additive features.
492 In *International Conference on Machine Learning*, pages 6995–7004. PMLR, 2019.
- 493 [61] Y. Yang and A. R. Barron. An asymptotic property of model selection criteria. *IEEE Transac-*
494 *tions on Information Theory*, 44(1):95–116, 1998.
- 495 [62] T. Zhang. From ϵ -entropy to KL-entropy: Analysis of minimum information complexity density
496 estimation. *The Annals of Statistics*, 34(5):2180–2210, 2006.
- 497 [63] D. Zhou, Q. Gu, and C. Szepesvari. Nearly minimax optimal reinforcement learning for linear
498 mixture markov decision processes. In *Conference on Learning Theory*, pages 4532–4576.
499 PMLR, 2021.
- 500 [64] A. Zimin and G. Neu. Online learning in episodic markovian decision processes by relative
501 entropy policy search. *Advances in neural information processing systems*, 26, 2013.

502 Checklist

503 The checklist follows the references. Please read the checklist guidelines carefully for information on
504 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
505 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
506 the appropriate section of your paper or providing a brief inline description. For example:

- 507 • Did you include the license to the code and datasets? **[Yes]**
- 508 • Did you include the license to the code and datasets? **[No]** The code and the data are
509 proprietary.
- 510 • Did you include the license to the code and datasets? **[N/A]**

511 Please do not modify the questions and only use the provided macros for your answers. Note that the
512 Checklist section does not count towards the page limit. In your paper, please delete this instructions
513 block and only keep the Checklist section heading above along with the questions/answers below.

- 514 1. For all authors...
 - 515 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
516 contributions and scope? **[Yes]**
 - 517 (b) Did you describe the limitations of your work? **[Yes]**
 - 518 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]**
 - 519 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
520 them? **[Yes]**
- 521 2. If you are including theoretical results...
 - 522 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
 - 523 (b) Did you include complete proofs of all theoretical results? **[Yes]**
- 524 3. If you ran experiments...
 - 525 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
526 mental results (either in the supplemental material or as a URL)? **[N/A]**
 - 527 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
528 were chosen)? **[N/A]**
 - 529 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
530 ments multiple times)? **[N/A]**
 - 531 (d) Did you include the total amount of compute and the type of resources used (e.g., type
532 of GPUs, internal cluster, or cloud provider)? **[N/A]**
- 533 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 534 (a) If your work uses existing assets, did you cite the creators? **[N/A]**
 - 535 (b) Did you mention the license of the assets? **[N/A]**
 - 536 (c) Did you include any new assets either in the supplemental material or as a URL? **[N/A]**
537
 - 538 (d) Did you discuss whether and how consent was obtained from people whose data you're
539 using/curating? **[N/A]**
 - 540 (e) Did you discuss whether the data you are using/curating contains personally identifiable
541 information or offensive content? **[N/A]**
- 542 5. If you used crowdsourcing or conducted research with human subjects...
 - 543 (a) Did you include the full text of instructions given to participants and screenshots, if
544 applicable? **[N/A]**
 - 545 (b) Did you describe any potential participant risks, with links to Institutional Review
546 Board (IRB) approvals, if applicable? **[N/A]**
 - 547 (c) Did you include the estimated hourly wage paid to participants and the total amount
548 spent on participant compensation? **[N/A]**

549	Contents of Appendix	
550	A Detailed Discussion of Related Work	15
551	B Preliminaries	16
552	C Technical Tools	17
553	C.1 Tail Bounds	17
554	C.2 Minimax Theorem	17
555	C.3 Information Theory	17
556	C.4 Online Learning	19
557	D Examples	19
558	D.1 Structured Bandits	19
559	D.2 Reinforcement Learning	20
560	D.3 Proofs for Examples	21
561	E Structural Results	24
562	E.1 Background on Complexity Measures	25
563	E.2 Decision-Estimation Coefficient and Information Ratio (Theorem 3.1)	26
564	E.3 High-Probability Exploration-By-Optimization and Information Ratio (Theorem 3.2)	28
565	F Proofs for Main Results (Section 2)	30
566	F.1 Proof of Theorem 2.1	30
567	F.2 Proof of Theorem 2.2	32
568	F.3 Proof of Theorem 2.3	36
569	F.4 Sub-Chebychev Algorithms	37

570 **A Detailed Discussion of Related Work**

571 Beyond Foster et al. [18], which was the starting point for this work, our results build on a long line
 572 of research on partial monitoring and the information ratio [51, 52, 35, 31, 32, 34, 33, 25, 26]; most
 573 closely related are the works of Lattimore and György [34] and Lattimore [33]. Below we
 574 discuss and compare to these results in greater detail.

575 **Comparison to partial monitoring setting.** Lattimore and György [34], Lattimore [33] and other
 576 works in this sequence consider a general partial monitoring setting in which each outcome $z^{(t)}$ is
 577 directly chosen by an adversary, and need not contain a reward signal.

- 578 • In terms of reward signal, our setting is more restrictive because we assume that $r^{(t)}$ is
 579 observed. Our upper bounds in fact paper encompass the more general setting where rewards
 580 are not observed by the learner, thus subsuming the partial monitoring problem, but our
 581 lower bounds that require that rewards are observed.
- 582 • In terms of data generation process, our setting is more general because we restrict to models
 583 in a known class in \mathcal{M} . This setup recovers the case where $z^{(t)}$ is fully adversarial because
 584 we can take \mathcal{M} to consist of point masses over \mathcal{Z} as a special case. However, the model also
 585 allows for semi-stochastic adversaries, and for settings like (structured) adversarial MDPs.
 586 For example, if all models in \mathcal{M} place ε probability mass on a particular outcome z , any
 587 adversary in our model must place ε mass on this outcome as well.

588 **Upper bounds.** On the upper bound side, our results build on the Exploration-by-Optimization
 589 technique, which was introduced in Lattimore and Szepesvári [37] and generalized significantly in
 590 Lattimore and György [34]. The latter result shows that for a general family of mirror descent-based
 591 Exploration-by-Optimization algorithms parameterized by Bregman divergences, the regret can be
 592 bounded by a certain generalized information ratio based on the associated Bregman divergence (cf.
 593 [Appendix E](#)). This approach yields bounds on expected regret with a similar form to [Theorem 2.1](#)
 594 (with $\text{dec}_\gamma(\text{co}(\mathcal{M}))$ replaced by the generalized information ratio), but does not appear to yield
 595 high-probability bounds (in general, in-expectation regret bounds do not imply high-probability
 596 regret bounds; for example, even for multi-armed bandits, the EXP3 algorithm can experience linear
 597 regret with constant probability [36]). To develop high-probability regret bounds which complement
 598 our lower bounds, we depart from the Bregman divergence-based framework and exploit refined
 599 properties of Hellinger distance. We note that the work of Lattimore and Szepesvári [37] also
 600 proposes a high-probability Exploration-by-Optimization objective, but it is unclear whether this
 601 objective (which precedes the information ratio-based results of Lattimore and György [34]) can be
 602 related to the information ratio or Decision-Estimation Coefficient for general models.¹¹

603 **Lower bounds.** On the lower bound side, we build on the proof strategy from Foster et al. [18].
 604 Our most important technical result is [Theorem F.1](#), which improves upon Theorem 3.1 from Foster
 605 et al. [18] even in the stochastic setting, by using a more refined change of measure argument. In
 606 particular, Theorem 3.1 of Foster et al. [18] gives a lower bound based on the DEC that holds with
 607 low probability, and therefore only provides a meaningful converse to algorithms with sub-Gaussian
 608 or sub-exponential tail behavior. Our result provides a meaningful converse to any upper bound
 609 with sub-Chebychev tail behavior, which is a significantly weaker assumption. We note that while
 610 Theorem 3.2 of Foster et al. [18] provides lower bounds on expected regret without algorithmic
 611 assumptions, this result requires a stronger notion of localization than the one we consider here, and
 612 it is not clear whether this notion can be achieved algorithmically in general. Of course, proving a
 613 lower bound on expected regret that matches our lower bound remains an interesting open problem.

614 Lastly, we mention recent work of Lattimore [33], which provides lower bounds on regret in a general
 615 partial monitoring setting based on a generalized information ratio (cf. [Appendix E](#)). This result is
 616 somewhat complementary to our lower bound ([Theorem 2.2](#)):

- 617 • On the positive side, it leads to lower bounds on *expected regret* that are always tight in
 618 terms of dependence on T , while our result only leads to tight dependence on T if one
 619 restricts to sub-Chebychev algorithms.

¹¹In particular, this objective is based on a Bernstein-type tail bound, which leads to a requirement of boundedness for the estimation functions. We avoid explicitly requiring boundedness using a more specialized tail bound based on [Lemma C.1](#).

620 • On the negative side, the lower bound is loose in $\text{poly}(|\Pi|)$ factors, while our lower bound
 621 is essentially only loose in $\text{poly}(\log|\Pi|)$ factors. As a result, only our lower bound leads to
 622 meaningful dependence on problem-dependence parameters such as dimension for models
 623 with large action spaces.

624 In addition, the lower bound in Lattimore [33] applies to the general partial monitoring setting, while
 625 our lower bound requires that rewards are observed. An interesting question for future work is to
 626 investigate whether the techniques of Lattimore [33] can be combined with our own to get the best of
 627 both worlds.

628 Finally, we mention in passing that the results of Lattimore [33] also imply a learnability characteri-
 629 zation similar to [Theorem 2.3](#). However, because these results are polynomially loose in $|\Pi|$, they
 630 cannot handle the case in which $\log|\Pi|$ grows polynomially in T .

631 B Preliminaries

632 **Basic notation.** For a set \mathcal{X} , we let $\Delta(\mathcal{X})$ denote the set of all Radon probability measures over \mathcal{X} .
 633 We let $\text{co}(\mathcal{X})$ denote the set of all finitely supported convex combinations of elements in \mathcal{X} . We use
 634 the shorthand $x \vee y = \max\{x, y\}$ and $x \wedge y = \min\{x, y\}$.

635 We adopt non-asymptotic big-oh notation: For functions $f, g : \mathcal{X} \rightarrow \mathbb{R}_+$, we write $f = O(g)$ (resp.
 636 $f = \Omega(g)$) if there exists a constant $C > 0$ such that $f(x) \leq Cg(x)$ (resp. $f(x) \geq Cg(x)$) for all
 637 $x \in \mathcal{X}$. We write $f = \tilde{O}(g)$ if $f = O(g \cdot \text{polylog}(T))$, $f = \tilde{\Omega}(g)$ if $f = \Omega(g/\text{polylog}(T))$, and
 638 $f = \tilde{\Theta}(g)$ if $f = \tilde{O}(g)$ and $f = \tilde{\Omega}(g)$. We write $f \propto g$ if $f = \tilde{\Theta}(g)$.

639 **Probability spaces.** We formalize the probability spaces for the DMSO framework in the same fash-
 640 ion as Foster et al. [18], which we briefly summarize here. decisions are associated with a measurable
 641 space (Π, \mathcal{P}) , rewards are associated with the space $(\mathcal{R}, \mathcal{R})$, and observations are associated with
 642 the space $(\mathcal{O}, \mathcal{O})$. The history up to time t is denoted by $\mathcal{H}^{(t)} = (\pi^{(1)}, r^{(1)}, o^{(1)}), \dots, (\pi^{(t)}, r^{(t)}, o^{(t)})$.
 643 We define

$$\Omega^{(t)} = \prod_{i=1}^t (\Pi \times \mathcal{R} \times \mathcal{O}), \quad \text{and} \quad \mathcal{F}^{(t)} = \bigotimes_{i=1}^t (\mathcal{P} \otimes \mathcal{R} \otimes \mathcal{O})$$

644 so that $\mathcal{H}^{(t)}$ is associated with the space $(\Omega^{(t)}, \mathcal{F}^{(t)})$.

645 Formally, a model $M = M(\cdot, \cdot | \cdot) \in \mathcal{M}$ is a probability kernel from (Π, \mathcal{P}) to $(\mathcal{R} \times \mathcal{O}, \mathcal{R} \otimes \mathcal{O})$;
 646 we use the convention $M(\pi) = M(\cdot, \cdot | \pi)$ throughout the paper.¹² An *algorithm* for horizon T is a
 647 sequence $p^{(1)}, \dots, p^{(T)}$, where $p^{(t)}(\cdot | \cdot)$ is a probability kernel from $(\Omega^{(t-1)}, \mathcal{F}^{(t-1)})$ to (Π, \mathcal{P}) .

648 Divergences.

649 For probability distributions \mathbb{P} and \mathbb{Q} over a measurable space (Ω, \mathcal{F}) with a common dominating
 650 measure, we define the total variation distance as

$$D_{\text{TV}}(\mathbb{P}, \mathbb{Q}) = \sup_{A \in \mathcal{F}} |\mathbb{P}(A) - \mathbb{Q}(A)| = \frac{1}{2} \int |d\mathbb{P} - d\mathbb{Q}|.$$

651 Hellinger distance is defined as

$$D_{\text{H}}^2(\mathbb{P}, \mathbb{Q}) = \int (\sqrt{d\mathbb{P}} - \sqrt{d\mathbb{Q}})^2,$$

652 and Kullback-Leibler divergence is defined as

$$D_{\text{KL}}(\mathbb{P} \parallel \mathbb{Q}) = \begin{cases} \int \log\left(\frac{d\mathbb{P}}{d\mathbb{Q}}\right) d\mathbb{P}, & \mathbb{P} \ll \mathbb{Q}, \\ +\infty, & \text{otherwise.} \end{cases}$$

653 For a convex function $f : (0, \infty) \rightarrow \mathbb{R}$, the associated f -divergence for measures \mathbb{P} and \mathbb{Q} with
 654 $\mathbb{P} \ll \mathbb{Q}$ is given by

$$D_f(\mathbb{P} \parallel \mathbb{Q}) := \mathbb{E}_{\mathbb{Q}} \left[f \left(\frac{d\mathbb{P}}{d\mathbb{Q}} \right) \right] \quad (20)$$

¹²For measurable spaces $(\mathcal{X}, \mathcal{X})$ and $(\mathcal{Y}, \mathcal{Y})$ a probability kernel $P(\cdot | \cdot)$ from $(\mathcal{X}, \mathcal{X})$ to $(\mathcal{Y}, \mathcal{Y})$ has the property that (i) For all $x \in \mathcal{X}$, $P(\cdot | x)$ is a probability measure, (ii) for all $Y \in \mathcal{Y}$, $x \mapsto P(Y | x)$ is measurable.

655 whenever $\mathbb{P} \ll \mathbb{Q}$. More generally, defining $p = \frac{d\mathbb{P}}{d\nu}$ and $q = \frac{d\mathbb{Q}}{d\nu}$ for a common dominating measure
 656 ν , we have

$$D_f(\mathbb{P} \parallel \mathbb{Q}) := \int_{q>0} qf\left(\frac{p}{q}\right) d\nu + \mathbb{P}(q=0) \cdot f'(\infty), \quad (21)$$

657 where $f'(\infty) := \lim_{x \rightarrow 0^+} xf(1/x)$.

658 C Technical Tools

659 C.1 Tail Bounds

660 **Lemma C.1** (e.g., Lemma A.4 of Foster et al. [18]). *For any sequence of real-valued random*
 661 *variables $(X_t)_{t \leq T}$ adapted to a filtration $(\mathcal{F}_t)_{t \leq T}$, we have that with probability at least $1 - \delta$,*

$$\sum_{t=1}^T X_t \leq \sum_{t=1}^T \log(\mathbb{E}[e^{X_t} \mid \mathcal{F}_{t-1}]) + \log(\delta^{-1}). \quad (22)$$

662 C.2 Minimax Theorem

663 **Lemma C.2** (Sion's Minimax Theorem [55]). *Let \mathcal{X} and \mathcal{Y} be convex sets in linear topological*
 664 *spaces, and assume \mathcal{X} is compact. Let $F : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ be such that (i) $F(x, \cdot)$ is concave and upper*
 665 *semicontinuous over \mathcal{Y} for all $x \in \mathcal{X}$ and (ii) $F(\cdot, y)$ is convex and lower semicontinuous over \mathcal{X} for*
 666 *all $y \in \mathcal{Y}$. Then*

$$\inf_{x \in \mathcal{X}} \sup_{y \in \mathcal{Y}} F(x, y) = \sup_{y \in \mathcal{Y}} \inf_{x \in \mathcal{X}} F(x, y). \quad (23)$$

667 C.3 Information Theory

668 C.3.1 Basic Results

669 **Proposition C.1.** *For any f -divergence $D_f(\cdot \parallel \cdot)$, one has that for any pair of random variables*
 670 *(X, Y) with joint law $\mathbb{P}_{X,Y}$,*

$$\mathbb{E}_{X \sim \mathbb{P}_X} [D_f(\mathbb{P}_{Y|X} \parallel \mathbb{P}_Y)] = \mathbb{E}_{Y \sim \mathbb{P}_Y} [D_f(\mathbb{P}_{X|Y} \parallel \mathbb{P}_X)].$$

671 **Proof of Proposition C.1.** Recalling that $D_f(\mathbb{P} \parallel \mathbb{Q}) = \mathbb{E}_{\mathbb{Q}} \left[f\left(\frac{d\mathbb{P}}{d\mathbb{Q}}\right) \right]$ for $\mathbb{P} \ll \mathbb{Q}$, we have

$$\begin{aligned} \mathbb{E}_{X \sim \mathbb{P}_X} [D_f(\mathbb{P}_{Y|X} \parallel \mathbb{P}_Y)] &= \mathbb{E}_{X \sim \mathbb{P}_X} \mathbb{E}_{Y \sim \mathbb{P}_Y} \left[f\left(\frac{d\mathbb{P}_{Y|X}}{d\mathbb{P}_Y}\right) \right] \\ &= \mathbb{E}_{X \sim \mathbb{P}_X} \mathbb{E}_{Y \sim \mathbb{P}_Y} \left[f\left(\frac{d\mathbb{P}_{X,Y}}{d(\mathbb{P}_X \otimes \mathbb{P}_Y)}\right) \right] \\ &= \mathbb{E}_{Y \sim \mathbb{P}_Y} \mathbb{E}_{X \sim \mathbb{P}_X} \left[f\left(\frac{d\mathbb{P}_{X|Y}}{d\mathbb{P}_X}\right) \right] = \mathbb{E}_{Y \sim \mathbb{P}_Y} [D_f(\mathbb{P}_{X|Y} \parallel \mathbb{P}_X)], \end{aligned}$$

672 where we have used that $\mathbb{P}_{Y|X} \ll \mathbb{P}_Y$, $\mathbb{P}_{X|Y} \ll \mathbb{P}_X$, and $\mathbb{P}_{X,Y} \ll \mathbb{P}_X \otimes \mathbb{P}_Y$. \square

673 C.3.2 Change of Measure

674 **Lemma C.3** (Donsker-Varadhan (e.g., Polyanskiy and Wu [48])). *Let \mathbb{P} and \mathbb{Q} be probability*
 675 *measures on $(\mathcal{X}, \mathcal{F})$. Then*

$$D_{\text{KL}}(\mathbb{P} \parallel \mathbb{Q}) = \sup_{h: \mathcal{X} \rightarrow \mathbb{R}} \{\mathbb{E}_{\mathbb{P}}[h(X)] - \log(\mathbb{E}_{\mathbb{Q}}[\exp(h(X))])\}. \quad (24)$$

676 **Lemma C.4.** *Let \mathbb{P} and \mathbb{Q} be probability distributions over a measurable space $(\mathcal{X}, \mathcal{F})$. Then for*
 677 *all functions $h : \mathcal{X} \rightarrow \mathbb{R}$,*

$$|\mathbb{E}_{\mathbb{P}}[h(X)] - \mathbb{E}_{\mathbb{Q}}[h(X)]| \leq \sqrt{2^{-1}(\mathbb{E}_{\mathbb{P}}[h^2(X)] + \mathbb{E}_{\mathbb{Q}}[h^2(X)]) \cdot D_{\text{H}}^2(\mathbb{P}, \mathbb{Q})}. \quad (25)$$

678 **Proof of Lemma C.4.** From Polyanskiy and Wu [48], we have that for all functions $h : \mathcal{X} \rightarrow \mathbb{R}$, if
679 $\mathbb{P} \ll \mathbb{Q}$,

$$|\mathbb{E}_{\mathbb{P}}[h(X)] - \mathbb{E}_{\mathbb{Q}}[h(X)]| \leq \sqrt{\mathbb{V}_{\mathbb{Q}}[h(X)] \cdot D_{\chi^2}(\mathbb{P} \parallel \mathbb{Q})} \leq \sqrt{\mathbb{E}_{\mathbb{Q}}[h^2(X)] \cdot D_{\chi^2}(\mathbb{P} \parallel \mathbb{Q})}, \quad (26)$$

680 where $D_{\chi^2}(\mathbb{P} \parallel \mathbb{Q}) := \int \frac{(d\mathbb{P} - d\mathbb{Q})^2}{d\mathbb{Q}}$ and $\mathbb{V}_{\mathbb{Q}}$ denotes the variance under \mathbb{Q} . The result follows by using
681 that $D_{\chi^2}(\mathbb{P} \parallel \frac{\mathbb{P} + \mathbb{Q}}{2}) \leq D_{\mathbb{H}}^2(\mathbb{P}, \mathbb{Q})$. \square

682 **Lemma 2.1.** Let \mathbb{P} and \mathbb{Q} be probability distributions over a measurable space $(\mathcal{X}, \mathcal{F})$. Then

$$\frac{1}{2} D_{\mathbb{H}}^2(\mathbb{P}, \mathbb{Q}) \leq \sup_{g: \mathcal{X} \rightarrow \mathbb{R}} \{1 - \mathbb{E}_{\mathbb{P}}[e^g] \cdot \mathbb{E}_{\mathbb{Q}}[e^{-g}]\} \leq D_{\mathbb{H}}^2(\mathbb{P}, \mathbb{Q}). \quad (12)$$

683 **Proof of Lemma 2.1.** We first show that Hellinger distance is lower bounded by the quantity in (12).
684 Recall that Hellinger distance is the f -divergence associated with $f(x) = (1 - \sqrt{x})^2$ (cf. (21)). Let
685 $f^*(y) := \sup_{x \geq 0} \{xy - f(x)\}$ be the Fenchel dual of f , which has the form

$$f^*(y) = \begin{cases} \frac{y}{1-y}, & y < 1, \\ \infty, & y \geq 1. \end{cases}$$

687 Using Theorem 7.14 of Polyanskiy [47], we express Hellinger distance as a following variational
688 problem based on the dual:

$$D_{\mathbb{H}}^2(\mathbb{P}, \mathbb{Q}) = \sup_{h: \mathcal{X} \rightarrow (-\infty, 1)} \{\mathbb{E}_{\mathbb{P}}[h(X)] - \mathbb{E}_{\mathbb{Q}}[f^*(h(X))]\} = \sup_{h: \mathcal{X} \rightarrow (-\infty, 1)} \left\{ \mathbb{E}_{\mathbb{P}}[h(X)] - \mathbb{E}_{\mathbb{Q}} \left[\frac{h(X)}{1-h(X)} \right] \right\}.$$

689 Reparameterizing via $h(X) = 1 - h'(X)$ for $h' : \mathcal{X} \rightarrow (0, \infty)$, this gives

$$D_{\mathbb{H}}^2(\mathbb{P}, \mathbb{Q}) = \sup_{h: \mathcal{X} \rightarrow (0, \infty)} \left\{ 2 - \mathbb{E}_{\mathbb{P}}[h(X)] - \mathbb{E}_{\mathbb{Q}} \left[\frac{1}{h(X)} \right] \right\}.$$

690 To conclude, we observe that for any test function $g : \mathcal{X} \rightarrow \mathbb{R}$, by setting $h(x) = e^{g(x)} \cdot \mathbb{E}_{\mathbb{Q}}[e^{-g}]$, we
691 have

$$\begin{aligned} 2 - \mathbb{E}_{\mathbb{P}}[h(X)] - \mathbb{E}_{\mathbb{Q}} \left[\frac{1}{h(X)} \right] &= 2 - \mathbb{E}_{\mathbb{P}}[e^g] \cdot \mathbb{E}_{\mathbb{Q}}[e^{-g}] - \mathbb{E}_{\mathbb{Q}}[e^{-g}] / \mathbb{E}_{\mathbb{Q}}[e^{-g}] \\ &= 1 - \mathbb{E}_{\mathbb{P}}[e^g] \cdot \mathbb{E}_{\mathbb{Q}}[e^{-g}], \end{aligned}$$

692 so that

$$D_{\mathbb{H}}^2(\mathbb{P}, \mathbb{Q}) \geq \sup_{g: \mathcal{X} \rightarrow \mathbb{R}} \{1 - \mathbb{E}_{\mathbb{P}}[e^g] \cdot \mathbb{E}_{\mathbb{Q}}[e^{-g}]\}.$$

693 We now prove the other direction of the inequality in (12). Let ν be a common dominating measure
694 for \mathbb{P} and \mathbb{Q} , and set $p = \frac{d\mathbb{P}}{d\nu}$ and $q = \frac{d\mathbb{Q}}{d\nu}$. We first consider the case where $p, q > 0$ everywhere.
695 Set $g(x) = \frac{1}{2} \log(q(x)/p(x))$. Then we have $\mathbb{E}_{\mathbb{P}}[e^g] = \int \sqrt{pq} d\nu = 1 - \frac{1}{2} D_{\mathbb{H}}^2(\mathbb{P}, \mathbb{Q})$, and likewise,
696 $\mathbb{E}_{\mathbb{Q}}[e^{-g}] = \int \sqrt{pq} d\nu = 1 - \frac{1}{2} D_{\mathbb{H}}^2(\mathbb{P}, \mathbb{Q})$. As a result,

$$\sup_{g: \mathcal{X} \rightarrow \mathbb{R}} \{1 - \mathbb{E}_{\mathbb{P}}[e^g] \cdot \mathbb{E}_{\mathbb{Q}}[e^{-g}]\} \geq 1 - (1 - \frac{1}{2} D_{\mathbb{H}}^2(\mathbb{P}, \mathbb{Q}))^2 \geq \frac{1}{2} D_{\mathbb{H}}^2(\mathbb{P}, \mathbb{Q}),$$

697 where we have used that $D_{\mathbb{H}}^2(\mathbb{P}, \mathbb{Q}) \in [0, 2]$. For the general case, one can appeal to Lemma C.5
698 below and take $\varepsilon \rightarrow 0$. \square

699 The following result is generalization of Lemma 2.1 which shows that up to small approximation
700 error, the lower bound in (12) can be obtained using test functions with small magnitude.

701 **Lemma C.5.** Let \mathbb{P} and \mathbb{Q} be probability distributions over a measurable space $(\mathcal{X}, \mathcal{F})$. Then for
702 any $\alpha \geq 1$, we have

$$\frac{1}{2} D_{\mathbb{H}}^2(\mathbb{P}, \mathbb{Q}) \leq \sup_{g \in \mathcal{G}_{\alpha}} \{1 - \mathbb{E}_{\mathbb{P}}[e^g] \cdot \mathbb{E}_{\mathbb{Q}}[e^{-g}]\} + 4e^{-\alpha}, \quad (27)$$

703 where $\mathcal{G}_{\alpha} := \{g : \mathcal{X} \rightarrow \mathbb{R} \mid \|g\|_{\infty} \leq \alpha\}$.

704 **Proof of Lemma C.5.** Fix $\alpha \geq 1$ and let $\varepsilon := e^{-2\alpha}$. Note that $\varepsilon \in (0, e^{-2})$. Given measures \mathbb{P} and
705 \mathbb{Q} , set $\mathbb{P}_\varepsilon = (1 - \varepsilon)\mathbb{P} + \varepsilon\mathbb{Q}$ and $\mathbb{Q}_\varepsilon = (1 - \varepsilon)\mathbb{Q} + \varepsilon\mathbb{P}$. Consider the test function $g = \frac{1}{2} \log(\frac{d\mathbb{Q}_\varepsilon}{d\mathbb{P}_\varepsilon})$,
706 which has the following properties:

- 707 • $\|g\|_\infty \leq \frac{1}{2} \log\left(\frac{1-\varepsilon}{\varepsilon} + \frac{\varepsilon}{1-\varepsilon}\right) \leq \frac{1}{2} \log(\varepsilon^{-1})$, where we have used that $\varepsilon \leq 1/2$. This
708 establishes that $g \in \mathcal{G}_\alpha$.
- 709 • $\mathbb{E}_\mathbb{P}[e^g] \leq (1 - \varepsilon)^{-1/2} \int \sqrt{d\mathbb{P}d\mathbb{Q}_\varepsilon} = (1 - \varepsilon)^{-1/2} (1 - \frac{1}{2}D_{\text{H}}^2(\mathbb{P}, \mathbb{Q}_\varepsilon))$.
- 710 • $\mathbb{E}_\mathbb{Q}[e^{-g}] \leq (1 - \varepsilon)^{-1/2} \int \sqrt{d\mathbb{P}_\varepsilon d\mathbb{Q}} = (1 - \varepsilon)^{-1/2} (1 - \frac{1}{2}D_{\text{H}}^2(\mathbb{P}_\varepsilon, \mathbb{Q}))$.

711 Using these bounds, we have

$$\begin{aligned} \sup_{g: \mathcal{X} \rightarrow \mathbb{R}} \{1 - \mathbb{E}_\mathbb{P}[e^g] \cdot \mathbb{E}_\mathbb{Q}[e^{-g}]\} &\geq 1 - (1 - \varepsilon)^{-1} (1 - \frac{1}{2}D_{\text{H}}^2(\mathbb{P}_\varepsilon, \mathbb{Q})) (1 - \frac{1}{2}D_{\text{H}}^2(\mathbb{P}, \mathbb{Q}_\varepsilon)) \\ &\geq 1 - (1 - \varepsilon)^{-1} (1 - \frac{1}{2}D_{\text{H}}^2(\mathbb{P}_\varepsilon, \mathbb{Q})) \\ &\geq (1 - \varepsilon)^{-1} \cdot \frac{1}{2}D_{\text{H}}^2(\mathbb{P}_\varepsilon, \mathbb{Q}) - 2\varepsilon. \end{aligned}$$

712 Finally, we note that by the triangle inequality for Hellinger distance and convexity of squared
713 Hellinger distance,

$$D_{\text{H}}(\mathbb{P}, \mathbb{Q}) \leq D_{\text{H}}(\mathbb{P}_\varepsilon, \mathbb{Q}) + D_{\text{H}}(\mathbb{P}, \mathbb{P}_\varepsilon) \leq D_{\text{H}}(\mathbb{P}_\varepsilon, \mathbb{Q}) + \varepsilon^{1/2}D_{\text{H}}(\mathbb{P}, \mathbb{Q}),$$

714 so that $D_{\text{H}}^2(\mathbb{P}_\varepsilon, \mathbb{Q}) \geq (1 - \varepsilon^{1/2})^2 D_{\text{H}}^2(\mathbb{P}, \mathbb{Q})$, and

$$\sup_{g: \mathcal{X} \rightarrow \mathbb{R}} \{1 - \mathbb{E}_\mathbb{P}[e^g] \cdot \mathbb{E}_\mathbb{Q}[e^{-g}]\} \geq \frac{(1 - \varepsilon^{1/2})^2}{1 - \varepsilon} \frac{1}{2}D_{\text{H}}^2(\mathbb{P}, \mathbb{Q}) - 2\varepsilon \geq \frac{1}{2}D_{\text{H}}^2(\mathbb{P}, \mathbb{Q}) - 4\varepsilon^{1/2},$$

715 where we have used that $\varepsilon \in (0, 1)$ and $D_{\text{H}}^2(\mathbb{P}, \mathbb{Q}) \in [0, 2]$.

716

□

717 C.4 Online Learning

718 **Lemma C.6** (e.g., Cesa-Bianchi and Lugosi [10]). *Let Π be a finite set. Consider the exponential
719 weights method with learning rate $\eta > 0$ and initial point $q^{(1)} = \text{unif}(\Pi)$, which has the update:*

$$q^{(t+1)}(\pi) = \frac{\exp(\eta \sum_{i \leq t} f^{(i)}(\pi))}{\sum_{\pi'} \exp(\eta \sum_{i \leq t} f^{(i)}(\pi'))},$$

720 *for an arbitrary (potentially adaptively selected) sequence of reward vectors $f^{(1)}, \dots, f^{(T)}$ in \mathbb{R}^Π .*
721 *This strategy ensures that with probability 1,*

$$\sum_{t=1}^T \langle q - q^{(t)}, f^{(t)} \rangle \leq \sum_{t=1}^T \langle q^{(t+1)} - q^{(t)}, f^{(t)} \rangle - \frac{1}{\eta} \sum_{t=1}^T D_{\text{KL}}(q^{(t+1)} \| q^{(t)}) + \frac{D_{\text{KL}}(q \| q^{(1)})}{\eta},$$

722 *for all $q \in \Delta(\Pi)$.*

723 D Examples

724 D.1 Structured Bandits

725 In this section we consider adversarial (structured) bandit problems, which correspond to the special
726 case of the adversarial DMSO setting in which there are no observations (i.e., $\mathcal{O} = \{\emptyset\}$). We
727 consider three examples: finite-armed bandits, linear bandits, and convex bandits. For each example,
728 we take $\mathcal{R} = [0, 1]$, fix a *reward function class* $\mathcal{F} \subseteq (\Pi \rightarrow [0, 1])$, and take $\mathcal{M}_\mathcal{F} = \{M \mid f^M \in \mathcal{F}\}$
729 to be the induced model class. Conceptually, $\mathcal{M}_\mathcal{F}$ should be thought of as the set of all reward
730 distributions over $[0, 1]$ with mean rewards in \mathcal{F} .

731 **Example D.1** (Finite-armed bandit). In the finite-armed bandit problem, we take $\Pi = \{1, \dots, A\}$ as
 732 the decision space, where $A \in \mathbb{N}$, then let $\mathcal{F} = [0, 1]^A$ and take $\mathcal{M} = \mathcal{M}_{\mathcal{F}}$ as the induced model
 733 class. For this setting, whenever $A \geq 2$, it holds that

$$\text{dec}_{\gamma}(\text{co}(\mathcal{M})) \leq \frac{A}{\gamma} \quad \forall \gamma > 0, \quad \text{and} \quad \text{dec}_{\gamma, \varepsilon_{\gamma}}(\text{co}(\mathcal{M})) \geq 2^{-6} \cdot \frac{A}{\gamma} \quad \forall \gamma \geq \frac{A}{3}, \quad (28)$$

734 where $\varepsilon_{\gamma} = \frac{A}{12\gamma}$. ◁

735 This result follows from Foster et al. [18, Proposition 5.2 and 5.3], noting that $\text{co}(\mathcal{M}) = \mathcal{M}$. Plugging
 736 (28) into Theorem 2.1 yields a $O(\sqrt{AT \log A})$ upper bound on regret, and plugging into Theorem 2.2
 737 gives a $\tilde{\Omega}(\sqrt{AT})$ lower bound for sub-Chebychev algorithms.¹³

738 **Example D.2** (Linear bandit). In the linear bandit problem, we have $\Pi \subseteq \mathbb{R}^d$. We take

$$\mathcal{F} = \{f : \Pi \rightarrow [0, 1] \mid f \text{ is linear}\},$$

739 and take $\mathcal{M} = \mathcal{M}_{\mathcal{F}}$ as the induced model class. For this setting, it holds that¹⁴

$$\text{dec}_{\gamma}(\text{co}(\mathcal{M})) \leq \frac{d}{4\gamma} \quad \forall \gamma > 0, \quad \text{and} \quad \text{dec}_{\gamma, \varepsilon_{\gamma}}(\text{co}(\mathcal{M})) \geq \frac{d}{12\gamma} \quad \forall \gamma \geq \frac{2d}{3}, \quad (29)$$

740 where $\varepsilon_{\gamma} := \frac{d}{3\gamma}$. ◁

741 This result follows from Foster et al. [18, Proposition 6.1 and 6.2], again noting that $\text{co}(\mathcal{M}) = \mathcal{M}$.
 742 Plugging (28) into Theorem 2.1 yields a $O(\sqrt{dT \log |\Pi|})$ upper bound on regret, and plugging into
 743 Theorem 2.2 gives a $\tilde{\Omega}(\sqrt{dT})$ lower bound for sub-Chebychev algorithms.

744 **Example D.3** (Convex bandit). In the convex bandit problem, we have $\Pi \subseteq \mathbb{R}^d$. We take

$$\mathcal{F} = \{f : \Pi \rightarrow [0, 1] \mid f \text{ is convex}\},$$

745 and take $\mathcal{M} = \mathcal{M}_{\mathcal{F}}$ as the induced model class. For this setting, it holds that for all $\gamma > 0$,

$$\text{dec}_{\gamma}(\text{co}(\mathcal{M})) \leq O\left(\frac{d^4}{\gamma} \cdot \text{polylog}(d, \text{diam}(\Pi), \gamma)\right). \quad (30)$$

746 ◁

747 This result follows from Foster et al. [18, Proposition 6.3] (which itself is a restatement of Lattimore
 748 and Szepesvári [36, Theorem 3]), and by noting once more that $\text{co}(\mathcal{M}) = \mathcal{M}$.

749 **Remark D.1.** The adversarial bandit literature [5, 4, 20, 11, 1, 8, 27, 17, 9, 31, 27, 7] typically
 750 considers a slightly different formulation in which the adversary selects a deterministic reward
 751 function. This can be captured by restricting \mathcal{M} to deterministic models. It is clear that the upper
 752 bounds on $\text{dec}_{\gamma}(\text{co}(\mathcal{M}))$ in the examples above lead to upper bounds for this model. The lower
 753 bounds in Examples D.1 and D.2 easily extend as well.

754 D.2 Reinforcement Learning

755 We now consider examples in reinforcement learning. We begin by recalling how to view the episodic
 756 reinforcement learning problem under the DMSO framework.

757 **Model class.** For episodic reinforcement learning, we fix a *horizon* H and let the model class
 758 \mathcal{M} consist of a set of non-stationary Markov Decision Processes (MDP). Each model $M \in \mathcal{M}$ is
 759 specified by

$$M = \{\{\mathcal{S}_h\}_{h=1}^{H+1}, \mathcal{A}, \{P_h^M\}_{h=1}^H, \{R_h^M\}_{h=1}^H, d_1\},$$

760 where \mathcal{S}_h is the state space for layer h , \mathcal{A} is the action space, $P_h^M : \mathcal{S}_h \times \mathcal{A} \mapsto \Delta(\mathcal{S}_{h+1})$ is the
 761 probability transition kernel for layer h , $R_h^M : \mathcal{S}_h \times \mathcal{A} \mapsto \Delta([0, 1])$ is the reward distribution for
 762 layer h and $d_1 \in \Delta(\mathcal{S}_1)$ is the initial state distribution. This formulation allows reward distribution

¹³For this example and Example D.2, the lower bound on $\text{dec}_{\gamma, \varepsilon_{\gamma}}(\text{co}(\mathcal{M}))$ in Foster et al. [18] is witnessed
 by a subfamily $\mathcal{M}' \subseteq \mathcal{M}$ with $V(\mathcal{M}') = O(1)$. As a result, we can take $C(T) = O(1)$ in Theorem 2.2.

¹⁴The upper bound here holds for all Π , while the lower bound holds for a specific choice for Π .

763 and transition kernel to vary across models in \mathcal{M} , but keeps the initial state distribution is fixed. We
 764 adopt the convention that $\bar{S}_{H+1} = \{s_{H+1}\}$ where s_{H+1} is a deterministic terminal state.

765 Before an episode, the learner selects a non-stationary policy, $\pi = (\pi_1, \dots, \pi_H)$ where $\pi_h : \mathcal{S}_h \mapsto \mathcal{A}$;
 766 we let Π_{NS} denote the set of all such policies. For a given MDP $M \in \mathcal{M}$, an episode proceeds by
 767 first sampling $s_1 \sim d_1$, then for $h = 1, \dots, H$:

- 768 • $a_h = \pi_h(s_h)$.
- 769 • $r_h \sim R_h^M(s_h, a_h)$ and $s_{h+1} \sim P_h^M(\cdot | s_h, a_h)$.

770 The value of the policy π under M is given by $f^M(\pi) := \mathbb{E}^{M, \pi}[\sum_{h=1}^H r_h]$, where $\mathbb{E}^{M, \pi}[\cdot]$ denotes
 771 expectation under the process above.

772 **Adversarial protocol.** Within the adversarial DMSO framework, model classes above lead to the
 773 following adversarial reinforcement learning protocol. At each time t , the learner plays selects a
 774 policy $\pi \in \Pi_{\text{NS}}$ and the adversary chooses an MDP $M^{(t)} \in \mathcal{M}$. The policy $\pi^{(t)}$ is then executed in
 775 the MDP $M^{(t)}$, resulting in a trajectory $\tau^{(t)} = (s_1^{(t)}, r_1^{(t)}, r_1^{(t)}), \dots, (s_H^{(t)}, r_H^{(t)}, r_H^{(t)})$. The learner then
 776 observes feedback $(r^{(t)}, o^{(t)})$, where $r^{(t)} := \sum_{h=1}^H r_h^{(t)}$ is the cumulative reward of the episode, and
 777 $o^{(t)} = \tau^{(t)}$ is the trajectory.

778 With this setting in mind, we give our main example.

779 **Example D.4 (Tabular MDP).** Let \mathcal{M} be the class of finite-state/action (tabular) MDPs with horizon
 780 $H, S \geq 2$ states, $A \geq 2$ actions, and $\sum_{h=1}^H r_h \in [0, 1]$. Then, for any $\gamma \geq A^{\min\{S-1, H\}}/6$,

$$\text{dec}_{\gamma, \varepsilon_\gamma}(\text{co}(\mathcal{M})) \geq \frac{A^{\min\{S-1, H\}}}{24\gamma},$$

781 where $\varepsilon_\gamma := A^{\min\{S-1, H\}}/24\gamma$. ◁

782 Using this result with [Theorem 2.2](#) leads to a lower bound on regret that scales with $\Omega(A^{\min\{S-1, H\}})$,
 783 which recovers existing intractability results for this setting [[30](#), [41](#)]. Note that we have $\text{dec}_\gamma(\mathcal{M}) =$
 784 $\text{poly}(S, A, H)/\gamma$ for this setting [[18](#)], so this is a case where there is a separation between the
 785 stochastic and adversarial setting.

786 We briefly mention that the set $\text{co}(\mathcal{M})$ can be interpreted as the set of *latent MDPs* [[30](#)]. In the latent
 787 MDP setting, each model is a mixture of MDPs. At the beginning of each episode, the underlying
 788 MDP from the mixture (the identity is not observed), and then run the MDP for the duration of the
 789 episode. This setting is also known to be intractable.

790 D.3 Proofs for Examples

791 D.3.1 Preliminaries

792 Our lower bounds on the Decision-Estimation Coefficient involve a constructing hard sub-family of
 793 models. Recall the following definition from [[18](#)].

794 **Definition D.1** ((α, β, δ) -family). *A reference model $\bar{M} \in \mathcal{M}$ and collection $\{M_1, \dots, M_N\}$ with
 795 $N \geq 2$ are said to be an (α, β, δ) -family if the following properties hold:*

- 796 1. *Regret property.* There exist functions $u^M : \Pi \mapsto [0, 1]$, with $\sum_{M \in \mathcal{M}} u^M(\pi) \leq \frac{N}{2}$ for all π
 797 such that

$$f^M(\pi_M) - f^M(\pi) \geq \alpha \cdot (1 - u^M(\pi))$$

798 for all $M \in \mathcal{M}$.

- 799 2. *Information property.* There exist functions $v^M : \Pi \mapsto [0, 1]$, with $\sum_{M \in \mathcal{M}} v^M(\pi) \leq 1$ for
 800 all π , such that

$$D_{\mathbb{H}}^2(M(\pi), \bar{M}(\pi)) \leq \beta \cdot v^M(\pi) + \delta.$$

801 Any (α, β, δ) -family leads to a difficult decision making problem because a given decision can have
 802 low regret or large information gain on (roughly) one model in the family. This is formalized through
 803 the following lemma.

804 **Lemma D.1** (Lemma 5.1, [18]). Let $\mathcal{M} = \{M_1, \dots, M_N\}$ be an (α, β, δ) -family with respect to \bar{M} .
 805 Then, for all $\gamma \geq 0$,

$$\text{dec}_\gamma(\mathcal{M}, \bar{M}) \geq \frac{\alpha}{2} - \gamma \left(\frac{\beta}{N} + \delta \right).$$

806 The following technical lemma bounds Hellinger distance for Bernoulli distributions.

807 **Lemma D.2** (Lemma A.7, [18]). For any $\Delta \in (0, 1/2)$,

$$D_{\text{H}}^2 \left(\text{Ber} \left(\frac{1}{2} + \Delta \right), \text{Ber} \left(\frac{1}{2} \right) \right) \leq 3\Delta^2.$$

808 D.3.2 Proof for Example D.4 (Tabular MDP)

809 In this section, we prove the lower bound in [Example D.4](#). We first derive an intermediate result
 810 which gives a lower bound on the Decision-Estimation Coefficient when the model class \mathcal{M} consists
 811 of mixtures of K MDPs; this is equivalent to the subset of $\text{co}(\mathcal{M})$ where we restrict to support size
 812 K , as well as the so-called latent MDP setting [30].

813 **Lemma D.3.** Let $K \geq 1$ be given. Let \mathcal{M} be the class of mixtures of K MDPs with horizon H ,
 814 $S \geq 2$ states, $A \geq 2$ actions, and $\sum_{h=1}^H r_h \in [0, 1]$. Then there exists $\bar{M} \in \mathcal{M}$ such that for all
 815 $\gamma \geq A^{\min\{S-1, H, K\}}/6$,

$$\text{dec}_\gamma(\mathcal{M}_{\varepsilon_\gamma}(\bar{M}), \bar{M}) \geq \frac{A^{\min\{S-1, H, K\}}}{24\gamma},$$

816 where $\varepsilon_\gamma := \frac{A^{\min\{S-1, H, K\}}}{24\gamma}$.

817 The proof of this result proceeds by constructing a hard sub-family of models and appealing to
 818 [Lemma D.1](#). Our construction is based of the lower bound for latent MDPs in Kwon et al. [30].

819 **Proof of Lemma D.3.** Let \mathcal{S} and \mathcal{A} be arbitrary sets with $|\mathcal{S}| = S$ and $|\mathcal{A}| = A$. Let $\Delta \in (0, 1/2)$
 820 be a parameter to be chosen later, and define $\bar{K} := \min\{S-1, K, H\}$. Partition the state space
 821 \mathcal{S} into sets \mathcal{S}' and $\mathcal{S} \setminus \mathcal{S}'$ such that $|\mathcal{S}'| = \bar{K} + 1$, and label the states in \mathcal{S}' as $\{s^{(1)}, \dots, s^{(\bar{K}+1)}\}$.
 822 Additionally, define sets via $\mathcal{S}_h = \{s^{(h)}, s^{(\bar{K}+1)}\}$ for $h \leq \bar{K}$ and $\mathcal{S}_h = \{s^{(\bar{K}+1)}\} \cup (\mathcal{S} \setminus \mathcal{S}')$ for
 823 $\bar{K} < h \leq H+1$. Recall that the decision space Π_{NS} is the set of all deterministic non-stationary
 824 policies $\pi = (\pi_1, \dots, \pi_H)$ where $\pi_h : \mathcal{S}_h \mapsto \mathcal{A}$.

825 We construct a class $\mathcal{M}' \subseteq \mathcal{M}$ in which each model $M \in \mathcal{M}'$ is specified by

$$M = \left\{ \{\mathcal{S}_h\}_{h=1}^{H+1}, \mathcal{A}, \{\mathbb{M}_k^M\}_{k=1}^{\bar{K}}, \{a_k^M\}_{k=1}^K \right\},$$

826 where for each $k \in [\bar{K}]$, $a_k^M \in \mathcal{A}$, and where \mathbb{M}_k^M is a tabular MDP specified by

$$\mathbb{M}_k^M = \left\{ \{\mathcal{S}_h\}_{h=1}^{H+1}, \mathcal{A}, \{P_{h,k}^M\}_{h=1}^H, \{R_{h,k}^M\}_{h=1}^H, \delta_{s^{(1)}} \right\}.$$

827 Here, $d_1 = \delta_{s^{(1)}}$, so that the initial state s_1 is $s^{(1)}$ deterministically. The transitions $P_{h,k}^M$ and rewards
 828 $R_{h,k}^M$ are constructed as follows.

829 • Construction of \mathbb{M}_1^M .

830 (i) For all $h \leq H$, the dynamics $P_{h,k}^M$ are deterministic. For an action a_h in the state s_h , the
 831 next state s_{h+1} is

$$s_{h+1} = \begin{cases} s^{(h+1)}, & \text{if } h \leq \bar{K}, s_h = s^{(h)}, \text{ and } a_h = a_i^M, \\ s^{(\bar{K}+1)}, & \text{if } h \leq \bar{K}, s_h = s^{(h)}, \text{ and } a_h \neq a_i^M, \\ s_h, & \text{otherwise.} \end{cases}$$

832 (ii) The reward distribution is given by

$$R_{h,k}^M(s_h, a_h) = \begin{cases} \text{Ber}(\frac{1}{2} + \Delta), & \text{if } h = \bar{K}, s_h = s^{(\bar{K})}, \text{ and } a_h = a_{\bar{K}}^M, \\ \text{Ber}(\frac{1}{2}), & \text{if } h = \bar{K}, s_h = s^{(\bar{K})}, \text{ and } a_h \neq a_{\bar{K}}^M, \\ 0, & \text{otherwise.} \end{cases}$$

833

- Construction of \mathbb{M}_j^M for $2 \leq j \leq \bar{K}$.

834

- (i) For each $h \leq H$, the dynamics $P_{h,k}^M$ are deterministic. For action a_h in state s_h , the next state s_{h+1} is

835

$$s_{h+1} = \begin{cases} s^{(h+1)} & \text{if } s_h = s^{(h)} \text{ and } h < j \\ s^{(\bar{K}+1)} & \text{if } s_h = s^{(h)}, h = j \text{ and } a_h = a_h^M \\ s^{(h+1)} & \text{if } s_h = s^{(h)}, h = j \text{ and } a_h \neq a_h^M \\ s^{(h+1)} & \text{if } s_h = s^{(h)}, h > j \text{ and } a_h = a_h^M \\ s^{(\bar{K}+1)} & \text{if } s_h = s^{(h)}, h > j \text{ and } a_h \neq a_h^M \\ s^{(\bar{K}+1)} & \text{if } h = \bar{K} - 1 \text{ or } h = \bar{K} \\ s_h & \text{otherwise} \end{cases}.$$

836

- (ii) The reward distribution is given by

$$R_{h,k}^M(s_h, a_h) = \begin{cases} \text{Ber}(\frac{1}{2}), & \text{if } h = \bar{K}, \\ 0, & \text{otherwise.} \end{cases}$$

837

Each model $M \in \mathcal{M}'$ is a uniform mixture of \bar{K} MDPs $\{\mathbb{M}_1^M, \dots, \mathbb{M}_{\bar{K}}^M\}$ as described above,

838

parameterized by the action sequence $a_{1:\bar{K}}^M$. The model class \mathcal{M}' is defined as the set of all such

839

mixture models (one for each sequence in $\mathcal{A}^{\bar{K}}$, so that $|\mathcal{M}'| = A^{\bar{K}}$).

840

At the start of each episode, an MDP \mathbb{M}_z^M is chosen by sampling $z \sim \text{Unif}([\bar{K}])$. The trajectory is then drawn by setting $s_1 = s^{(1)}$, and for $h = 1, \dots, H$:

841

842

- $a_h = \pi_h(s_h)$.

843

- $r_h \sim R_{h,z}^M(s_h, a_h)$ and $s_{h+1} \sim P_{h,z}^M(\cdot | s_h, a_h)$.

844

Note that rewards can be non-zero only at layer $h = \bar{K}$. We receive a reward from $\text{Ber}(\frac{1}{2} + \Delta)$ only

845

when $z = 1$ and the first \bar{K} actions match $a_{1:\bar{K}}^M$, i.e. $a_{1:\bar{K}} = a_{1:\bar{K}}^M$. For every other action sequence,

846

the reward is sampled from $\text{Ber}(\frac{1}{2})$. Thus, for any policy π ,

$$f^M(\pi) = \frac{1}{2} + \Delta \mathbb{I}\{\pi(s_{1:\bar{K}}) = a_{1:\bar{K}}^M\},$$

847

which implies that

$$f^M(\pi_M) - f^M(\pi) = \Delta(1 - \mathbb{I}\{\pi(s_{1:\bar{K}}) = a_{1:\bar{K}}^M\}). \quad (31)$$

848

Finally, we define the reference model \bar{M} . The model \bar{M} is specified by $\{\{\mathcal{S}_h\}_{h=1}^{H+1}, \mathcal{A}, \mathbb{M}^{\bar{M}}\}$ where

849

$\mathbb{M}^{\bar{M}}$ is a tabular MDP given by

$$\mathbb{M}^{\bar{M}} = \{\{\mathcal{S}_h\}_{h=1}^{H+1}, \mathcal{A}, P_h^{\bar{M}}, R_h^{\bar{M}}, \delta_{s^{(1)}}\}.$$

850

Here, the initial state s_1 is $s^{(1)}$ deterministically, and the transitions $P_{h,k}^{\bar{M}}$ and rewards $R_{h,k}^{\bar{M}}$ are as follows:

851

852

- (i) Transitions are stochastic and independent of the chosen action. In particular, for each $h \leq H$, the dynamics $P_h^{\bar{M}}$ are given by

853

$$P_h^{\bar{M}}(s_{h+1} | s_h, a_h) = \begin{cases} \frac{\bar{K}-h}{\bar{K}-h+1} & \text{if } h \leq \bar{K}, s_h = s^{(h)} \text{ and } s_{h+1} = s^{(h+1)} \\ \frac{1}{\bar{K}-h+1} & \text{if } h \leq \bar{K}, s_h = s^{(h)} \text{ and } s_{h+1} = s^{(\bar{K}+1)} \\ 1 & \text{if } h \leq \bar{K}, s_h \neq s^{(h)} \text{ and } s_h = s_{h+1} \\ 1 & \text{if } h > \bar{K} \text{ and } s_h = s_{h+1} \\ 0 & \text{otherwise} \end{cases}.$$

854

- (ii) The reward distribution is given by

$$R_h^{\bar{M}}(s_h, a_h) = \begin{cases} \text{Ber}(\frac{1}{2}), & \text{if } h = \bar{K}, \\ 0, & \text{otherwise.} \end{cases}$$

855 Note that \bar{M} can be thought of as a mixture of \bar{K} identical tabular MDPs each given by $\mathbb{M}^{\bar{M}}$. Note that
 856 for any policy π , the rewards for any trajectory in \bar{M} are sampled from $\text{Ber}(\frac{1}{2})$, and thus $f^{\bar{M}}(\pi) = \frac{1}{2}$
 857 which implies that

$$f^{\bar{M}}(\pi_{\bar{M}}) - f^{\bar{M}}(\pi) = 0. \quad (32)$$

858 We define $\mathcal{M}'' = \mathcal{M}' \cup \{\bar{M}\} \subseteq \mathcal{M}$, and note that for any policy π , the distribution over the
 859 trajectories is identical in all mixture models in \mathcal{M}'' . However, as mentioned before, the rewards in
 860 \bar{M} are sampled from $\text{Ber}(\frac{1}{2})$ and for any $M \in \mathcal{M}'$, the rewards in M are sampled from $\text{Ber}(\frac{1}{2} +$
 861 $\frac{\Delta}{\bar{K}} \mathbb{I}\{\pi(s_{1:\bar{K}}) = a_{1:\bar{K}}^M\})$. Thus, for any policy π and $M \in \mathcal{M}'$,

$$\begin{aligned} D_{\text{H}}^2(M(\pi), \bar{M}(\pi)) &= D_{\text{H}}^2\left(\text{Ber}\left(\frac{1}{2} + \frac{\Delta}{\bar{K}} \mathbb{I}\{\pi(s_{1:\bar{K}}) = a_{1:\bar{K}}^M\}\right), \text{Ber}\left(\frac{1}{2}\right)\right) \\ &\leq 3 \frac{\Delta^2}{\bar{K}^2} \cdot \mathbb{I}\{\pi(s_{1:\bar{K}}) = a_{1:\bar{K}}^M\}, \end{aligned} \quad (33)$$

862 where the last line uses [Lemma D.2](#).

863 The bounds in (31), (32) and (33) together imply that the model class \mathcal{M}'' is a $(\frac{\Delta}{\bar{K}}, 3\frac{\Delta^2}{\bar{K}^2}, 0)$ -family in
 864 the sense of [Definition D.1](#), where for each $\pi \in \Pi$ and $M \in \mathcal{M}''$ we take

$$u^M(\pi) := \mathbb{I}\{\pi(s_{1:\bar{K}}) = a_{1:\bar{K}}^M\} \quad \text{and} \quad v^M(\pi) := \mathbb{I}\{\pi(s_{1:\bar{K}}) = a_{1:\bar{K}}^M\},$$

865 with $u^{\bar{M}}(\pi) := 1$ and $v^{\bar{M}}(\pi) := 0$. As a result, [Lemma D.1](#) implies that

$$\text{dec}_{\gamma}(\mathcal{M}, \bar{M}) \geq \frac{\Delta}{2\bar{K}} - \frac{3\gamma\Delta^2}{\bar{K}^2 N},$$

866 for $N := A^{\bar{K}} + 1$. Setting $\Delta = \frac{\bar{K}N}{12\gamma}$ leads to the lower bound $\text{dec}_{\gamma}(\mathcal{M}, \bar{M}) \geq \frac{N}{24\gamma}$. We conclude by
 867 noting that all $M \in \mathcal{M}''$ have $M \in \mathcal{M}_{\varepsilon_{\gamma}}(\bar{M})$ with $\varepsilon_{\gamma} = \frac{N}{24\gamma}$, and thus the lower bound on the DEC
 868 also applies to the class $\mathcal{M}_{\varepsilon_{\gamma}}(\bar{M})$. \square

869 **Proof for Example D.4.** let \mathcal{M} be the class of all tabular MDPs, and let $\mathcal{M}^{(K)}$ denote the set of all
 870 mixture models in which each $M \in \mathcal{M}^{(K)}$ is a mixture of K MDPs from \mathcal{M} . Additionally, define
 871 $\widetilde{\mathcal{M}} = \text{co}(\mathcal{M})$, and note that $\mathcal{M}^{(K)} \subseteq \widetilde{\mathcal{M}}$ for all $K \geq 1$. For any $\varepsilon > 0$ and $\bar{M} \in \mathcal{M}^{(K)}$, we have
 872 that $\mathcal{M}_{\varepsilon}^{(K)}(\bar{M}) \subseteq \widetilde{\mathcal{M}}_{\varepsilon}(\bar{M})$, which implies that

$$\text{dec}_{\gamma}(\widetilde{\mathcal{M}}_{\varepsilon}(\bar{M}), \bar{M}) \geq \text{dec}_{\gamma}(\mathcal{M}_{\varepsilon}^{(K)}(\bar{M}), \bar{M}),$$

873 because $\text{dec}_{\gamma}(\cdot, \bar{M})$ is a non-decreasing function with respect to inclusion. Using [Lemma D.3](#), we
 874 have that for any $K \geq 1$ and $\gamma \geq A^{\min\{S-1, H, K\}}/6$, with $\varepsilon_{\gamma} := A^{\min\{S-1, H, K\}}/24\gamma$,

$$\text{dec}_{\gamma}(\widetilde{\mathcal{M}}_{\varepsilon}(\bar{M}), \bar{M}) \geq \text{dec}_{\gamma}(\mathcal{M}_{\varepsilon}^{(K)}(\bar{M}), \bar{M}) \geq \frac{A^{\min\{S-1, H, K\}}}{24\gamma}.$$

875 Setting $K = S$ above gives the desired lower bound. \square

876 E Structural Results

877 This section is organized as follows.

- 878 • In [Appendix E.1](#), we recall existing variants of the information ratio and state some basic
 879 properties.
- 880 • In [Appendix E.2](#), we prove equivalence of the Decision-Estimation Coefficient and the
 881 parameterized information ratio with Hellinger distance ([Theorem 3.1](#)), as well as a general-
 882 ization of this result ([Theorem E.1](#)).
- 883 • In [Appendix E.3](#), we prove equivalence of the parameterized information ratio with Hellinger
 884 distance and the high-probability exploration-by-optimization objective.

885 **E.1 Background on Complexity Measures**

886 For a measurable space $(\mathcal{X}, \mathcal{F})$, let us call any function $D : \Delta(\mathcal{X}) \times \Delta(\mathcal{X}) \rightarrow \mathbb{R}_+$ a *divergence-like*
887 function.

888 **Generalized information ratio.** Below we recall two notions of *generalized information ratio*
889 introduced by Lattimore and György [34] and Lattimore [33], which extend the original definition of
890 Russo and Van Roy [51, 52].

891 For a given prior $\mu \in \Delta(\mathcal{M} \times \Pi)$, define $\mu_{\text{pr}}(\pi') := \mathbb{P}(\pi^* = \pi')$ and $\mu_{\text{po}}(\pi'; \pi, z) := \mathbb{P}(\pi^* = \pi' \mid$
892 $(\pi, z))$ under the process $(M, \pi^*) \sim \mu, \pi \sim p, z \sim M(\pi)$.

893 1. Lattimore and György [34] define a class \mathcal{M} to have generalized information ratio (α, β, λ)
894 (where $\alpha, \beta \geq 0, \lambda > 1$) if for each prior $\mu \in \Delta(\mathcal{M} \times \Pi)$, there exists a distribution
895 $p \in \Delta(\Pi)$ such that

$$\mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)] \leq \alpha + \beta^{1-1/\lambda} \left(\mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D(\mu_{\text{po}}(\cdot; \pi, z) \parallel \mu_{\text{pr}})] \right)^{1/\lambda}. \quad (34)$$

896 2. Lattimore [33] define the generalized information ratio for a class \mathcal{M} (for $\lambda > 1$) via

$$\Psi_\lambda(\mathcal{M}) = \sup_{\mu \in \Delta(\mathcal{M} \times \Pi)} \inf_{p \in \Delta(\Pi)} \left\{ \frac{(\mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)])^\lambda}{\mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D(\mu_{\text{po}}(\cdot; \pi, z) \parallel \mu_{\text{pr}})]} \right\}. \quad (35)$$

897 As mentioned in Section 3, the formulations above slightly generalize the original versions in
898 Lattimore and György [34], Lattimore [33] by incorporating models $M \in \mathcal{M}$ and considering
899 general distances.

900 The following proposition shows that boundedness of the generalized information ratio implies
901 boundedness of the parameterized information ratio (Definition 3.1).

902 **Proposition E.1.** Fix $\alpha, \beta \geq 0$ and $\lambda > 1$. If a class \mathcal{M} has generalized information ratio (α, β, λ)
903 in the sense of (34), then

$$\inf_\gamma^D(\mathcal{M}) \leq \alpha + \frac{\beta}{\gamma^{\frac{1}{\lambda-1}}} \quad \forall \gamma > 0.$$

904 Likewise, the generalized information ratio in (35) satisfies

$$\inf_\gamma^D(\mathcal{M}) \leq (\Psi_\lambda(\mathcal{M})/\gamma)^{\frac{1}{\lambda-1}} \quad \forall \gamma > 0.$$

905 **Proof of Proposition E.1.** Suppose \mathcal{M} has generalized information ratio (α, β, λ) . Then there exists
906 $p \in \Delta(\Pi)$ such that for all $\mu \in \Delta(\mathcal{M} \times \Pi)$, we have

$$\begin{aligned} \mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)] &\leq \alpha + \beta^{1-1/\lambda} \left(\mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D(\mu_{\text{po}}(\cdot; \pi, z) \parallel \mu_{\text{pr}})] \right)^{1/\lambda} \\ &\leq \alpha + \frac{\beta}{\gamma^{\frac{1}{\lambda-1}}} + \gamma \cdot \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D(\mu_{\text{po}}(\cdot; \pi, z) \parallel \mu_{\text{pr}})], \end{aligned}$$

907 where we have applied Young's inequality, which gives that $xy \leq \frac{\lambda-1}{\lambda} x^{\frac{\lambda}{\lambda-1}} + \frac{1}{\lambda} y^\lambda$ for $x, y \geq 0$.

908 For the second result, we use that the definition of $\Psi_\lambda(\mathcal{M})$ implies generalized information ratio
909 $(0, (\Psi_\lambda(\mathcal{M}))^{\frac{1}{\lambda-1}}, \lambda)$. \square

910 This results show that an upper bound in terms of the parameterized information ratio in Definition 3.1
911 implies an upper bound in terms of either version of the generalized information ratio. It is also
912 straightforward to see that generalized information ratio $(0, \beta, \lambda)$ in (34) implies that $\Psi_\lambda(\mathcal{M}) \leq \beta^{\lambda-1}$
913 and vice-versa. Note that $\alpha = 0$ is the most interesting regime, as the regret bounds in Lattimore and
914 György [34] scale with $\alpha \cdot T$ when $\alpha > 0$.

915 Another important property of the parameterized information ratio (as well both generalized informa-
916 tion ratios) is that it is invariant under convexification.

917 **Proposition 3.1.** For any divergence-like function $D(\cdot \parallel \cdot) : \Delta(\Pi) \times \Delta(\Pi) \rightarrow \mathbb{R}_+$, we have

$$\inf_\gamma^D(\mathcal{M}) = \inf_\gamma^D(\text{co}(\mathcal{M})), \quad \forall \gamma > 0.$$

918 **Proof of Proposition 3.1.** Fix $\mu \in \Delta(\text{co}(\mathcal{M}) \times \Pi)$. We can represent any $\bar{M} \in \text{co}(\mathcal{M})$ as a mixture
919 $\nu \in \Delta(\mathcal{M})$, so that $\bar{M} = \mathbb{E}_{M \sim \nu}[M]$. Let $\tilde{\mu} \in \Delta(\Delta(\mathcal{M}) \times \Pi)$ be such that the process $(\nu, \pi^*) \sim \tilde{\mu}$,
920 $\bar{M} = \mathbb{E}_{M \sim \nu}[M]$ has the same law as $(\bar{M}, \pi^*) \sim \tilde{\mu}$. Finally, let $\mu' \in \Delta(\mathcal{M} \times \Pi)$ be the law of
921 (M, π^*) induced by sampling $(\nu, \pi^*) \sim \tilde{\mu}$ and $M \sim \nu$.

922 We observe that for any distribution $p \in \Delta(\Pi)$,

$$\begin{aligned} & \mathbb{E}_{(\bar{M}, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^{\bar{M}}(\pi^*) - f^{\bar{M}}(\pi)] \\ &= \mathbb{E}_{(\nu, \pi^*) \sim \tilde{\mu}} \mathbb{E}_{\pi \sim p} \mathbb{E}_{M \sim \nu} [f^M(\pi^*) - f^M(\pi)] \\ &= \mathbb{E}_{(M, \pi^*) \sim \mu'} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)]. \end{aligned}$$

923 Next, observe that (π, π^*, z) are identically distributed under the processes $\pi \sim p$, $(\bar{M}, \pi^*) \sim \mu$,
924 $z \sim \bar{M}(\pi)$ and $\pi \sim p$, $(M, \pi^*) \sim \mu'$, $z \sim M(\pi)$. As a result, we have $\mu_{\text{pr}} = \mu'_{\text{pr}}$ and $\mu_{\text{po}} = \mu'_{\text{po}}$, so

$$\mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D(\mu_{\text{po}}(\cdot; \pi, z) \parallel \mu_{\text{pr}})] = \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D(\mu'_{\text{po}}(\cdot; \pi, z) \parallel \mu'_{\text{pr}})].$$

925 This establishes that $\inf_{\gamma}^D(\text{co}(\mathcal{M})) \leq \inf_{\gamma}^D(\mathcal{M})$; the other direction is trivial. \square

926 E.2 Decision-Estimation Coefficient and Information Ratio (Theorem 3.1)

927 **Theorem 3.1.** For all $\gamma > 0$, $\inf_{\gamma}^H(\mathcal{M}) \leq \text{dec}_{\gamma}(\text{co}(\mathcal{M})) \leq \inf_{\gamma/4}^H(\mathcal{M})$.

928 **Theorem 3.1** is a special case of the following theorem, which concerns general divergence-like
929 functions.

930 **Theorem E.1.** Let $\Delta(\Pi) \times \Delta(\Pi) \rightarrow \mathbb{R}_+$ be any divergence-like function for which there exist
931 constants $c_1, c_2 \geq 1$ such that:

- 932 1. For all $\mathbb{Q} \in \Delta(\Pi)$, $\mathbb{P} \mapsto D(\mathbb{P} \parallel \mathbb{Q})$ is convex.
933 2. For all pairs of random variables (X, Y) ,

$$\mathbb{E}_{X \sim \mathbb{P}_X} [D(\mathbb{P}_{Y|X} \parallel \mathbb{P}_Y)] \leq c_1 \cdot \mathbb{E}_{Y \sim \mathbb{P}_Y} [D(\mathbb{P}_{X|Y} \parallel \mathbb{P}_X)]$$

- 934 3. For all pairs of random variables (X, Y) ,

$$\mathbb{E}_{X \sim \mathbb{P}_X} [D(\mathbb{P}_{Y|X} \parallel \mathbb{P}_Y)] \leq c_2 \cdot \inf_{\mathbb{Q}} \mathbb{E}_{X \sim \mathbb{P}_X} [D(\mathbb{P}_{Y|X} \parallel \mathbb{Q})].$$

- 935 4. For all $\varepsilon > 0$ sufficiently small, and all $\mathbb{Q} \in \Delta(\Pi)$, there exists $\mathbb{Q}' \in \Delta(\Pi)$ such that
936 $D(\mathbb{P} \parallel \mathbb{Q}) \geq D(\mathbb{P} \parallel \mathbb{Q}') - \varepsilon$ and $\sup_{\mathbb{P} \in \Delta(\Pi)} D(\mathbb{P} \parallel \mathbb{Q}') < \infty$.

937 Then we have

$$\inf_{c_1 \gamma}^D(\mathcal{M}) \leq \text{dec}_{\gamma}^D(\text{co}(\mathcal{M})) \leq \inf_{(c_1 c_2)^{-1} \gamma}^D(\mathcal{M}). \quad (36)$$

938 All f -divergences satisfy Property 2 with $c_1 = 1$, but may not satisfy Property 3. On the other hand,
939 Bregman divergences¹⁵ satisfy Property 3 with $c_2 = 1$, but may not satisfy Property 2 (consider
940 squared euclidean distance). KL-divergence, being both an f -divergence and a Bregman divergence,
941 satisfies both properties with $c_1 = c_2 = 1$ (this fact has been used tacitly in many prior works).
942 Squared Hellinger distance is an f -divergence but not a Bregman divergence, yet satisfies Property 3
943 with $c_2 = 4$ as a consequence of the triangle inequality.

944 **Proof of Theorem E.1.** We first bound the DEC by the information ratio, then proceed to bound the
945 information ratio by the DEC.

946 **Bounding the DEC by the information ratio.** Fix $M' \in \mathcal{M}$, and $\varepsilon > 0$ and let M'' be such that
947 $D_{\text{H}}^2(\cdot, M'(\pi)) \geq D_{\text{H}}^2(\cdot, M''(\pi)) - \varepsilon$ and $D_{\text{H}}^2(\cdot, M''(\pi)) < \infty$ (as guaranteed by Property 4). Using

¹⁵Recall that for a convex set \mathcal{X} and regularizer $\mathcal{R} : \mathcal{X} \rightarrow \mathbb{R}$, $D_{\mathcal{R}}(x \parallel y) := \mathcal{R}(x) - \mathcal{R}(y) - \langle \nabla \mathcal{R}(y), x - y \rangle$ is the associated Bregman divergence.

948 the minimax theorem (Lemma C.2), we have

$$\begin{aligned}
\text{dec}_\gamma^D(\mathcal{M}, M') &\leq \inf_{p \in \Delta(\Pi)} \sup_{M \in \mathcal{M}} \mathbb{E}_{\pi \sim p} [f^M(\pi_M) - f^M(\pi) - \gamma \cdot D(M(\pi) \| M''(\pi))] + \gamma\varepsilon \\
&= \inf_{p \in \Delta(\Pi)} \sup_{\nu \in \Delta(\mathcal{M})} \mathbb{E}_{\pi \sim p} \mathbb{E}_{M \sim \nu} [f^M(\pi_M) - f^M(\pi) - \gamma \cdot D(M(\pi) \| M''(\pi))] + \gamma\varepsilon \\
&= \sup_{\nu \in \Delta(\mathcal{M})} \inf_{p \in \Delta(\Pi)} \mathbb{E}_{\pi \sim p} \mathbb{E}_{M \sim \nu} [f^M(\pi_M) - f^M(\pi) - \gamma \cdot D(M(\pi) \| M''(\pi))] + \gamma\varepsilon.
\end{aligned}$$

949 Note that the application of the minimax theorem is admissible here, since $\Delta(\Pi)$ is compact (a
950 consequence of finiteness of Π) and the objective value is bounded (a consequence of the choice of
951 M'' and the fact that $f^M \in [0, 1]$).

952 Fix $\nu \in \Delta(\mathcal{M})$, and let $\mu \in \Delta(\mathcal{M} \times \Pi)$ be the induced law of (M, π_M) . Let $\bar{M}_{\pi'}(\pi) =$
953 $\mathbb{E}_{M \sim \nu} [M(\pi) | \pi_M = \pi']$ and $\bar{M}(\pi) = \mathbb{E}_{M \sim \mu} [M(\pi)] = \mathbb{E}_{\pi^* \sim \mu} [M_{\pi^*}(\pi)]$. Then for any $p \in \Delta(\Pi)$,
954 we have

$$\begin{aligned}
&\mathbb{E}_{M \sim \nu} \mathbb{E}_{\pi \sim p} [f^M(\pi_M) - f^M(\pi) - \gamma \cdot D(M(\pi) \| M''(\pi))] \\
&= \mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi) - \gamma \cdot D(M(\pi) \| M''(\pi))] \\
&\leq \mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi) - \gamma \cdot D(\bar{M}_{\pi^*}(\pi) \| M'; (\pi))] \\
&\leq \mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi) - \gamma c_2^{-1} \cdot D(\bar{M}_{\pi^*}(\pi) \| \bar{M}(\pi))],
\end{aligned}$$

955 where the first inequality uses convexity of $\mathbb{P} \mapsto D(\mathbb{P} \| \mathbb{Q})$ (Property 1), and the second inequality
956 uses Property 3. To proceed, let \mathbb{P} be the law of the process $\pi \sim p$, $(M, \pi^*) \sim \mu$, $z \sim M(\pi)$.
957 Observe that $\bar{M}_{\pi^*}(\pi) = \mathbb{P}_{z|\pi, \pi^*}$ and $\bar{M}(\pi) = \mathbb{P}_{z|\pi}$. Hence, using Property 2, we have that for all π ,

$$\mathbb{E}_{\pi^* \sim \mu} [D(\bar{M}_{\pi^*}(\pi) \| \bar{M}(\pi))] \geq c_1^{-1} \mathbb{E}_{z|\pi} [D(\mathbb{P}_{\pi^*|\pi, z} \| \mathbb{P}_{\pi^*|\pi})] = c_1^{-1} \mathbb{E}_{z|\pi} [D(\mathbb{P}_{\pi^*|\pi, z} \| \mathbb{P}_{\pi^*})],$$

958 where the last equality uses that π and π^* are independent (marginally). Since $D(\mathbb{P}_{\pi^*|\pi, z} \| \mathbb{P}_{\pi^*}) =$
959 $D(\mu_{\text{po}}(\cdot; \pi, z) \| \mu_{\text{pr}})$, if we choose p to attain the minimum in (19) for μ we are guaranteed that

$$\begin{aligned}
&\mathbb{E}_{M \sim \nu} \mathbb{E}_{\pi \sim p} [f^M(\pi_M) - f^M(\pi) - \gamma \cdot D(M(\pi) \| M'(\pi))] \\
&\leq \mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)] - \gamma (c_1 c_2)^{-1} \cdot \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D(\mu_{\text{po}}(\cdot; \pi, z) \| \mu_{\text{pr}})] + \gamma\varepsilon \\
&\leq \inf_{(c_1 c_2)^{-1} \gamma}^D(\mathcal{M}) + \gamma\varepsilon.
\end{aligned}$$

960 Taking $\varepsilon \rightarrow 0$, we conclude that $\text{dec}_\gamma^D(\mathcal{M}) \leq \inf_{(c_1 c_2)^{-1} \gamma}^D(\mathcal{M})$. By Proposition 3.1, $\inf_\gamma^D(\mathcal{M}) =$
961 $\inf_\gamma^D(\text{co}(\mathcal{M}))$, so applying the result to $\text{co}(\mathcal{M})$ yields

$$\text{dec}_\gamma^D(\text{co}(\mathcal{M})) \leq \inf_{(c_1 c_2)^{-1} \gamma}^D(\mathcal{M}).$$

962 **Bounding the information ratio by the DEC.** We now consider the opposite direction. Fix a prior
963 $\mu \in \Delta(\mathcal{M} \times \Pi)$ and consider the value for the parameterized information ratio:

$$\mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)] - \gamma \cdot \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D(\mu_{\text{po}}(\cdot; \pi, z) \| \mu_{\text{pr}})].$$

964 Define $\bar{M}_{\pi'}(\pi) := \mathbb{E}_\mu [M(\pi) | \pi^* = \pi']$ and $\bar{M}(\pi) = \mathbb{E}_{M \sim \mu} [M(\pi)]$. Using that (π^*, π) are inde-
965 pendent, along with Property 3, we have

$$\begin{aligned}
\mathbb{E}_{z|\pi} [D(\mu_{\text{po}}(\cdot; \pi, z) \| \mu_{\text{pr}})] &= \mathbb{E}_{z|\pi} [D(\mathbb{P}_{\pi^*|\pi, z} \| \mathbb{P}_{\pi^*})] \\
&= \mathbb{E}_{z|\pi} [D(\mathbb{P}_{\pi^*|\pi, z} \| \mathbb{P}_{\pi^*|\pi})] \geq c_1^{-1} \mathbb{E}_{\pi^* \sim \mu} [D(\bar{M}_{\pi^*}(\pi) \| \bar{M}(\pi))].
\end{aligned}$$

966 Next, observe that we have

$$\begin{aligned}
\mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)] &= \mathbb{E}_{\pi \sim p} \mathbb{E}_{\pi^* \sim \mu} \mathbb{E} [f^M(\pi^*) - f^M(\pi) | \pi^*] \\
&= \mathbb{E}_{\pi \sim p} \mathbb{E}_{\pi^* \sim \mu} [f^{\bar{M}_{\pi^*}}(\pi^*) - f^{\bar{M}_{\pi^*}}(\pi)] \\
&\leq \mathbb{E}_{\pi^* \sim \mu} \mathbb{E}_{\pi \sim p} \left[\max_{\pi'} f^{\bar{M}_{\pi^*}}(\pi') - f^{\bar{M}_{\pi^*}}(\pi) \right].
\end{aligned}$$

967 Recall that the definition of $\text{dec}_\gamma(\text{co}(\mathcal{M}))$ implies the following: For any $\kappa \in \Delta(\mathcal{M})$ there exists
 968 a distribution $p \in \Delta(\Pi)$ such that for all $\nu \in \Delta(\mathcal{M})$, defining $\bar{M}_\kappa(\pi) := \mathbb{E}_{M \sim \kappa}[M(\pi)]$ and
 969 $\bar{M}_\nu(\pi) := \mathbb{E}_{M \sim \nu}[M(\pi)]$, we have

$$\mathbb{E}_{\pi \sim p} \left[\max_{\pi'} f^{\bar{M}_\nu}(\pi') - f^{\bar{M}_\nu}(\pi) - \gamma \cdot D(\bar{M}_\nu(\pi) \parallel \bar{M}_\kappa(\pi)) \right] \leq \text{dec}_\gamma(\text{co}(\mathcal{M})). \quad (37)$$

970 By invoking (37) with $\bar{M}_\kappa = \bar{M}$ and $\bar{M}_\nu = \bar{M}_{\pi^*}$, we are guaranteed that for every draw of π^*

$$\mathbb{E}_{\pi \sim p} \left[\max_{\pi'} f^{\bar{M}_{\pi^*}}(\pi') - f^{\bar{M}_{\pi^*}}(\pi) \right] \leq \gamma c_1^{-1} \cdot \mathbb{E}_{\pi \sim p} [D(\bar{M}_{\pi^*}(\pi) \parallel \bar{M}(\pi))] + \text{dec}_{c_1^{-1}\gamma}(\text{co}(\mathcal{M})).$$

971 Taking the expectation over $\pi^* \sim \mu$, we conclude that

$$\inf_\gamma^D(\mathcal{M}) \leq \text{dec}_{c_1^{-1}\gamma}(\text{co}(\mathcal{M})).$$

972

□

973 E.3 High-Probability Exploration-By-Optimization and Information Ratio (Theorem 3.2)

974 **Theorem 3.2.** For all $\eta > 0$, $\inf_{\eta-1}^H(\mathcal{M}) \leq \text{exo}_\eta(\mathcal{M}) \leq \inf_{(8\eta)-1}^H(\mathcal{M})$.

975 **Proof of Theorem 3.2.** We first state the following basic result, which is proven in the sequel.

976 **Lemma E.1.** For any fixed $M \in \mathcal{M}$ and $\pi^* \in \Pi$, the map $(p, g) \mapsto \Gamma_{q,\eta}(p, g; \pi^*, M)$ is jointly
 977 convex with respect to $(p, g) \in \Delta(\Pi) \times \mathcal{G}$, where $\mathcal{G} := (\Pi \times \Pi \times \mathcal{Z} \rightarrow \mathbb{R})$.

978 **Upper bound: Minimax theorem.** We first use the minimax theorem to move to a Bayesian coun-
 979 terpart to the Exploration-by-Optimization objective. This requires some care to ensure boundedness
 980 and compactness, but otherwise is conceptually straightforward. To begin, observe that we can write
 981 the Exploration-by-Optimization objective as

$$\begin{aligned} \text{exo}_\eta(\mathcal{M}) &= \sup_{q \in \Delta(\Pi)} \inf_{p \in \Delta(\Pi), g \in \mathcal{G}} \sup_{M \in \mathcal{M}, \pi^* \in \Pi} [\Gamma_{q,\eta}(p, g; \pi^*, M)] \\ &= \sup_{q \in \Delta(\Pi)} \inf_{p \in \Delta(\Pi), g \in \mathcal{G}} \sup_{\mu \in \Delta(\mathcal{M} \times \Pi)} \mathbb{E}_{(M, \pi^*) \sim \mu} [\Gamma_{q,\eta}(p, g; \pi^*, M)]. \end{aligned}$$

982 Fix $\alpha \geq 1 \vee \eta^{-1}$ and $\varepsilon \in (0, 1)$, and define

$$\mathcal{G}_\alpha = \{g \in \mathcal{G} \mid \|g\|_\infty \leq \alpha\}, \quad \text{and} \quad \mathcal{P}_\varepsilon = \{p \in \Delta(\Pi) \mid p(\pi) \geq \varepsilon |\Pi|^{-1} \forall \pi\}.$$

983 Then, by restricting to these classes, we have¹⁶

$$\text{exo}_\eta(\mathcal{M}) \leq \sup_{q \in \Delta(\Pi)} \inf_{p \in \mathcal{P}_\varepsilon, g \in \mathcal{G}_\alpha} \sup_{\mu \in \Delta(\mathcal{M} \times \Pi)} \mathbb{E}_{(M, \pi^*) \sim \mu} [\Gamma_{q,\eta}(p, g; \pi^*, M)]$$

984 We verify that the conditions required to apply the minimax theorem are satisfied.

- 985 • The map $\mu \mapsto \mathbb{E}_{(M, \pi^*) \sim \mu} [\Gamma_{q,\eta}(p, g; \pi^*, M)]$ is linear. Furthermore, by Lemma E.1, the map
 986 $(p, g) \mapsto \mathbb{E}_{(M, \pi^*) \sim \mu} [\Gamma_{q,\eta}(p, g; \pi^*, M)]$ is convex.
- 987 • Since we have restricted to $p \in \mathcal{P}_\varepsilon$ and $g \in \mathcal{G}_\alpha$, the value $\Gamma_{q,\eta}(p, g; \pi^*, M)$ is uniformly bounded,
 988 as well as continuous with respect to p and g (so long as $\varepsilon > 0$ and $\alpha < \infty$).
- 989 • The set $\Delta(\mathcal{M} \times \Pi)$ is convex. Since $|\Pi| < \infty$, the set $\mathcal{P}_\varepsilon \times \mathcal{G}_\alpha$ is convex and compact (for \mathcal{P}_ε
 990 equipped with the usual topology and \mathcal{G}_α equipped with the product topology; see Lattimore and
 991 György [34] for details).

992 Hence, using Lemma C.2 we can bound by the value of the Bayesian game as follows:

$$\text{exo}_\eta(\mathcal{M}) \leq \sup_{q \in \Delta(\Pi)} \sup_{\mu \in \Delta(\mathcal{M} \times \Pi)} \inf_{p \in \mathcal{P}_\varepsilon, g \in \mathcal{G}_\alpha} \mathbb{E}_{(M, \pi^*) \sim \mu} [\Gamma_{q,\eta}(p, g; \pi^*, M)]. \quad (38)$$

¹⁶Restricting to these sets allows us to enforce boundedness and continuity of the Exploration-by-Optimization objective, which is necessary to appeal to the minimax theorem. The parameters α and ε will not enter the final bound quantitatively.

993 **Upper bound: Moving to Hellinger distance.** For any $q \in \Delta(\Pi)$, $\mu \in \Delta(\mathcal{M} \times \Pi)$, and $p \in \mathcal{P}_\varepsilon$
 994 the value of the game in (38) is

$$\begin{aligned} & \mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)] \\ & + \eta^{-1} \inf_{g \in \mathcal{G}_\alpha} \mathbb{E}_{(M, \pi^*) \sim \mu} \left[\mathbb{E}_{\pi \sim p, z \sim M(\pi)} \mathbb{E}_{\pi' \sim q} \exp\left(\frac{\eta}{p(\pi)} (g(\pi'; \pi, z) - g(\pi^*; \pi, z))\right) - 1 \right]. \end{aligned}$$

995 Using Bayes' rule, we can rewrite the second term above as

$$\inf_{g \in \mathcal{G}_\alpha} \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} \left[\mathbb{E}_{\pi' \sim q} \left[\exp\left(\eta \frac{g(\pi'; \pi, z)}{p(\pi)}\right) \right] \cdot \mathbb{E}_{\pi^* \sim \mu_{\text{po}}(\cdot; \pi, z)} \left[\exp\left(-\eta \frac{g(\pi^*; \pi, z)}{p(\pi)}\right) \right] - 1 \right]$$

996 By reparameterizing via $g(\pi'; \pi, z) \leftarrow \frac{p(\pi)}{\eta} g(\pi'; \pi, z)$, the value is upper bounded by

$$\inf_{g \in \mathcal{G}_{\alpha\eta}} \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} \left[\mathbb{E}_{\pi' \sim q} [\exp(g(\pi'; \pi, z))] \cdot \mathbb{E}_{\pi^* \sim \mu_{\text{po}}(\cdot; \pi, z)} [\exp(-g(\pi^*; \pi, z))] - 1 \right].$$

997 Furthermore, by skolemizing, we can rewrite this as

$$V(p, q, \mu) := \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} \inf_{g: \Pi \rightarrow \mathbb{R}, \|g\|_\infty \leq \alpha\eta} \left\{ \mathbb{E}_{\pi' \sim q} [\exp(g(\pi'))] \cdot \mathbb{E}_{\pi^* \sim \mu_{\text{po}}(\cdot; \pi, z)} [\exp(-g(\pi^*))] - 1 \right\}.$$

998 We now appeal to [Lemma C.5](#), which grants that

$$V(p, q, \mu) \leq -\frac{1}{2} \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D_{\text{H}}^2(\mu_{\text{po}}(\cdot; \pi, z), q)] + 4e^{-\alpha\eta}. \quad (39)$$

999 Using (39), we have

$$\begin{aligned} & \text{exo}_\eta(\mathcal{M}) \\ & \leq \sup_{q \in \Delta(\Pi)} \sup_{\mu \in \Delta(\mathcal{M} \times \Pi)} \inf_{p \in \mathcal{P}_\varepsilon} \left\{ \mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)] - \frac{1}{2\eta} \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D_{\text{H}}^2(\mu_{\text{po}}(\cdot; \pi, z), q)] \right\} + 4\eta^{-1} e^{-\alpha\eta}. \end{aligned}$$

1000 In addition, since $f^M \in [0, 1]$ and $D_{\text{H}}^2(\cdot, \cdot) \in [0, 2]$, we can further upper bound by

$$\begin{aligned} & \sup_{q \in \Delta(\Pi)} \sup_{\mu \in \Delta(\mathcal{M} \times \Pi)} \inf_{p \in \Delta(\Pi)} \left\{ \mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)] - \frac{1}{2\eta} \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D_{\text{H}}^2(\mu_{\text{po}}(\cdot; \pi, z), q)] \right\} \\ & + O(\eta^{-1} e^{-\alpha\eta} + \varepsilon \cdot (1 + \eta^{-1})). \end{aligned}$$

1001 Since this expression only depends on α and ε through the additive approximation terms, taking the
 1002 limit as $\alpha \rightarrow \infty$ and $\varepsilon \rightarrow 0$ yields

$$\text{exo}_\eta(\mathcal{M}) \leq \sup_{q \in \Delta(\Pi)} \sup_{\mu \in \Delta(\mathcal{M} \times \Pi)} \inf_{p \in \Delta(\Pi)} \left\{ \mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)] - \frac{1}{2\eta} \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D_{\text{H}}^2(\mu_{\text{po}}(\cdot; \pi, z), q)] \right\}.$$

1003 Finally, recall that since Hellinger distance satisfies the triangle inequality, we have

$$\mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D_{\text{H}}^2(\mu_{\text{po}}(\cdot; \pi, z), \mu_{\text{pr}})] \leq 2 \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D_{\text{H}}^2(\mu_{\text{po}}(\cdot; \pi, z), q)] + 2D_{\text{H}}^2(\mu_{\text{pr}}, q).$$

1004 Using that $\mu_{\text{pr}}(\pi') = \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [\mu_{\text{po}}(\pi'; \pi, z)]$ and that squared Hellinger distance is convex, we
 1005 have $D_{\text{H}}^2(\mu_{\text{pr}}, q) \leq \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D_{\text{H}}^2(\mu_{\text{po}}(\cdot; \pi, z), q)]$, and so

$$\mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D_{\text{H}}^2(\mu_{\text{po}}(\cdot; \pi, z), \mu_{\text{pr}})] \leq 4 \cdot \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D_{\text{H}}^2(\mu_{\text{po}}(\cdot; \pi, z), q)].$$

1006 It follows that

$$\begin{aligned} \text{exo}_\eta(\mathcal{M}) & \leq \sup_{\mu \in \Delta(\mathcal{M} \times \Pi)} \inf_{p \in \Delta(\Pi)} \left\{ \mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)] - \frac{1}{8\eta} \mathbb{E}_{\pi \sim p} \mathbb{E}_{z|\pi} [D_{\text{H}}^2(\mu_{\text{po}}(\cdot; \pi, z), \mu_{\text{pr}})] \right\} \\ & = \inf_{(8\eta)^{-1}}^{\text{H}}(\mathcal{M}). \end{aligned}$$

1007 **Lower bound.** It is immediate (without having to invoke the minimax theorem) that

$$\begin{aligned} \text{exo}_\eta(\mathcal{M}) & = \sup_{q \in \Delta(\Pi)} \inf_{p \in \Delta(\Pi), g \in \mathcal{G}} \sup_{\mu \in \Delta(\mathcal{M} \times \Pi)} \mathbb{E}_{(M, \pi^*) \sim \mu} [\Gamma_{q, \eta}(p, g; \pi^*, M)] \\ & \geq \sup_{q \in \Delta(\Pi)} \sup_{\mu \in \Delta(\mathcal{M} \times \Pi)} \inf_{p \in \Delta(\Pi), g \in \mathcal{G}} \mathbb{E}_{(M, \pi^*) \sim \mu} [\Gamma_{q, \eta}(p, g; \pi^*, M)]. \end{aligned}$$

1008 Performing the same sequence of calculations as in the upper bound, we have that for any $q \in \Delta(\Pi)$,
 1009 $\mu \in \Delta(\mathcal{M} \times \Pi)$, and $p \in \Delta(\Pi)$,

$$\begin{aligned} & \inf_{g \in \mathcal{G}} \mathbb{E}_{(M, \pi^*) \sim \mu} [\Gamma_{q, \eta}(p, g; \pi^*, M)] \\ &= \mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)] \\ & \quad + \eta^{-1} \inf_{g \in \mathcal{G}} \mathbb{E}_{(M, \pi^*) \sim \mu} \left[\mathbb{E}_{\pi \sim p, z \sim M(\pi)} \mathbb{E}_{\pi' \sim q} \exp\left(\frac{\eta}{p(\pi)}(g(\pi'; \pi, z) - g(\pi^*; \pi, z))\right) - 1 \right] \\ &= \mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)] + \eta^{-1} \mathbb{E}_{\pi \sim p} \mathbb{E}_{z | \pi} \inf_{g \in \mathcal{G}} \left\{ \mathbb{E}_{\pi' \sim q} [\exp(g(\pi'))] \cdot \mathbb{E}_{\pi^* \sim \mu_{\text{po}}(\cdot; \pi, z)} [\exp(-g(\pi^*))] - 1 \right\}. \end{aligned}$$

1010 Using Lemma 2.1, we have

$$\mathbb{E}_{\pi \sim p} \mathbb{E}_{z | \pi} \inf_{g \in \mathcal{G}} \left\{ \mathbb{E}_{\pi' \sim q} [\exp(g(\pi'))] \cdot \mathbb{E}_{\pi^* \sim \mu_{\text{po}}(\cdot; \pi, z)} [\exp(-g(\pi^*))] - 1 \right\} \geq -\mathbb{E}_{\pi \sim p} \mathbb{E}_{z | \pi} [D_{\text{H}}^2(\mu_{\text{po}}(\cdot; \pi, z), q)].$$

1011 We conclude that

$$\begin{aligned} \text{exo}_{\eta}(\mathcal{M}) &\geq \sup_{q \in \Delta(\Pi)} \sup_{\mu \in \Delta(\mathcal{M} \times \Pi)} \inf_{p \in \Delta(\Pi)} \left\{ \mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)] - \frac{1}{\eta} \mathbb{E}_{\pi \sim p} \mathbb{E}_{z | \pi} [D_{\text{H}}^2(\mu_{\text{po}}(\cdot; \pi, z), q)] \right\} \\ &\geq \sup_{\mu \in \Delta(\mathcal{M} \times \Pi)} \inf_{p \in \Delta(\Pi)} \left\{ \mathbb{E}_{(M, \pi^*) \sim \mu} \mathbb{E}_{\pi \sim p} [f^M(\pi^*) - f^M(\pi)] - \frac{1}{\eta} \mathbb{E}_{\pi \sim p} \mathbb{E}_{z | \pi} [D_{\text{H}}^2(\mu_{\text{po}}(\cdot; \pi, z), \mu_{\text{pr}})] \right\} \\ &= \inf_{\eta^{-1}}^{\text{H}}(\mathcal{M}). \end{aligned}$$

1012

□

1013 **Proof of Lemma E.1.** Let $M \in \mathcal{M}$ and $\pi^* \in \Pi$ be fixed. The map $p \mapsto \mathbb{E}_{\pi \sim p} [f^M(\pi_M) - f^M(\pi)]$
 1014 is linear, so our main task is to show that the function

$$(p, g) \mapsto \sum_{\pi} p(\pi) \mathbb{E}_{z \sim M(\pi)} \left[\sum_{\pi'} q(\pi') \exp\left(\frac{\eta}{p(\pi)}(g(\pi'; \pi, z) - g(\pi^*; \pi, z))\right) \right]$$

1015 is jointly convex. We can rewrite this as

$$\sum_{\pi} q(\pi') \sum_{\pi} p(\pi) \mathbb{E}_{z \sim M(\pi)} \left[\exp\left(\frac{\eta}{p(\pi)}(g(\pi'; \pi, z) - g(\pi^*; \pi, z))\right) \right].$$

1016 Since convexity is preserved under summation with non-negative weights, it suffices to show that for
 1017 any fixed (π, π') , the map

$$(p(\pi), g) \mapsto p(\pi) \mathbb{E}_{z \sim M(\pi)} \left[\exp\left(\frac{\eta}{p(\pi)}(g(\pi'; \pi, z) - g(\pi^*; \pi, z))\right) \right] \quad (40)$$

1018 is convex. Since the function $g \mapsto \mathbb{E}_{z \sim M(\pi)} [\exp(\eta(g(\pi'; \pi, z) - g(\pi^*; \pi, z)))]$ is convex over \mathcal{G} ,
 1019 convexity for (40) follows from the following standard result.

1020 **Proposition E.2** (Convexity of perspective transformation). *Let $f : \mathbb{R}^d \rightarrow (-\infty, \infty)$ be a convex*
 1021 *function. Then the function*

$$(x, t) \mapsto t \cdot f(x/t)$$

1022 *is convex over $\mathbb{R}^d \times \mathbb{R}_+$.*

1023

□

1024 F Proofs for Main Results (Section 2)

1025 F.1 Proof of Theorem 2.1

1026 **Theorem 2.1** (Main upper bound). *For any choice of $\eta > 0$, Algorithm 1 ensures that for all $\delta > 0$,*
 1027 *with probability at least $1 - \delta$,*

$$\mathbf{Reg}_{\text{DM}} \leq \text{dec}_{(8\eta)^{-1}}(\text{co}(\mathcal{M})) \cdot T + 2\eta^{-1} \cdot \log(|\Pi|/\delta). \quad (8)$$

1028 *In particular, for any $\delta > 0$, with appropriate η , the algorithm has that with probability at least $1 - \delta$,*

$$\mathbf{Reg}_{\text{DM}} \leq O(1) \cdot \inf_{\gamma > 0} \{ \text{dec}_{\gamma}(\text{co}(\mathcal{M})) \cdot T + \gamma \cdot \log(|\Pi|/\delta) \}. \quad (9)$$

1029 **Proof of Theorem 2.1.** Let us adopt convention $\langle p, f \rangle = \sum_{\pi} p(\pi) \cdot f(\pi)$ and let e_{π} denote the π th
 1030 standard basis vector in \mathbb{R}^{Π} . For each $\pi^* \in \Pi$, we write regret as

$$\mathbf{Reg}_{\text{DM}}(\pi^*) = \sum_{t=1}^T \mathbb{E}_{\pi \sim p^{(t)}} \left[f^{M^{(t)}}(\pi^*) - f^{M^{(t)}}(\pi) \right] = \sum_{t=1}^T \langle e_{\pi^*} - p^{(t)}, f^{M^{(t)}} \rangle.$$

1031 Adding and subtracting $\sum_{t=1}^T \langle e_{\pi^*} - q^{(t)}, \hat{f}^{(t)} \rangle$, we rewrite this as

$$\sum_{t=1}^T \langle e_{\pi^*} - p^{(t)}, f^{M^{(t)}} \rangle = \sum_{t=1}^T \langle e_{\pi^*} - p^{(t)}, f^{M^{(t)}} \rangle + \sum_{t=1}^T \langle e_{\pi^*} - q^{(t)}, \hat{f}^{(t)} \rangle - \sum_{t=1}^T \langle e_{\pi^*} - q^{(t)}, \hat{f}^{(t)} \rangle. \quad (41)$$

1032 The exponential weights update ensures (Lemma C.6) that with probability 1,

$$\begin{aligned} \sum_{t=1}^T \langle e_{\pi^*} - q^{(t)}, \hat{f}^{(t)} \rangle &\leq \sum_{t=1}^T \langle q^{(t+1)} - q^{(t)}, \hat{f}^{(t)} \rangle - \frac{1}{\eta} \sum_{t=1}^T D_{\text{KL}}(q^{(t+1)} \| q^{(t)}) + \frac{D_{\text{KL}}(e_{\pi^*} \| q^{(1)})}{\eta} \\ &\leq \sum_{t=1}^T \langle q^{(t+1)} - q^{(t)}, \hat{f}^{(t)} \rangle - \frac{1}{\eta} \sum_{t=1}^T D_{\text{KL}}(q^{(t+1)} \| q^{(t)}) + \frac{\log|\Pi|}{\eta}. \end{aligned}$$

1033 In addition, using Lemma C.3, we have that for all t ,

$$\langle q^{(t+1)}, \hat{f}^{(t)} \rangle - \frac{1}{\eta} D_{\text{KL}}(q^{(t+1)} \| q^{(t)}) \leq \frac{1}{\eta} \log \left(\sum_{\pi} q^{(t)}(\pi) \exp(\eta \cdot \hat{f}^{(t)}(\pi)) \right).$$

1034 Hence, combining this with (41), we have

$$\mathbf{Reg}_{\text{DM}}(\pi^*) \leq \sum_{t=1}^T \langle e_{\pi^*} - p^{(t)}, f^{M^{(t)}} \rangle - \langle e_{\pi^*}, \hat{f}^{(t)} \rangle + \frac{1}{\eta} \sum_{t=1}^T \log \left(\sum_{\pi} q^{(t)}(\pi) \exp(\eta \cdot \hat{f}^{(t)}(\pi)) \right) + \frac{\log|\Pi|}{\eta}.$$

1035 Let $\mathcal{F}_t := \sigma(\pi^{(1)}, z^{(1)}, \dots, \pi^{(t)}, z^{(t)})$ be a filtration, and let $\mathbb{E}_t[\cdot] := \mathbb{E}[\cdot | \mathcal{F}_t]$. For each $\pi \in \Pi$,
 1036 define a sequence of random variables $\{X_t(\pi)\}_{t=1}^T$ via

$$X_t(\pi) = \frac{1}{\eta} \log \left(\sum_{\pi'} q^{(t)}(\pi') \exp(\eta \cdot \hat{f}^{(t)}(\pi')) \right) - \langle e_{\pi}, \hat{f}^{(t)} \rangle.$$

1037 Using Lemma C.1 and a union bound, we have that for any $\eta > 0$, with probability at least $1 - \delta$, for
 1038 all $\pi \in \Pi$

$$\sum_{t=1}^T X_t(\pi) \leq \frac{1}{\eta} \sum_{t=1}^T \log(\mathbb{E}_{t-1}[\exp(\eta X_t(\pi))]) + \frac{\log(|\Pi|/\delta)}{\eta}.$$

1039 Since this bounded holds uniformly for all π , we have that with probability at least $1 - \delta$, for all
 1040 $\pi^* \in \Pi$,

$$\mathbf{Reg}_{\text{DM}}(\pi^*) \leq \sum_{t=1}^T \langle e_{\pi^*} - p^{(t)}, f^{M^{(t)}} \rangle + \frac{1}{\eta} \sum_{t=1}^T \log(\mathbb{E}_{t-1}[\exp(\eta X_t(\pi^*))]) + 2 \frac{\log(|\Pi|/\delta)}{\eta}.$$

1041 We compute that for any $\pi^* \in \Pi$,

$$\begin{aligned} &\log(\mathbb{E}_{t-1}[\exp(\eta X_t(\pi^*))]) \\ &= \log \left(\mathbb{E}_{\pi \sim p^{(t)}} \mathbb{E}_{z \sim M^{(t)}(\pi)} \mathbb{E}_{\pi' \sim q^{(t)}} \left[\exp \left(\frac{\eta}{p^{(t)}(\pi)} \cdot (g^{(t)}(\pi'; \pi, z) - g^{(t)}(\pi^*; \pi, z)) \right) \right] \right) \\ &\leq \mathbb{E}_{\pi \sim p^{(t)}} \mathbb{E}_{z \sim M^{(t)}(\pi)} \mathbb{E}_{\pi' \sim q^{(t)}} \left[\exp \left(\frac{\eta}{p^{(t)}(\pi)} \cdot (g^{(t)}(\pi'; \pi, z) - g^{(t)}(\pi^*; \pi, z)) \right) \right] - 1, \end{aligned}$$

1042 where we have used that $\log(x) \leq x - 1$ for $x > 0$. Hence, with probability at least $1 - \delta$, for all
 1043 $\pi^* \in \Pi$,

$$\begin{aligned} \mathbf{Reg}_{\text{DM}}(\pi^*) &\leq \sum_{t=1}^T \langle e_{\pi^*} - p^{(t)}, f^{M^{(t)}} \rangle + 2 \frac{\log(|\Pi|/\delta)}{\eta} \\ &\quad + \frac{1}{\eta} \left(\mathbb{E}_{\pi \sim p^{(t)}} \mathbb{E}_{z \sim M^{(t)}(\pi)} \mathbb{E}_{\pi' \sim q^{(t)}} \left[\exp \left(\frac{\eta}{p^{(t)}(\pi)} \cdot (g^{(t)}(\pi'; \pi, z) - g^{(t)}(\pi^*; \pi, z)) \right) \right] - 1 \right) \\ &= \sum_{t=1}^T \Gamma_{q^{(t)}, \eta}(p^{(t)}, g^{(t)}; \pi^*, M^{(t)}) + 2 \frac{\log(|\Pi|/\delta)}{\eta} \\ &\leq \text{exo}_\eta(\mathcal{M}) \cdot T + 2 \frac{\log(|\Pi|/\delta)}{\eta}, \end{aligned}$$

1044 where the last line uses that $(p^{(t)}, g^{(t)})$ are chosen to minimize the Exploration-By-Optimization
 1045 objective. Finally, using [Corollary 3.1](#), we have that $\text{exo}_\eta(\mathcal{M}) \leq \text{dec}_{(8\eta)^{-1}}(\text{co}(\mathcal{M}))$.

1046 □

1047 F.2 Proof of Theorem 2.2

1048 In this section we prove [Theorem 2.2](#). Most of the work consists of proving an improved lower
 1049 bound for the *stochastic* setting in which $M^{(t)} = M^*$ is fixed across t ([Theorem F.1](#)). We then
 1050 appeal to this stochastic lower bound with the class $\text{co}(\mathcal{M})$. Since $\text{co}(\mathcal{M})$ is equivalent to the set of
 1051 mixtures of models in \mathcal{M} , this establishes existence of distribution $\mu \in \Delta(\mathcal{M})$ and mixture model
 1052 $M_\mu = \mathbb{E}_{M \sim \mu}[M]$ for which regret in the stochastic setting must scale with $\text{dec}_{\gamma, \varepsilon_\gamma}(\text{co}(\mathcal{M}))$. The
 1053 proof concludes by arguing that this yields a lower bound for the adversarial setting when we sample
 1054 $M^{(t)} \sim \mu$.

1055 Throughout this section, we define the *one-sided variance* for a random variable Z as

$$\mathbb{V}_+[Z] := \mathbb{E}[(Z - \mathbb{E}[Z])_+^2].$$

1056 **Theorem 2.2** (Main lower bound). *Let $C(T) := c \cdot \log(T \wedge V(\mathcal{M}))$ for a sufficiently large numerical
 1057 constant $c > 0$. Set $\varepsilon_\gamma := \frac{\gamma}{4C(T)T}$. For any algorithm, there exists an oblivious adversary for which*

$$\mathbb{E}[\mathbf{Reg}_{\text{DM}}] + \sqrt{\mathbb{E}[\mathbf{Reg}_{\text{DM}}]_+^2} \geq \Omega(1) \cdot \sup_{\gamma > \sqrt{2C(T)T}} \text{dec}_{\gamma, \varepsilon_\gamma}(\text{co}(\mathcal{M})) \cdot T - O(T^{1/2}). \quad (13)$$

1058 We also have the following slight variant of [Theorem 2.2](#).

1059 **Theorem 2.2a.** *Let $C(T) := c \cdot \log(T \wedge V(\mathcal{M}))$ for a sufficiently large numerical constant $c > 0$.
 1060 Set $\varepsilon_\gamma := \frac{\gamma}{4C(T)T}$. For any algorithm, there exists an oblivious adversary for which $\mathbb{E}[\mathbf{Reg}_{\text{DM}}] \geq 0$
 1061 and*

$$\mathbb{E}[\mathbf{Reg}_{\text{DM}}] + \sqrt{\mathbb{E}[\mathbf{Reg}_{\text{DM}}] \cdot T} \geq \Omega(1) \cdot \sup_{\gamma > \sqrt{2C(T)T}} \text{dec}_{\gamma, \varepsilon_\gamma}(\text{co}(\mathcal{M})) \cdot T, \quad (42)$$

1062 **Proof of Theorem 2.2.** We invoke [Theorem F.1](#) with the model class $\text{co}(\mathcal{M})$, which implies that
 1063 there exists a distribution $\mu \in \Delta(\mathcal{M})$ for which

$$\mathbb{E}[\widetilde{\mathbf{Reg}}_{\text{DM}}] + \sqrt{\mathbb{V}_+[\widetilde{\mathbf{Reg}}_{\text{DM}}]} \geq L := 8^{-1} \cdot \sup_{\gamma > \sqrt{2C(T)T}} \text{dec}_{\gamma, \varepsilon_\gamma}(\text{co}(\mathcal{M})) \cdot T,$$

1064 where

$$\widetilde{\mathbf{Reg}}_{\text{DM}} := \sum_{t=1}^T \mathbb{E}_{\pi^{(t)} \sim p^{(t)}} \mathbb{E}_{M \sim \mu}[f^M(\pi_\mu) - f^M(\pi^{(t)})],$$

1065 and $\pi_\mu := \arg \max_{\pi \in \Pi} \mathbb{E}_{M \sim \mu}[f^M(\pi)]$, with the data generating process is (for each $t = 1, \dots, T$):

- 1066 • The learner samples $\pi^{(t)} \sim p^{(t)}$.

1067 • Nature samples $z^{(t)} \sim \mathbb{E}_{M \sim \mu}[M(\pi^{(t)})]$.

1068 Observe that this is equivalent in law to the following data-generating process, which constitutes an
1069 admissible adversary (with $M^{(t)} \in \mathcal{M}$):

- 1070 • The learner samples $\pi^{(t)} \sim p^{(t)}$.
- 1071 • Nature samples $M^{(t)} \sim \mu$ and $z^{(t)} \sim M^{(t)}(\pi^{(t)})$.

1072 Likewise, we can equivalently write

$$\widetilde{\mathbf{Reg}}_{\text{DM}} = \sum_{t=1}^T \mathbb{E}_{M^{(t)} \sim \mu} \mathbb{E}_{\pi^{(t)} \sim p^{(t)}} \left[f^{M^{(t)}}(\pi_{\mu}) - f^{M^{(t)}}(\pi^{(t)}) \right].$$

1073 Hence, all that remains is to relate the quantity $\widetilde{\mathbf{Reg}}_{\text{DM}}$ to the realized regret \mathbf{Reg}_{DM} for the sequence
1074 $M^{(1)}, \dots, M^{(T)}$, which entails removing the conditional expectation over $M^{(t)} \sim \mu$. To this end, we
1075 first observe that

$$\begin{aligned} \mathbb{E}[\widetilde{\mathbf{Reg}}_{\text{DM}}] &= \mathbb{E} \left[\sum_{t=1}^T \mathbb{E}_{\pi^{(t)} \sim p^{(t)}} \left[f^{M^{(t)}}(\pi_{\mu}) - f^{M^{(t)}}(\pi^{(t)}) \right] \right] \\ &\leq \mathbb{E} \left[\max_{\pi^* \in \Pi} \sum_{t=1}^T \mathbb{E}_{\pi^{(t)} \sim p^{(t)}} \left[f^{M^{(t)}}(\pi^*) - f^{M^{(t)}}(\pi^{(t)}) \right] \right] = \mathbb{E}[\mathbf{Reg}_{\text{DM}}]. \end{aligned}$$

1076 Next, note that since $\widetilde{\mathbf{Reg}}_{\text{DM}}$ is non-negative, $\mathbb{V}_+[\widetilde{\mathbf{Reg}}_{\text{DM}}] \leq \mathbb{E}[(\widetilde{\mathbf{Reg}}_{\text{DM}})_+^2]$. Define

$$\widehat{\mathbf{Reg}}_{\text{DM}} := \sum_{t=1}^T \mathbb{E}_{\pi^{(t)} \sim p^{(t)}} \left[f^{M^{(t)}}(\pi_{\mu}) - f^{M^{(t)}}(\pi^{(t)}) \right].$$

1077 Then we have

$$\begin{aligned} \mathbb{E}[(\widetilde{\mathbf{Reg}}_{\text{DM}})_+^2] &\leq 2 \mathbb{E}[(\widehat{\mathbf{Reg}}_{\text{DM}})_+^2] + 2 \mathbb{E}[(\widetilde{\mathbf{Reg}}_{\text{DM}} - \widehat{\mathbf{Reg}}_{\text{DM}})^2] \\ &\leq 2 \mathbb{E}[(\mathbf{Reg}_{\text{DM}})_+^2] + 2 \mathbb{E}[(\widetilde{\mathbf{Reg}}_{\text{DM}} - \mathbf{Reg}_{\text{DM}})^2] \\ &\leq 2 \mathbb{E}[(\mathbf{Reg}_{\text{DM}})_+^2] + 2T, \end{aligned}$$

1078 where the first inequality uses that $\widehat{\mathbf{Reg}}_{\text{DM}} \leq \mathbf{Reg}_{\text{DM}}$ almost surely, and the second inequality uses (i)
1079 $f^M \in [0, 1]$, and (ii) for any sequence of random variables $(Z_t)_{t=1}^T$ with $\mathbb{E}[Z_t | Z_1, \dots, Z_{t-1}] = 0$,
1080 $\mathbb{E}[(\sum_{t=1}^T Z_t)^2] = \sum_{t=1}^T \mathbb{E}[Z_t^2]$. Putting everything together, we conclude that

$$\mathbb{E}[\mathbf{Reg}_{\text{DM}}] + \sqrt{2 \mathbb{E}[(\mathbf{Reg}_{\text{DM}})_+^2]} \geq L - \sqrt{2T}.$$

1081 This proves [Theorem 2.2](#). To prove [Theorem 2.2a](#), we use that since $\widetilde{\mathbf{Reg}}_{\text{DM}} \in [0, T]$,

$$\mathbb{V}_+[\widetilde{\mathbf{Reg}}_{\text{DM}}] \leq T \cdot \mathbb{E}[\widetilde{\mathbf{Reg}}_{\text{DM}}] \leq T \cdot \mathbb{E}[\mathbf{Reg}_{\text{DM}}].$$

1082 □

1083 The following result concerns the *stochastic setting* in Foster et al. [18]. Here, there is a (unknown)
1084 underlying model $M^* \in \mathcal{M}$. For $t = 1, \dots, T$, data is generated through the process:

- 1085 • Learner samples $\pi^{(t)} \sim p^{(t)}$.
- 1086 • Nature samples $z^{(t)} \sim M^*(\pi^{(t)})$.

1087 In addition, regret simplifies to

$$\mathbf{Reg}_{\text{DM}} = \sum_{t=1}^T \mathbb{E}_{\pi^{(t)} \sim p^{(t)}} \left[f^{M^*}(\pi_{M^*}) - f^{M^*}(\pi^{(t)}) \right] \quad (43)$$

1088 For a fixed algorithm, let \mathbb{P}^M denote the law of $\mathcal{H}^{(T)}$ when $M^* = M$, and let $\mathbb{E}^M[\cdot]$ and $\mathbb{V}_+ \sup M[\cdot]$
1089 denote the corresponding expectation non-negative variance. Our main lower bound for the stochastic
1090 setting is as follows.

1091 **Theorem F.1.** Let $C(T) := 2^9 \log(T \wedge V(\mathcal{M}))$, and set $\varepsilon_\gamma = \frac{\gamma}{4C(T)T}$. For any algorithm, there
 1092 exists a model in \mathcal{M} for which

$$\mathbb{E}^M[\mathbf{Reg}_{\text{DM}}] + \sqrt{\mathbb{V}_+^M[\mathbf{Reg}_{\text{DM}}]} \geq 8^{-1} \cdot \sup_{\gamma \geq 4\sqrt{C(T)T}} \sup_{\bar{M} \in \mathcal{M}} \text{dec}_\gamma(\mathcal{M}_{\varepsilon_\gamma}(\bar{M}), \bar{M}) \cdot T.$$

1093 The general structure of the lower bound follows that of Theorem 3.1 in Foster et al. [18], with
 1094 the main difference being that we use a more refined change-of-measure argument to move from a
 1095 “reference” model $\bar{M} \in \mathcal{M}$ to a worst-case alternative. Specifically, we replace Lemma A.11 in Foster
 1096 et al. [18], which requires an almost sure bound on the random variables under consideration (in our
 1097 case, regret), with Lemma C.4, which requires only boundedness of the second moment. Combining
 1098 this with a self-bounding argument that takes advantage of the localized model class yields the result.

1099 **Proof of Theorem F.1.** Throughout this proof we will use that \mathbf{Reg}_{DM} is non-negative in the
 1100 stochastic setting, which can be seen by inspecting (43) (in the general adversarial setting, it is
 1101 possible for \mathbf{Reg}_{DM} to be negative).

1102 Let us introduce some additional notation. For $M \in \mathcal{M}$, define $g^M(\pi) = f^M(\pi_M) - f^M(\pi)$, and for
 1103 $p \in \Delta(\Pi)$, let $g^M(p) = \mathbb{E}_{\pi \sim p}[g^M(\pi)]$. Let $\hat{p} := \frac{1}{T} \sum_{t=1}^T p^{(t)}$, and $p_M := \mathbb{E}^M\left[\frac{1}{T} \sum_{t=1}^T p^{(t)}\right]$.

1104 To begin, fix $\bar{M} \in \mathcal{M}$, $\gamma > 0$, and $\varepsilon > 0$, and set

$$M = \arg \max_{M \in \mathcal{M}_\varepsilon(\bar{M})} \mathbb{E}_{\pi \sim p_{\bar{M}}}[f^M(\pi_M) - f^M(\pi) - \gamma \cdot D_{\text{H}}^2(M(\pi), \bar{M}(\pi))].$$

1105 Abbreviate $\text{dec}_\gamma \equiv \text{dec}_\gamma(\mathcal{M}_\varepsilon(\bar{M}), \bar{M})$. The definition of the DEC implies that

$$\text{dec}_\gamma \leq \mathbb{E}_{p_{\bar{M}}}[g^M(\pi)] - \gamma \cdot \mathbb{E}_{p_{\bar{M}}}[D_{\text{H}}^2(M(\pi), \bar{M}(\pi))] = \mathbb{E}^{\bar{M}}[g^M(\hat{p})] - \gamma \cdot \mathbb{E}_{p_{\bar{M}}}[D_{\text{H}}^2(M(\pi), \bar{M}(\pi))]. \quad (44)$$

1106

1107 **Change of measure.** To proceed, we write

$$\begin{aligned} \mathbb{E}^{\bar{M}}[g^M(\hat{p})] &= \mathbb{E}^{\bar{M}}[g^M(\hat{p}) - g^{\bar{M}}(\hat{p}) - \mathbb{E}^M[g^M(\hat{p})]] + \mathbb{E}^{\bar{M}}[g^{\bar{M}}(\hat{p})] + \mathbb{E}^M[g^M(\hat{p})] \\ &\leq \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}) - \mathbb{E}^M[g^M(\hat{p})])_+] + \mathbb{E}^{\bar{M}}[g^{\bar{M}}(\hat{p})] + \mathbb{E}^M[g^M(\hat{p})]. \end{aligned} \quad (45)$$

1108 We recall the following technical lemma.

1109 **Lemma C.4.** Let \mathbb{P} and \mathbb{Q} be probability distributions over a measurable space $(\mathcal{X}, \mathcal{F})$. Then for
 1110 all functions $h : \mathcal{X} \rightarrow \mathbb{R}$,

$$|\mathbb{E}_{\mathbb{P}}[h(X)] - \mathbb{E}_{\mathbb{Q}}[h(X)]| \leq \sqrt{2^{-1}(\mathbb{E}_{\mathbb{P}}[h^2(X)] + \mathbb{E}_{\mathbb{Q}}[h^2(X)]) \cdot D_{\text{H}}^2(\mathbb{P}, \mathbb{Q})}. \quad (25)$$

1111 Defining $h(\hat{p}) = (g^M(\hat{p}) - g^{\bar{M}}(\hat{p}) - \mathbb{E}^M[g^M(\hat{p})])_+$, Lemma C.4 implies that

$$\begin{aligned} &\mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}) - \mathbb{E}^M[g^M(\hat{p})])_+] \\ &\leq \mathbb{E}^M[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}) - \mathbb{E}^M[g^M(\hat{p})])_+] + \sqrt{(\mathbb{E}^M[h(\hat{p})^2] + \mathbb{E}^{\bar{M}}[h(\hat{p})^2]) \cdot D_{\text{H}}^2(\mathbb{P}^M, \mathbb{P}^{\bar{M}})} \\ &\leq \mathbb{E}^M[g^M(\hat{p})] + \sqrt{(\mathbb{E}^M[h(\hat{p})^2] + \mathbb{E}^{\bar{M}}[h(\hat{p})^2]) \cdot D_{\text{H}}^2(\mathbb{P}^M, \mathbb{P}^{\bar{M}})}, \end{aligned} \quad (46)$$

1112 where we have used that $g^M, g^{\bar{M}} \geq 0$. We proceed to bound the second moment terms. First, we have

$$\begin{aligned} \mathbb{E}^M[h(\hat{p})^2] &= \mathbb{E}^M[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}) - \mathbb{E}^M[g^M(\hat{p})])_+^2] \\ &\leq \mathbb{E}^M[(g^M(\hat{p}) - \mathbb{E}^M[g^M(\hat{p})])_+^2] \\ &= \mathbb{V}_+^M[g^M(\hat{p})]. \end{aligned} \quad (47)$$

1113 where the first inequality uses that $g^{\bar{M}} \geq 0$. For the second variance term, we have

$$\mathbb{E}^{\bar{M}}[h(\hat{p})^2] = \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}) - \mathbb{E}^M[g^M(\hat{p})])_+^2] \leq \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}))_+^2].$$

1114 We have

$$\begin{aligned} & \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}))_+^2] \\ &= \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}))_+(f^M(\pi_M) - f^{\bar{M}}(\pi_{\bar{M}}) + f^{\bar{M}}(\hat{p}) - f^M(\hat{p}))_+] \\ &\leq \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}))_+(f^M(\pi_M) - f^{\bar{M}}(\pi_{\bar{M}}))_+] + \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}))_+ f^{\bar{M}}(\hat{p}) - f^M(\hat{p})_+]. \end{aligned}$$

1115 For the first term above, we have

$$\mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}))_+(f^M(\pi_M) - f^{\bar{M}}(\pi_{\bar{M}}))_+] \leq \varepsilon \cdot \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}))_+] \leq \varepsilon \cdot \mathbb{E}^{\bar{M}}[g^M(\hat{p})],$$

1116 where we have used the localization property and the fact that $g^M, g^{\bar{M}} \geq 0$. For the second term,
1117 using the AM-GM inequality, we have

$$\begin{aligned} & \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}))_+ f^{\bar{M}}(\hat{p}) - f^M(\hat{p})_+] \\ &\leq \frac{1}{2} \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}))_+^2] + \frac{1}{2} \mathbb{E}^{\bar{M}}[(f^{\bar{M}}(\hat{p}) - f^M(\hat{p}))^2] \\ &\leq \frac{1}{2} \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}))_+^2] + \frac{1}{2} \mathbb{E}_{\pi \sim p_{\bar{M}}}[(f^M(\pi) - f^{\bar{M}}(\pi))^2] \\ &\leq \frac{1}{2} \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}))_+^2] + \frac{1}{2} \mathbb{E}_{\pi \sim p_{\bar{M}}}[D_{\text{H}}^2(M(\pi), \bar{M}(\pi))], \end{aligned}$$

1118 where the last line uses that rewards are observed and bounded in $[0, 1]$. After combining these results
1119 and rearranging, we have

$$\mathbb{E}^{\bar{M}}[h(\hat{p})^2] \leq \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}))_+^2] \leq 2\varepsilon \cdot \mathbb{E}^{\bar{M}}[g^M(\hat{p})] + \mathbb{E}_{\pi \sim p_{\bar{M}}}[D_{\text{H}}^2(M(\pi), \bar{M}(\pi))]. \quad (48)$$

1120 From Lemma A.13 of Foster et al. [18], we have

$$D_{\text{H}}^2(\mathbb{P}^M, \mathbb{P}^{\bar{M}}) \leq C(T) \cdot T \cdot \mathbb{E}_{\pi \sim p_{\bar{M}}}[D_{\text{H}}^2(M(\pi), \bar{M}(\pi))], \quad (49)$$

1121 where $C(T) \leq 2^8 \cdot \log(T \wedge V(\mathcal{M}))$.

1122 Combining the variance bounds with (46), we have

$$\begin{aligned} & \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}) - \mathbb{E}^M[g^M(\hat{p})])_+] \\ &\leq \mathbb{E}^M[g^M(\hat{p})] + \sqrt{(\mathbb{V}_+^M[g^M(\hat{p})] + 2\varepsilon \cdot \mathbb{E}^{\bar{M}}[g^M(\hat{p})] + \mathbb{E}_{\pi \sim p_{\bar{M}}}[D_{\text{H}}^2(M(\pi), \bar{M}(\pi))]) \cdot D_{\text{H}}^2(\mathbb{P}^M, \mathbb{P}^{\bar{M}})} \\ &\leq \mathbb{E}^M[g^M(\hat{p})] + \sqrt{2\mathbb{V}_+^M[g^M(\hat{p})]} + \sqrt{(2\varepsilon \cdot \mathbb{E}^{\bar{M}}[g^M(\hat{p})] + \mathbb{E}_{\pi \sim p_{\bar{M}}}[D_{\text{H}}^2(M(\pi), \bar{M}(\pi))]) \cdot D_{\text{H}}^2(\mathbb{P}^M, \mathbb{P}^{\bar{M}})} \\ &\leq \mathbb{E}^M[g^M(\hat{p})] + \sqrt{2\mathbb{V}_+^M[g^M(\hat{p})]} + \sqrt{C(T)T} \cdot \mathbb{E}_{\pi \sim p_{\bar{M}}}[D_{\text{H}}^2(M(\pi), \bar{M}(\pi))] \\ &\quad + \sqrt{2\varepsilon \mathbb{E}^{\bar{M}}[g^M(\hat{p})] \cdot C(T)T \mathbb{E}_{\pi \sim p_{\bar{M}}}[D_{\text{H}}^2(M(\pi), \bar{M}(\pi))]}, \end{aligned}$$

1123 where the second inequality uses that $D_{\text{H}}^2(\cdot, \cdot) \leq 2$ and the last inequality uses (49). where the second
1124 inequality uses that $D_{\text{H}}^2(\mathbb{P}^M, \mathbb{P}^{\bar{M}}) \leq 2$.

1125 Now, suppose we restrict to $\varepsilon \leq \frac{\gamma}{4TC(T)}$. Then we have

$$\begin{aligned} \sqrt{2\varepsilon \cdot \mathbb{E}^{\bar{M}}[g^M(\hat{p})] \cdot C(T)T \mathbb{E}_{\pi \sim p_{\bar{M}}}[D_{\text{H}}^2(M(\pi), \bar{M}(\pi))]} &\leq \sqrt{\mathbb{E}^{\bar{M}}[g^M(\hat{p})] \cdot \frac{\gamma}{2} \cdot \mathbb{E}_{\pi \sim p_{\bar{M}}}[D_{\text{H}}^2(M(\pi), \bar{M}(\pi))]} \\ &\leq \frac{1}{2} \mathbb{E}^{\bar{M}}[g^M(\hat{p})] + \frac{\gamma}{4} \cdot \mathbb{E}_{\pi \sim p_{\bar{M}}}[D_{\text{H}}^2(M(\pi), \bar{M}(\pi))]. \end{aligned}$$

1126 Altogether, we have

$$\begin{aligned} & \mathbb{E}^{\bar{M}}[(g^M(\hat{p}) - g^{\bar{M}}(\hat{p}) - \mathbb{E}^M[g^M(\hat{p})])_+] \\ &\leq \mathbb{E}^M[g^M(\hat{p})] + \sqrt{2\mathbb{V}_+^M[g^M(\hat{p})]} + (\sqrt{C(T)T} + \gamma/4) \cdot \mathbb{E}_{\pi \sim p_{\bar{M}}}[D_{\text{H}}^2(M(\pi), \bar{M}(\pi))] + \frac{1}{2} \mathbb{E}^{\bar{M}}[g^M(\hat{p})] \end{aligned}$$

1127 and, using (45),

$$\begin{aligned} \mathbb{E}^{\bar{M}}[g^M(\hat{p})] &\leq 2\mathbb{E}^M[g^M(\hat{p})] + \mathbb{E}^{\bar{M}}[g^{\bar{M}}(\hat{p})] + \sqrt{2\mathbb{V}_+^M[g^M(\hat{p})]} \\ &\quad + (\sqrt{C(T)T} + \gamma/4) \cdot \mathbb{E}_{\pi \sim p_{\bar{M}}}[D_{\text{H}}^2(M(\pi), \bar{M}(\pi))] + \frac{1}{2} \mathbb{E}^{\bar{M}}[g^M(\hat{p})]. \end{aligned}$$

1128 After rearranging, this implies that

$$\mathbb{E}^{\overline{M}}[g^M(\widehat{p})] \leq 4\mathbb{E}^M[g^M(\widehat{p})] + 2\mathbb{E}^{\overline{M}}[g^{\overline{M}}(\widehat{p})] + \sqrt{8\mathbb{V}_+^M[g^M(\widehat{p})]} + 2(\sqrt{C(T)T} + \gamma/4) \cdot \mathbb{E}_{\pi \sim p_{\overline{M}}} [D_{\text{H}}^2(M(\pi), \overline{M}(\pi))]. \quad (50)$$

1129 **Completing the proof.** Combining (50) with (44), we have

$$\text{dec}_\gamma \leq 4\mathbb{E}^M[g^M(\widehat{p})] + 2\mathbb{E}^{\overline{M}}[g^{\overline{M}}(\widehat{p})] + \sqrt{8\mathbb{V}_+^M[g^M(\widehat{p})]} + \left(2(\sqrt{C(T)T} + \gamma/4) - \gamma\right) \cdot \mathbb{E}_{\pi \sim p_{\overline{M}}} [D_{\text{H}}^2(M(\pi), \overline{M}(\pi))].$$

1130 In particular, whenever $\gamma \geq 4\sqrt{C(T)T}$, this implies that there exists an instance $M' \in \{M, \overline{M}\}$ for
1131 which

$$\mathbb{E}^{M'}[g^{M'}(\widehat{p})] + \sqrt{\mathbb{V}_+^{M'}[g^{M'}(\widehat{p})]} \geq 8^{-1} \cdot \text{dec}_\gamma.$$

1132 Finally, we observe that $g^{M'}(\widehat{p})$ is identical in law to \mathbf{Reg}_{DM} under $\mathbb{P}^{M'}$.

1133

□

1134 F.3 Proof of Theorem 2.3

1135 **Theorem 2.3.** Suppose there exists $M_0 \in \mathcal{M}$ such that f^{M_0} is a constant function, and that $|\Pi| < \infty$.

1136 1. If there exists $\rho > 0$ s.t. $\lim_{\gamma \rightarrow \infty} \text{dec}_\gamma(\text{co}(\mathcal{M})) \cdot \gamma^\rho = 0$, then $\lim_{T \rightarrow \infty} \frac{\mathfrak{M}(\mathcal{M}, T)}{T^p} = 0$ for $p < 1$.

1137 2. If $\lim_{\gamma \rightarrow \infty} \text{dec}_\gamma(\text{co}(\mathcal{M})) \cdot \gamma^\rho > 0$ for all $\rho > 0$, then $\lim_{T \rightarrow \infty} \frac{\mathfrak{M}(\mathcal{M}, T)}{T^p} = \infty$ for all $p < 1$.

1138 The same conclusion holds when $\Pi = \Pi_T$ grows with T , but has $\log|\Pi_T| = O(T^q)$ for any $q < 1$.

1139 **Proof of Theorem 2.3.** This proof closely follows that of Theorem 3.5 in Foster et al. [18].

1140 **Upper bound.** Assume that $\lim_{\gamma \rightarrow \infty} \text{dec}_\gamma(\text{co}(\mathcal{M})) \cdot \gamma^\rho = 0$ for some $\rho > 0$, and that $\log|\Pi_T| =$
1141 $\widetilde{O}(T^q)$ for some $q < 1$. Using Theorem 2.1 with $\delta = 1/T$, we have that for each T , for all
1142 adversaries,

$$\mathbb{E}n[\mathbf{Reg}_{\text{DM}}(T)] \leq \widetilde{O}(\text{dec}_\gamma(\text{co}(\mathcal{M})) \cdot T + \gamma \cdot \log|\Pi_T|) \leq \widetilde{O}(\text{dec}_\gamma(\text{co}(\mathcal{M})) \cdot T + \gamma \cdot T^q),$$

1143 with $\widetilde{O}(\cdot)$ hiding factors logarithmic in T . For each T , we set $\gamma = \gamma_T := T^{\frac{1-q}{1+\rho}}$; recall that $1 - q > 0$.
1144 The assumption that $\lim_{\gamma \rightarrow \infty} \text{dec}_\gamma(\text{co}(\mathcal{M})) \cdot \gamma^\rho = 0$, implies that for all $\varepsilon > 0$, there exists $\gamma' > 0$
1145 such that $\text{dec}_\gamma(\text{co}(\mathcal{M})) \leq \varepsilon/\gamma^\rho$ for all $\gamma \geq \gamma'$. For T sufficiently large, this implies that for all
1146 adversaries

$$\mathbb{E}[\mathbf{Reg}_{\text{DM}}] \leq \widetilde{O}\left(\frac{T}{\gamma_T^\rho} + \gamma_T \cdot T^q\right) = \widetilde{O}(T^{\frac{1+\rho q}{1+\rho}}).$$

1147 Defining $p' := \frac{1}{2}(p + 1) < 1$, this establishes that

$$\lim_{T \rightarrow \infty} \frac{\mathfrak{M}(\mathcal{M}, T)}{T^{p'}} = 0.$$

1148 **Lower bound.** Assume that $\lim_{\gamma \rightarrow \infty} \text{dec}_\gamma(\text{co}(\mathcal{M})) \cdot \gamma^\rho = \infty$ for all $\rho > 0$ (this is equivalent
1149 to assuming that $\lim_{\gamma \rightarrow \infty} \text{dec}_\gamma(\text{co}(\mathcal{M})) \cdot \gamma^\rho > 0$ for all $\rho > 0$, as in the theorem statement). Let
1150 $\rho \in (0, 1/2)$ be fixed. Using Theorem 2.2a, we are guaranteed that for any algorithm, there exists an
1151 adversary for which $\mathbb{E}[\mathbf{Reg}_{\text{DM}}] \geq 0$ and

$$\mathbb{E}[\mathbf{Reg}_{\text{DM}}] + \sqrt{\mathbb{E}[\mathbf{Reg}_{\text{DM}}] \cdot T} = \widetilde{\Omega}(\text{dec}_{\gamma, \varepsilon(\gamma, T)}(\text{co}(\mathcal{M})) \cdot T),$$

1152 for all $\gamma = \omega(\sqrt{T \log(T)})$, where $\varepsilon(\gamma, T) := c \cdot \frac{\gamma}{T \log(T)}$ for a sufficiently small numerical constant
1153 $c \leq 1$. Since there exists $M_0 \in \mathcal{M}$ such that the function f^{M_0} is constant, Lemma B.1 of Foster et al.
1154 [18] further implies that

$$\mathbb{E}[\mathbf{Reg}_{\text{DM}}] + \sqrt{\mathbb{E}[\mathbf{Reg}_{\text{DM}}] \cdot T} = \widetilde{\Omega}(\varepsilon(\gamma, T) \cdot \text{dec}_\gamma(\text{co}(\mathcal{M})) \cdot T).$$

1155 For each T , set $\gamma = \gamma_T := T$. By the assumption that $\lim_{\gamma \rightarrow \infty} \text{dec}_\gamma(\text{co}(\mathcal{M})) \cdot \gamma^\rho = \infty$, we have
 1156 that for T sufficiently large, $\text{dec}_{\gamma_T}(\text{co}(\mathcal{M})) \geq \gamma_T^{-\rho}$, which implies that and

$$\mathbb{E}[\mathbf{Reg}_{\text{DM}}] + \sqrt{\mathbb{E}[\mathbf{Reg}_{\text{DM}}] \cdot T} = \tilde{\Omega}\left(\frac{T}{\gamma_T^\rho}\right),$$

1157 where we have used that $\varepsilon(\gamma_T, T) \propto \frac{1}{\log(T)}$. Rearranging, this implies that

$$\mathbb{E}[\mathbf{Reg}_{\text{DM}}] = \tilde{\Omega}(T^{1-2\rho}).$$

1158 Hence, for any $p \in (0, 1)$, by setting $\rho = \frac{1-p}{2} \in (0, 1/2)$, we have

$$\mathbb{E}[\mathbf{Reg}_{\text{DM}}] = \tilde{\Omega}(T^p).$$

1159 Applying this argument with $p' = \frac{1}{2}(p+1) \in (1/2, 1)$ yields

$$\lim_{T \rightarrow \infty} \frac{\mathfrak{M}(\mathcal{M}, T)}{T^p} = \infty.$$

1160 □

1161 F.4 Sub-Chebychev Algorithms

1162 **Proposition F.1.** Any random variable with $\mathbb{E}[X_+^2] \leq R$ has

$$\mathbb{P}(X_+ > t) \leq \frac{R^2}{t^2}, \quad \forall t > 0.$$

1163 Conversely, if $X \in (-\infty, B)$ and has $\mathbb{P}(X_+ > t) \leq \frac{R^2}{t^2} \forall t > 0$, then

$$\mathbb{E}[X_+^2] \leq R^2(\log(B/R) + 1).$$

1164 **Proof of Proposition F.1.** For the first direction, note that if $\mathbb{E}[X_+^2] \leq R$, Chebychev's inequality
 1165 implies that for all $t > 0$,

$$\mathbb{P}(X_+^2 > t) \leq \frac{R^2}{t^2}. \quad (51)$$

1166 For the other direction, since $X_+ \in [0, B]$ almost surely, we have

$$\mathbb{E}[X_+^2] = \int_0^B \mathbb{P}(X_+ > t) t dt \leq R^2 + \int_R^B \mathbb{P}(X_+ > t) t dt \leq R^2 + R^2 \int_R^B \frac{1}{t} dt \leq R^2 + R^2 \log(B/R).$$

1167 □

1168 **Proposition F.2.** Suppose that for any $\delta > 0$, an algorithm (with δ as a parameter) ensures that with
 1169 probability at least $1 - \delta$,

$$\mathbf{Reg}_{\text{DM}} \leq R \log^\rho(\delta^{-1})$$

1170 for some $R \geq 1$ and $\rho > 0$. Then the algorithm, when invoked with parameter $\delta = 1/T^2$, is
 1171 sub-Chebychev with parameter $5^{1/2} R \log^\rho(T)$.

1172 **Proof of Proposition F.2.** Set $\delta = 1/T^2$. Then, since $|\mathbf{Reg}_{\text{DM}}| \leq T$, the law of total expectation
 1173 implies that

$$\mathbb{E}[(\mathbf{Reg}_{\text{DM}})_+^2] \leq R^2 \log^{2\rho}(T^2) + T^2/T^2 \leq 5R^2 \log^2(T),$$

1174 where we have used that $R \geq 1$. Chebychev's inequality now implies that for all $t > 0$

$$\mathbb{P}((\mathbf{Reg}_{\text{DM}})_+ \geq t) \leq \frac{\mathbb{E}[(\mathbf{Reg}_{\text{DM}})_+^2]}{t^2} \leq \frac{5R^2 \log^{2\rho}(T)}{t}.$$

1175 □

1176 **Corollary 2.1.** Any regret minimization algorithm with sub-Chebychev parameter $R > 0$ must have

$$R \geq \tilde{\Omega}(1) \cdot \sup_{\gamma > \sqrt{2C(T)T}} \text{dec}_{\gamma, \varepsilon_\gamma}(\text{co}(\mathcal{M})) \cdot T - O(T^{1/2}). \quad (15)$$

1177 **Proof of Corollary 2.1.** This result immediately follows from [Proposition F.1](#), [Proposition F.2](#), and
 1178 [Theorem 2.2](#). □