

Instance-Optimality in Interactive Decision Making: Toward a Non-Asymptotic Theory

Andrew Wagenmaker

University of Washington, Seattle, WA

AJWAGEN@CS.WASHINGTON.EDU

Dylan J. Foster

Microsoft Research, New England

DYLANFOSTER@MICROSOFT.COM

Editors: Gergely Neu and Lorenzo Rosasco

Abstract

We consider the development of adaptive, instance-dependent algorithms for interactive decision making (bandits, reinforcement learning, and beyond) that, rather than only performing well in the worst case, adapt to favorable properties of real-world instances for improved performance. We aim for *instance-optimality*, a strong notion of adaptivity which asserts that, on any particular problem instance, the algorithm under consideration outperforms all consistent algorithms. Instance-optimality enjoys a rich asymptotic theory originating from the work of [Lai and Robbins \(1985\)](#) and [Graves and Lai \(1997\)](#), but *non-asymptotic* guarantees have remained elusive outside of certain special cases. Even for problems as simple as tabular reinforcement learning, existing algorithms do not attain instance-optimal performance until the number of rounds of interaction is *doubly exponential* in the number of states.

In this paper, we take the first step toward developing a non-asymptotic theory of instance-optimal decision making with general function approximation. We introduce a new complexity measure, the Allocation-Estimation Coefficient (AEC), and provide a new algorithm, AE^2 , which attains non-asymptotic instance-optimal performance at a rate controlled by the AEC. Our results recover the best known guarantees for well-studied problems such as finite-armed and linear bandits and, when specialized to tabular reinforcement learning, attain the first instance-optimal regret bounds with polynomial dependence on all problem parameters, improving over prior work exponentially. We complement these results with lower bounds that show that i) existing notions of statistical complexity are insufficient to derive non-asymptotic guarantees, and ii) under certain technical conditions, boundedness of the Allocation-Estimation Coefficient is *necessary* to learn an instance-optimal allocation of decisions in finite time.

1. Introduction

We consider the development of adaptive, sample-efficient algorithms for *interactive decision making*, encompassing bandit problems and reinforcement learning with general function approximation. For decision making in high-dimensional spaces with a long horizon, existing approaches ([Lillicrap et al., 2015](#); [Mnih et al., 2015](#); [Silver et al., 2016](#)) are sample-hungry, which presents an obstacle for real-world deployment in settings where data is scarce or high-quality simulators are not available. To overcome this challenge, algorithms should both i) flexibly incorporate users' domain knowledge, as expressed via modeling and function approximation, and ii) explore the environment in a deliberate, adaptive fashion, taking advantage of favorable structure whenever possible.

Toward achieving these goals, a major area of research aims to develop algorithms with optimal sample complexity and understand the fundamental limits for such algorithms ([Russo and Van Roy, 2013](#); [Jiang et al., 2017](#); [Sun et al., 2019](#); [Wang et al., 2020](#); [Du et al., 2021](#); [Jin et al., 2021](#); [Foster](#)

et al., 2021), and the foundations are beginning to fall into place. In particular, focusing on *minimax regret* (that is, the best regret that can be achieved for a worst-case problem instance in a given class of problems), Foster et al. (2021, 2022b, 2023) provide unified algorithm design principles and measures of statistical complexity that are both necessary and sufficient for low regret. However, minimax regret and other notions of worst-case performance are inherently pessimistic, and may not be sufficient to close the gap between theory and practice. For example, recent work has shown that algorithms that are optimal in the worst-case can be arbitrarily suboptimal on “easier” instances (Wagenmaker et al., 2022b). To overcome these challenges and develop algorithms that perform well on *every* instance, a promising approach is to develop algorithms that *adapt* to the difficulty of the problem instance under consideration.

The performance of such adaptive algorithms can be quantified through *instance-dependent* regret bounds, which become smaller (leading to low regret) when the underlying problem instance is favorable. Algorithms with such guarantees have been studied throughout the literature on bandits and reinforcement learning; basic examples include adapting to large gaps in value between alternative actions (Lai and Robbins, 1985) or low noise or variance in bandit problems (Allenberg et al., 2006; Hazan and Kale, 2011; Foster et al., 2016; Wei and Luo, 2018; Bubeck et al., 2018), and adapting to the difficulty of reaching certain states in Markov Decision Processes (Zanette and Brunskill, 2019; Simchowitz and Jamieson, 2019; Dann et al., 2021; Wagenmaker et al., 2022b).

While there are many notions of adaptivity and instance-dependence, they are generally incomparable. A stronger notion of adaptivity is *instance-optimality*, which asserts that the performance of the algorithm on a problem instance of interest exceeds that of *any consistent algorithm* (that is, any algorithm with sublinear regret for all problem instances). Instance-optimality enjoys a rich theory originating with the work of Lai and Robbins (1985) and Graves and Lai (1997), with a celebrated line of research developing sharp guarantees for the special case of finite-armed bandits (Burnetas and Katehakis, 1996; Garivier et al., 2016; Kaufmann et al., 2016; Lattimore, 2018; Garivier et al., 2019). Beyond the finite-armed bandit setting, however, development has been largely *asymptotic* in nature, and existing algorithms either:

1. achieve instance-optimality only as $T \rightarrow \infty$ (or, to the extent that they are non-asymptotic, require T to be exponentially large in problem-dependent parameters) (Graves and Lai, 1997; Komiyama et al., 2015; Combes et al., 2017; Degenne et al., 2020b; Dong and Ma, 2022), or
2. achieve non-asymptotic guarantees, but require restrictive modeling assumptions such as linear function approximation (Tirinzi et al., 2020; Kirschner et al., 2021).

Indeed, even for the simple problem of tabular (finite-state/action) reinforcement learning, existing algorithms do not attain instance-optimal performance until the number of rounds of interaction is *doubly exponential* in the number of states (Ok et al., 2018; Dong and Ma, 2022). In this paper, we address these challenges, providing algorithms that i) accommodate flexible, general-purpose function approximation, and ii) attain instance-optimality in finite time, in a sense which is itself optimal.

Contributions. We take the first steps toward building a non-asymptotic theory of instance-optimal decision making. We observe that asymptotic characterizations for instance-optimality:

1. reflect the regret incurred by an allocation of decisions designed to optimally distinguish the ground truth problem instance from a set of alternatives, but
2. do not capture the statistical complexity required to *learn* such an allocation.

To address this, we introduce a new complexity measure, the Allocation-Estimation Coefficient (AEC), which aims to capture the statistical complexity of learning an optimal Graves-Lai allocation. We provide a new algorithm, AE^2 , which attains non-asymptotic instance-optimal regret at a rate controlled by the AEC. We complement this result with lower bounds that show that under certain technical conditions, boundedness of the Allocation-Estimation Coefficient is not just sufficient, but *necessary* to learn an instance-optimal allocation in finite time.

Our algorithm is simple, and can be applied to any hypothesis class in a generic fashion. It recovers the best known guarantees for standard problems such as finite-armed and linear bandits and, when specialized to tabular reinforcement learning, achieves the first instance-optimal regret bounds with polynomial dependence on all problem parameters. We believe that our approach clarifies and elucidates many tradeoffs and statistical considerations left implicit in prior work, and hope that it will serve as a foundation for further development of instance-optimal algorithms.

1.1. Interactive Decision Making

We adopt the *Decision Making with Structured Observations* (DMSO) framework of Foster et al. (2021), which is a general setting for interactive decision making that encompasses bandit problems (structured, contextual, and so forth) and reinforcement learning with function approximation.

The DMSO framework is specified by a *decision space* Π , *reward space* $\mathcal{R} \subseteq \mathbb{R}$, and *observation space* \mathcal{O} . The learner is given access to a (known) *model class* $\mathcal{M} \subset (\Pi \rightarrow \Delta_{\mathcal{R} \times \mathcal{O}})$, and it is assumed there exists some true model $M^* \in \mathcal{M}$, unknown to the learner, which represents the underlying environment. Formally, we make the following assumption.

Assumption 1.1 (Realizability). *We have that $M^* \in \mathcal{M}$.*

The learning protocol consists of T rounds. For each round $t = 1, \dots, T$:

1. The learner selects a *decision* $\pi^t \in \Pi$.
2. The learner receives a reward $r^t \in \mathcal{R}$ and observation $o^t \in \mathcal{O}$ sampled $(r^t, o^t) \sim M^*(\pi^t)$, and observes (r^t, o^t) .

We can think of the model class \mathcal{M} as representing the learner’s prior knowledge about the decision making problem, and it allows one to appeal to estimation and function approximation. For structured bandit problems, for example, models correspond to reward distributions, and \mathcal{M} encodes structure in the reward landscape. For reinforcement learning problems, models correspond to Markov decision processes (MDPs), and \mathcal{M} typically encodes structure in value functions or transition probabilities. See Appendix A.6 and Appendix A.7 for concrete examples of how standard decision-making settings can be instantiated within the DMSO framework, and Foster et al. (2021) for further background. For a model $M \in \mathcal{M}$, $\mathbb{E}^{M, \pi}[\cdot]$ denotes the expectation under the process $(r, o) \sim M(\pi)$, $f^M(\pi) := \mathbb{E}^{M, \pi}[r]$ denotes the mean reward function, and $\pi_M := \arg \max_{\pi \in \Pi} f^M(\pi)$ denotes the optimal decision. When the algorithm is clear from context, $\mathbb{E}^M[\cdot]$ and $\mathbb{P}^M[\cdot]$ refer to the expectation and probability measure, respectively, induced over histories under M . When the context is clear, we overload notation somewhat and use $\mathbb{P}^{M, \pi}[\cdot]$ to refer to the conditional density over $\mathcal{R} \times \mathcal{O}$ induced by playing π on M . We make the following assumptions.

Assumption 1.2 (Bounded Reward Means). *For all $M \in \mathcal{M}$, $\pi \in \Pi$, we have $f^M(\pi) \in [0, 1]$.*

Assumption 1.3 (Unique Optimal Action). *For the ground truth model $M^* \in \mathcal{M}$, the optimal action π_{M^*} is unique.*

Note that the latter assumption is standard in the literature on instance-optimality. We measure performance in terms of regret, which is given by

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

where p^t is the learner’s randomization distribution for round t . In addition, we define $\Delta^M(\pi) = f^M(\pi_M) - f^M(\pi)$ as the *suboptimality gap* function for model M and decision π , and the *minimum suboptimality gap* as

$$\Delta_{\min}^M := \begin{cases} \inf_{\pi \in \Pi: \Delta^M(\pi) > 0} \Delta^M(\pi), & \pi_M \text{ is unique,} \\ 0, & \text{otherwise.} \end{cases} \quad (2)$$

Since by assumption π_{M^*} is unique, we have $\Delta_{\min}^{M^*} > 0$. Throughout, we replace dependence on M^* with “ * ” when the meaning is clear from context, for example: $\Delta_{\min}^* := \Delta_{\min}^{M^*}$, or $f^*(\pi) := f^{M^*}(\pi)$.

Further Notation. We let $\mathcal{M}^+ = \{M : \Pi \rightarrow \Delta_{\mathcal{R} \times \mathcal{O}} \mid f^M(\pi) \in [0, 1]\}$ denote the space of all possible models M with rewards in \mathcal{R} and $f^M(\pi) \in [0, 1]$. We use $\Delta_{\mathcal{X}}$ to refer to the set of probability distributions over any \mathcal{X} . Throughout, we often abbreviate $\mathbb{E}_{\pi \sim p}[\cdot]$ with $\mathbb{E}_p[\cdot]$.

1.2. Background: Asymptotic Instance-Optimality

Our aim is to develop algorithms that are *instance-optimal* in a strong sense: for *every model* $M^* \in \mathcal{M}$, the regret of the algorithm under M^* is at least as good as that of any *consistent* algorithm; here, an algorithm is said to be “consistent” if it ensures that $\mathbb{E}^M[\mathbf{Reg}(T)] = o(T)$ for all $M \in \mathcal{M}$. Instance-optimality is a powerful notion of performance: no algorithm—even one designed specifically with M^* in mind—can achieve lower regret on M^* without giving up consistency. For multi-armed bandits, a long line of work initiated by [Lai and Robbins \(1985\)](#) characterizes the instance-optimal regret as a function of the instance M^* , and provides efficient algorithms that attain instance-optimality in finite time ([Garivier et al., 2016](#); [Kaufmann et al., 2016](#); [Lattimore, 2018](#); [Garivier et al., 2019](#)). For the general decision making setting we consider, the forward-looking work of [Graves and Lai \(1997\)](#) (see [Dong and Ma \(2022\)](#) for a contemporary treatment) introduced a complexity measure we refer to as the *Graves-Lai Coefficient*, which asymptotically characterizes the instance-optimal performance as a function of the instance M^* and model class \mathcal{M} . For any class \mathcal{M} and model $M \in \mathcal{M}$, the Graves-Lai Coefficient is defined as

$$\text{glc}(\mathcal{M}, M) := \inf_{\eta \in \mathbb{R}_+^{\Pi}} \left\{ \sum_{\pi \in \Pi} \eta_{\pi} \Delta^M(\pi) \mid \forall M' \in \mathcal{M}^{\text{alt}}(M) : \sum_{\pi \in \Pi} \eta_{\pi} D_{\text{KL}}(M(\pi) \parallel M'(\pi)) \geq 1 \right\}, \quad (3)$$

where, for M with unique optimal decision π_M , we define $\mathcal{M}^{\text{alt}}(M) := \{M' \in \mathcal{M} \mid \pi_M \notin \pi_{M'}\}$ the set of “alternative” models—the models $M' \in \mathcal{M}$ that disagree with M on the optimal decision¹—and $D_{\text{KL}}(\cdot \parallel \cdot)$ denotes the Kullback-Leibler divergence. When \mathcal{M} is clear from context, we will abbreviate $\mathbf{g}^M := \text{glc}(\mathcal{M}, M)$ (and $\mathbf{g}^* := \text{glc}(\mathcal{M}, M^*)$). We also denote any solution to [Eq. \(3\)](#) by η^M —note that this is not in general unique. The characterization of [Graves and Lai \(1997\)](#) is as follows.

1. For $M \in \mathcal{M}$ such that π_M is not unique, we define $\mathcal{M}^{\text{alt}}(M) := \{M' \in \mathcal{M} \mid \pi_M \cap \pi_{M'} = \emptyset\}$, and define $\text{glc}(\mathcal{M}, M)$ as in [Eq. \(3\)](#), with respect to this $\mathcal{M}^{\text{alt}}(M)$. We also define $\mathcal{M}^{\text{alt}}(\pi) = \{M \in \mathcal{M} \mid \pi \notin \pi_M\}$.

Proposition 1.1 (Graves and Lai (1997); Dong and Ma (2022)). *For any model class \mathcal{M} with $|\Pi| < \infty$, any algorithm that is consistent with respect to \mathcal{M} must have*

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \geq \text{glc}(\mathcal{M}, M^*) \cdot \log(T) - o(\log(T)) \quad (4)$$

for any $M^* \in \mathcal{M}$, and there exists an algorithm which achieves, for all $M^* \in \mathcal{M}$ satisfying Assumption 1.3,²

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq \text{glc}(\mathcal{M}, M^*) \cdot \log(T) + o(\log(T)). \quad (5)$$

The interpretation of the Graves-Lai Coefficient of M^* with respect to \mathcal{M} , $\text{glc}(\mathcal{M}, M^*)$, is simple. It asks, if M^* is known to the learner (to be clear, M^* is not known a-priori), what is the minimum regret that must be incurred to gather enough information to rule out all possible alternatives $M' \in \mathcal{M}$ which do not have π_{M^*} as an optimal decision (i.e., $\pi_{M^*} \notin \pi_{M'}$)? In other words, it aims to *certify* that π_{M^*} is indeed the optimal decision while incurring the minimum regret possible.

The Graves-Lai Coefficient characterization is appealing in its simplicity, but the catch—at least when one moves beyond finite-armed bandits—is hiding in the lower-order terms, particularly for the upper bound (5). For general model classes, the best known finite-time regret bounds (Dong and Ma, 2022) take the form

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq \text{glc}(\mathcal{M}, M^*) \cdot \log(T) + \text{poly}(|\Pi|, (\Delta_{\min}^{M^*})^{-1}) \cdot \log^{1-c}(T) \quad (6)$$

where $c > 0$ is a universal constant. While this indeed leads to instance-optimality as $T \rightarrow \infty$, the “lower-order” term in Eq. (6) scales with the size of the decision space, which is intractably large for most problems of interest. As an example, consider the problem of tabular reinforcement learning in an episodic MDP with S states, A actions, and horizon H . Here, we typically have $\text{glc}(\mathcal{M}, M^*) = \text{poly}(S, A, H)$, yet $|\Pi| = A^{HS}$. Consequently, the Graves-Lai Coefficient does not become the dominant term in (6) until $\geq \exp(\exp(S))$. That is, for realistic time horizons, asymptotic instance-optimality does not tell the full story.

Learning an optimal allocation. Given knowledge of an optimal Graves-Lai allocation η^{M^*} solving Eq. (3), a learner could simply take actions as specified by η^{M^*} , and would achieve the instance-optimal rate given in Proposition 1.1. However, this is typically infeasible, as the optimal allocation itself depends strongly upon the ground truth model M^* , and is therefore unknown to the learner. In light of this challenge, the approach taken by essentially all existing algorithms (Burnetas and Katehakis, 1996; Graves and Lai, 1997; Magureanu et al., 2014; Komiyama et al., 2015; Lattimore and Szepesvari, 2017; Combes et al., 2017; Hao et al., 2019, 2020; Van Parys and Golrezaei, 2020; Degenne et al., 2020b; Tirinzoni et al., 2020; Kirschner et al., 2021; Dong and Ma, 2022) is to first learn an estimate for a Graves-Lai allocation η^{M^*} , and then take actions as specified by this estimate. In addition to being natural, this approach is *necessary* in a certain (weak) sense: for any algorithm that achieves instance-optimality, the expected decision frequencies must converge to an approximately optimal allocation as T grows (cf. Lemma B.1).³

The presence of the lower-order term scaling with $|\Pi|$ in Eq. (5) (and in similar regret bounds from most existing work) reflects the sample complexity required to learn an optimal Graves-Lai allocation

2. To be precise, rather than scaling directly with $\text{glc}(\mathcal{M}, M^*)$, the upper bound given by Dong and Ma (2022) scales with a quantity $\text{glc}_T(\mathcal{M}, M^*)$ such that $\text{glc}_T(\mathcal{M}, M^*) \rightarrow_T \text{glc}(\mathcal{M}, M^*)$.

3. The connection between instance-optimal regret and learning an optimal allocation has many subtleties; we refer ahead to Appendix B for extensive discussion.

through uniform exploration. Specifically, one can estimate an allocation by uniformly exploring the decision space to gather data, and then solving an empirical approximation to the Graves-Lai program (3). Naive exploration of this type inevitably results in $\Omega(|\Pi|)$ sample complexity, and it is natural to ask whether a more deliberate exploration strategy, perhaps by exploiting the structure of \mathcal{M} , could lead to better finite-time regret bounds. For the setting of linear bandits, where $\Pi \subseteq \mathbb{R}^d$ and the mean reward function $\pi \mapsto f^M(\pi)$ is linear, this is indeed the case: a recent line of work (Tirinzoni et al., 2020; Kirschner et al., 2021) provides regret bounds of the form

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq \text{glc}(\mathcal{M}, M^*) \cdot \log(T) + \text{poly}(d, (\Delta_{\min}^{M^*})^{-1}) \cdot \log^{1-c}(T).$$

This bound replaces the size of the decision space in the lower-order term by the dimension d , reflecting the fact that there are only d “effective” directions in which exploration is required. While this is an encouraging start, the techniques used in these works are specialized to linear bandits, and it is unclear how to generalize them beyond this setting.

1.3. A Motivating Example

As discussed in the prequel, for finite-armed bandits and linear bandits, it is possible to achieve instance-optimal regret bounds where the lower-order terms scale with the number of actions A , or dimension d , respectively. Extrapolating, one might be tempted to ask whether we can always learn a near-optimal allocation with sample complexity no larger than, say, the minimax rate for \mathcal{M} . The starting point for our work is to recognize that in general, the answer is no: existing notions of statistical complexity are insufficient to capture the complexity of learning the Graves-Lai allocation in finite time, as illustrated in the following simple example.

Example 1.1 (Searching for an informative arm). Let $A, N \geq 2$ and $\beta \in (0, 1)$ be parameters, and consider the class \mathcal{M} of all models defined as follows. First, $\Pi = [A] \cup \{\pi_i^\circ\}_{i \in [N]}$; decisions in $[A]$ are “bandit arms”, and decisions in $\{\pi_i^\circ\}_{i \in [N]}$ are “informative” (or, revealing) arms. Each model M has a unique optimal decision π_M , and the following structure, with $\mathcal{O} = [A] \cup \{\perp\}$.

- For each bandit arm $k \in [A]$ we have $r \sim \mathcal{N}(f^M(k), 1)$ for $f^M \in [0, 1]$. There are no observations, i.e. $o = \perp$ almost surely.
- All informative arms π_k° give 0 reward almost surely. There exists a unique informative arm $\pi_M^\circ \in \{\pi_i^\circ\}_{i \in [N]}$ associated with M , so that if we play any π_k° , we receive an observation

$$o \sim \begin{cases} \text{Unif}([A]), & \pi_k^\circ \neq \pi_M^\circ, \\ \beta \mathbb{I}_{\pi_M} + (1 - \beta) \text{Unif}([A]), & \pi_k^\circ = \pi_M^\circ. \end{cases}$$

We take \mathcal{M} to consist of all possible models with this structure. The interpretation here is as follows. If one were to ignore the revealing arms $\{\pi_i^\circ\}_{i \in [N]}$, this would be a standard finite-armed bandit problem. In particular, if we were to consider a model M^* with $f^{M^*}(\pi) = \frac{1}{2} + \Delta \mathbb{I}\{\pi = i\}$ for $i \in [A]$, a standard calculation would yield $\text{g}^{M^*} \propto \frac{A}{\Delta}$. However, the presence of the informative arms makes the problem substantially easier. With $\beta = 9/10$ (for concreteness), one can see that for the model M , pulling the informative arm $\pi_{M^*}^\circ$ will give $o = \pi_{M^*}$ with probability at least $9/10$, meaning that we can identify that π_{M^*} is optimal with high probability by pulling $\pi_{M^*}^\circ$ a constant number of times. It follows that the optimal allocation is to ignore the bandit arms and set

$\eta^{M^*}(\pi) \propto \mathbb{I}\{\pi = \pi_{M^*}^\circ\}$. This gives $\mathbf{g}^{M^*} \leq O(1)$, which is substantially better than $\mathbf{g}^{M^*} \propto \frac{A}{\Delta}$ if Δ is small or A is large.

If one only is only concerned with asymptotic rates, this is the end of the story, but for non-asymptotic rates, we need to consider the amount of exploration required to *learn the optimal allocation*. In particular, in order to identify the informative arm $\pi_{M^*}^\circ$, which is necessary to learn the optimal allocation, it is clear that in the worst case, any algorithm needs to try all of the revealing arms, leading to $\mathbb{E}[\mathbf{Reg}(T)] = \Omega(N)$. While the complexity of learning the optimal allocation is washed away by an asymptotic analysis with $T \rightarrow \infty$, it cannot be ignored for finite T . In addition, the $\Omega(N)$ factor cannot be explained away by standard complexity measures. As we have seen, $\sup_{M \in \mathcal{M}} \mathbf{g}^M = O(1)$, so the Graves-Lai Coefficient is not sufficient to explain it. Furthermore, the minimax rate for this problem is always bounded by $O(\sqrt{AT})$, which does not scale with N ; yet $\Omega(A)$ sample complexity does not suffice to learn an optimal allocation. In addition, it can be shown that existing complexity measures such as the Decision-Estimation Coefficient (Foster et al., 2021) and information ratio (Russo and Van Roy, 2018) also do not scale with N . \triangleleft

Example 1.1 shows that if we want to achieve instance-optimality in finite time, new notions of problem complexity for the class \mathcal{M} are required, motivating the following central questions:

1. Can we develop algorithms for general model classes \mathcal{M} that achieve non-asymptotic instance-optimal regret bounds of the form

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq \text{glc}(\mathcal{M}, M^*) \cdot \log(T) + \text{comp}(\mathcal{M}) \cdot \log^{1-c}(T), \quad (7)$$

where $\text{comp}(\mathcal{M})$ is a complexity measure that reflects the intrinsic difficulty of exploring in order to learn a Graves-Lai optimal allocation for \mathcal{M} ?

2. Can we understand when the presence of such lower-order terms is *necessary*?

1.4. Organization

The remainder of the paper is organized as follows. In Section 2 we introduce a novel complexity measure, the Allocation-Estimation Coefficient, which captures the complexity of learning a Graves-Lai optimal allocation (Section 2.1), present our main upper and lower bounds (Section 2.2), and instantiate our bounds on several examples (Section 2.3). In Section 3 we present an overview of our main algorithm, AE², and in Section 4 offer directions for future work. Due to space constraints, results in the main body are presented informally—see Part I of the appendix for full statements.

2. Overview of Results

To capture the statistical complexity of learning an optimal Graves-Lai allocation in finite time, we provide a new complexity measure, the *Allocation-Estimation Coefficient* (AEC).

2.1. The Allocation-Estimation Coefficient

For a model $M \in \mathcal{M}$ and parameter $\varepsilon \in [0, 1]$, we define

$$\Lambda(M; \varepsilon) = \left\{ \lambda \in \Delta_\Pi \mid \exists n \in \mathbb{R}_+ \text{ s.t. } \mathbb{E}_{\pi \sim \lambda}[\Delta^M(\pi)] \leq \frac{(1 + \varepsilon)\mathbf{g}^M}{n}, \right. \\ \left. \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \mathbb{E}_{\pi \sim \lambda}[D_{\text{KL}}(M(\pi) \parallel M'(\pi))] \geq \frac{1 - \varepsilon}{n} \right\} \quad (8)$$

the set of (normalized) allocations $\lambda \in \Delta_\Pi$ which are ε -optimal for the Graves-Lai program $\text{glc}(\mathcal{M}, \bar{M})$ in Eq. (3)—both in terms of achieving the optimal objective value and satisfying the information constraint. In addition, for a distribution $\lambda \in \Delta_\Pi$, we define

$$\mathcal{M}_\varepsilon^{\text{gl}}(\lambda) = \{M \in \mathcal{M} \mid \lambda \in \Lambda(M; \varepsilon)\}. \quad (9)$$

Informally, $\mathcal{M}_\varepsilon^{\text{gl}}(\lambda)$ represents the set of models for which the (normalized allocation) $\lambda \in \Delta_\Pi$ is ε -optimal for the Graves-Lai program $\text{glc}(\mathcal{M}, \bar{M})$.

For a *reference model* $\bar{M} : \Pi \rightarrow \Delta_{\mathcal{R} \times \mathcal{O}}$ (not necessarily in \mathcal{M}) and parameter $\varepsilon > 0$, the Allocation-Estimation Coefficient is given by

$$\text{aec}_\varepsilon(\mathcal{M}, \bar{M}) = \inf_{\lambda, \omega \in \Delta_\Pi} \sup_{M \in \mathcal{M} \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)} \left\{ \frac{1}{\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))]} \right\}, \quad (10)$$

where we adopt the convention that the value is 0 if $\mathcal{M}_\varepsilon^{\text{gl}}(\lambda) = \mathcal{M}$. In addition, letting $\text{co}(\mathcal{M})$ denote the convex hull for \mathcal{M} , we define $\text{aec}_\varepsilon(\mathcal{M}) := \sup_{\bar{M} \in \text{co}(\mathcal{M})} \text{aec}_\varepsilon(\mathcal{M}, \bar{M})$.

The Allocation-Estimation Coefficient is a game between a min-player choosing $\lambda, \omega \in \Delta_\Pi$ and a max-player choosing a model $M \in \mathcal{M}$ (with the restriction that $M \notin \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)$). The distribution $\lambda \in \Delta_\Pi$ represents a normalized Graves-Lai allocation, while $\omega \in \Delta_\Pi$ is an *exploration* distribution used to gather information. The reference model \bar{M} should be interpreted as a guess for the underlying $M^* \in \mathcal{M}$. When $\lambda \in \Delta_\Pi$ is fixed, $\inf_{\omega \in \Delta_\Pi} \sup_{M \in \mathcal{M} \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)} (\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))])^{-1}$ represents the time required to gather enough information to distinguish between the reference model \bar{M} and all alternative models $M \notin \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)$ for which λ is not an ε -optimal Graves-Lai allocation—provided that we explore optimally by minimizing over $\omega \in \Delta_\Pi$. For intuition, consider the case when $\bar{M} \in \mathcal{M}$. In this case, λ *must* be chosen so that $\bar{M} \in \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)$ (i.e., λ must be a Graves-Lai optimal allocation for \bar{M}), as otherwise the value of the AEC will be infinite, since $\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \parallel \bar{M}(\pi))] = 0$. Therefore, in such cases, the AEC reflects the *difficulty of distinguishing \bar{M} from models that have different Graves-Lai optimal allocations*. Such models might have the *same optimal decision* π_M as \bar{M} (cf. [Example 1.1](#)) but, if our goal is to play a Graves-Lai optimal allocation for \bar{M} , we must still distinguish \bar{M} from such models.

The Allocation-Estimation Coefficient plays a natural role for deriving both upper and lower bounds on the time required to learn an optimal allocation. For lower bounds, the significance of the AEC is somewhat immediate: it precisely quantifies the time required to acquire enough information to learn an ε -optimal allocation for the *best possible exploration strategy*, and thus leads to a lower bound on time required to learn such an allocation for any algorithm. Notably, the AEC serves as a lower bound for *all possible model classes* \mathcal{M} , and hence may be thought of as an intrinsic structural property of the class \mathcal{M} . For upper bounds, the Allocation-Estimation Coefficient acts as a mechanism to drive exploration; see [Section 3](#) for further explanation.

Generalized Allocation-Estimation Coefficient. For certain results, we make use of the following, slightly more general variant of the AEC. For a reference model $\bar{M} : \Pi \rightarrow \Delta_{\mathcal{R} \times \mathcal{O}}$ and *subset of models* $\mathcal{M}_0 \subseteq \mathcal{M}$, we define

$$\text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}_0, \bar{M}) = \inf_{\lambda, \omega \in \Delta_\Pi} \sup_{M \in \mathcal{M}_0 \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)} \left\{ \frac{1}{\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))]} \right\}, \quad (11)$$

where we adopt the convention that the value is 0 if $\mathcal{M}_0 \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda) = \emptyset$. Here \mathcal{M}_0 denotes the set we take the supremum over, while \mathcal{M} denotes the set that $\mathcal{M}_\varepsilon^{\text{gl}}(\lambda)$ is defined with respect to (i.e., the set with respect to which the Graves-Lai allocation is defined). When $\mathcal{M}_0 = \mathcal{M}$, we recover the AEC as defined in Eq. (10): $\text{aec}_\varepsilon(\mathcal{M}, \bar{M}) = \text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}, \bar{M})$.

2.2. Main Results

Building on the intuition above, our main results show that boundedness of the Allocation-Estimation Coefficient is sufficient to achieve instance-optimality in finite time, and is also necessary in order to learn a near-optimal allocation. Formal statements of our upper bounds are given in Appendix A and formal statements for our lower bounds in Appendix B.

Upper Bound. Our upper bounds are based on a new algorithm, AE^2 (*Allocation Estimation via Adaptive Exploration*), achieves instance-optimality by using the Allocation-Estimation Coefficient to drive exploration.

Theorem 2.1 (Upper Bound—Informal Version of Theorem A.1). *For any model class \mathcal{M} satisfying certain regularity conditions, the AE^2 algorithm ensures that for all $\varepsilon > 0$, $M^* \in \mathcal{M}$, and $T \in \mathbb{N}$:*

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq (1 + \varepsilon) \cdot \text{glc}(\mathcal{M}, M^*) \cdot \log(T) + \tilde{O}^+(\text{aec}_{\varepsilon/12}(\mathcal{M}) + \text{aec}_{\varepsilon/12}^{1/2}(\mathcal{M}) \cdot \log^{1/2}(T)),$$

where $\tilde{O}^+(\cdot)$ suppresses polynomial dependence on ε^{-1} , the log-covering number of \mathcal{M} , $\sup_{M \in \mathcal{M}} 1/\Delta_{\min}^M$, $\log \log T$, and several other measures of the regularity for the class \mathcal{M} .

Theorem 2.1 shows that it is therefore possible to achieve instance-optimality in finite time with lower-order terms scaling (primarily) as the cost of learning the optimal allocation, as captured by the AEC. For multi-armed bandits with $\Pi = [A]$, we have $\text{aec}_\varepsilon(\mathcal{M}) = \tilde{O}^+(\text{poly}(A))$, and for linear bandits with $\Pi \subseteq \mathbb{R}^d$, we have $\text{aec}_\varepsilon(\mathcal{M}) = \tilde{O}^+(\text{poly}(d))$, so that the regret bound in Theorem 2.1 enjoys similar scaling as existing non-asymptotic approaches (Tirinzoni et al., 2020; Kirschner et al., 2021). For tabular reinforcement learning, we have $\text{aec}_\varepsilon(\mathcal{M}) = \tilde{O}^+(\text{poly}(H, S, A))$, which leads to exponential improvement over prior work (Dong and Ma, 2022). Finally, for the instance in Example 1.1, in cases when $N \gg A$, $1/\Delta_{\min}$, $\text{aec}_\varepsilon(\mathcal{M}) = O(N)$, so $\text{aec}_\varepsilon(\mathcal{M})$ captures the intuitive difficulty of learning a Graves-Lai allocation in this setting.

Remark 2.1 (Technical Conditions). *The technical conditions under which Theorem 2.1 is proven are relatively mild, and include certain smoothness of the KL divergences, sub-Gaussian tail behavior for log-likelihood ratios, and bounded covering number for \mathcal{M} with standard parametric growth (note that \mathcal{M} may be infinite), all of which can be shown to hold for standard classes. In addition, we require that the amount of information that can be gained by playing the optimal decision for M^* is bounded (see Appendix A.1 for precise statements of our conditions).*

Remark 2.2 (Asymptotic Performance). *Asymptotically as $T \rightarrow \infty$, the regret bound given in Theorem 2.1 scales as $(1 + \varepsilon) \cdot \text{glc}(\mathcal{M}, M^*) \cdot \log T$, which is a factor of $(1 + \varepsilon)$ off of the lower bound. For all standard classes, $\text{aec}_{\varepsilon/12}(\mathcal{M})$ scales polynomially in $1/\varepsilon$ so, to obtain an optimal leading-order constant, it suffices to choose $\varepsilon = 1/\log^a T$, for small enough $a > 0$.*

Adapting to Minimum Gap. Note that the lower-order term given in [Theorem 2.1](#) scales with $\sup_{M \in \mathcal{M}} 1/\Delta_{\min}^M$, the minimum gap of the entire model class. In [Appendix A.5](#), we give a refinement of the AE^2 algorithm (AE_*^2) which attains an improved regret bound which replaces the term $\sup_{M \in \mathcal{M}} 1/\Delta_{\min}^M$ with $1/\Delta_{\min}^*$, the minimum gap of the underlying model; notably AE_*^2 requires no prior knowledge of Δ_{\min}^* (i.e., it is able to adapt to the minimum gap of the underlying model). In addition, rather than scaling with $\text{aec}_{\varepsilon/12}(\mathcal{M})$, the lower-order term now scales with $\text{aec}_{\varepsilon/12}^{\mathcal{M}}(\mathcal{M}^*)$ for a subset $\mathcal{M}^* \subset \mathcal{M}$ which, informally, restricts to models in \mathcal{M} for which the minimum gap is at least Δ_{\min}^* .

Theorem 2.2 (Upper Bound—Informal Version of [Theorem A.2](#)). *For any model class \mathcal{M} satisfying certain regularity conditions, the AE_*^2 algorithm ensures that for all $\varepsilon > 0$, $M^* \in \mathcal{M}$, and $T \in \mathbb{N}$:*

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq (1 + \varepsilon) \cdot \text{glc}(\mathcal{M}, M^*) \cdot \log(T) + \tilde{O}^+((\text{aec}_{\varepsilon/12}^{\mathcal{M}}(\mathcal{M}^*))^3 + \log^{6/7}(T)), \quad (12)$$

where $\tilde{O}^+(\cdot)$ suppresses polynomial dependence on ε^{-1} , the log-covering number of \mathcal{M} , $1/\Delta_{\min}^*$, $\log \log T$, and several other measures of the regularity for the class \mathcal{M} .

Lower Bound. To provide lower bounds, we adopt a novel minimax framework which asks, for the model class \mathcal{M} under consideration, what is the least value of $T \in \mathbb{N}$ for which it is possible to learn an ε -optimal Graves-Lai allocation for any model in \mathcal{M} . To state our result, we introduce the following notation, defined with respect to any $\bar{M} \in \mathcal{M}^+$:

$$\mathcal{M}^{\text{opt}}(\bar{M}) = \{M \in \mathcal{M} \mid \pi_M \subseteq \pi_{\bar{M}}, D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi)) = 0 \ \forall \pi \in \pi_{\bar{M}}\}.$$

The set $\mathcal{M}^{\text{opt}}(\bar{M})$ represents the set of models where 1) the optimal decision coincides with that of \bar{M} and 2) \bar{M} and $M \in \mathcal{M}$ cannot be distinguished by playing the optimal decision.

Our main lower bound provides a sort of converse to the upper bound in [Theorem 2.1](#).

Theorem 2.3 (Lower Bound—Informal Version of [Theorem B.2](#)). *For any model class \mathcal{M} and $\varepsilon > 0$, it holds that unless*

$$\log(T) \geq \sup_{\bar{M} \in \mathcal{M}^+} \tilde{\Omega}^+(\text{aec}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}^{\text{opt}}(\bar{M}), \bar{M})), \quad (13)$$

no algorithm can simultaneously achieve the following for all instances $M \in \mathcal{M}$:

1. attain Graves-Lai optimality on M within a constant factor (i.e., ensure $\mathbb{E}^M[\mathbf{Reg}(T)] \leq 2 \cdot \text{glc}(\mathcal{M}, M) \log(T)$).
2. discover an ε -optimal allocation for M (i.e., find λ with $M \in \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda)$) with probability greater than $\tilde{\Omega}^+(1)$.

Here, $\tilde{\Omega}^+(\cdot)$ hides polynomial dependence on regularity parameters of \mathcal{M} .

Observe that the Graves-Lai Coefficient becomes the dominant term in the upper bound [Theorem 2.1](#) as soon as $\log(T) \geq \tilde{\Omega}^+(\text{aec}_{\varepsilon/12}(\mathcal{M}))$. The lower bound (13) shows that for any algorithm that aims to estimate the Graves-Lai allocation (in particular, AE^2), such scaling is necessary, and therefore the lower-order term in [Theorem 2.1](#) is in some sense unimprovable. To the best of our knowledge, this is the first general approach to quantifying the lower-order terms necessary in order to achieve instance-optimality. We make several remarks on the lower bound.

Remark 2.3 (Scaling in AEC). *Our upper and lower bounds scale with a slightly different version of the AEC, as the lower bound restricts the AEC to $\mathcal{M}^{\text{opt}}(\bar{M})$. In [Appendix B](#), we show an additional lower bound that scales directly with $\text{aec}_\varepsilon(\mathcal{M})$, matching our upper bound, but which only provides a lower bound on T rather than $\log(T)$ (see [Theorem B.1](#)).*

Remark 2.4 (Asymptotic Optimality and Learning Optimal Allocations). *[Theorem 2.3](#) gives a lower bound on the time needed to learn a near-optimal Graves-Lai allocation, but does not directly imply that it is necessary that an asymptotically optimal algorithm learn such an allocation. As we have noted, the allocations played by any asymptotically optimal algorithm must converge to an optimal allocation in expectation. However, showing that this convergence is necessary with even constant probability (the condition under which [Theorem 2.3](#) is proved) is rather subtle. As we show in [Appendix B](#) ([Theorem B.3](#)), if one assumes that, in addition to being asymptotically optimal in expectation, the algorithm under consideration also has regret with appropriately bounded second moment, then if $\mathbb{E}^M[\mathbf{Reg}(T)] \leq (1 + \varepsilon) \cdot \text{glc}(\mathcal{M}, M) \log(T)$ for all $M \in \mathcal{M}$, it is indeed necessary that a burn-in time analogous to [Eq. \(13\)](#) is satisfied.*

Together, our upper and lower bounds represent an initial step toward building a sharp non-asymptotic theory of instance-optimality, and lead to a number of new conceptual insights. Our results open the door for further-development, and to this end we highlight a number of opportunities for improvement ([Appendix B.4](#)), as well as open problems ([Section 4](#)).

2.3. Concrete Examples

We next present several examples illustrating our upper and lower bounds. All results in this section are informal—see [Appendices A.3, A.6, A.7](#) and [B.3](#) for formal results and additional examples.

Example 2.1 (Searching for an Informative Arm (revisited)). We return to the example of [Section 1.3](#). Some calculation shows that, for the choice of \mathcal{M} in [Example 1.1](#), as long as β is constant and $N \geq A/\Delta^2$, we have

$$\Omega(N) \leq \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}^{\text{opt}}(\bar{M}), \bar{M}) \quad \text{and} \quad \text{aec}_\varepsilon(\mathcal{M}) \leq O(N).$$

[Theorem 2.1](#) then implies that AE^2 has regret on [Example 1.1](#) of

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq (1 + \varepsilon) \cdot g^* \log(T) + N \cdot \text{poly}\left(A, \frac{1}{\Delta}, \frac{1}{\varepsilon}, \log N, \log \log T\right) \cdot \log^{1/2}(T).$$

Furthermore, [Theorem 2.3](#) shows that a scaling of $\log(T) \geq \tilde{\Omega}^+(N)$ is necessary for any algorithm to learn a Graves-Lai optimal allocation. It follows that, on this example, the AEC reflects a notion of problem difficulty not captured by any existing complexity measure, matching our intuitive understanding of what the correct scaling should be. We remark that the scaling $\log(T) \geq \tilde{\Omega}^+(N)$ is natural (as compared to $T \geq \tilde{\Omega}^+(N)$) since, if an algorithm is instance-optimal as required by [Theorem 2.3](#), it can allocate at most $O^+(\log(T))$ pulls to suboptimal decisions. To pull every informative arm (each of which is suboptimal) while achieving instance-optimality, it follows that we must have $\log(T) \geq \tilde{\Omega}^+(N)$. ◀

Example 2.2 (Tabular Reinforcement Learning). Consider the setting of tabular reinforcement learning. Here we take M to be a (tabular) episodic Markov Decision Processes (MDP) with S

states, A actions, horizon H , probability transition kernels $\{P_h^M\}_{h=1}^H$, and Gaussian rewards; see [Appendix A.7](#) for a full definition of this setting. Let \mathcal{M} denote the set of all such tabular MDPs which, for each state-action-state triple (s, a, s') and $h \in [H]$, have $P_h^M(s' | s, a) \geq P_{\min} > 0$; that is, each transition can occur with some minimum probability. Then it can be shown that:

$$\text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}^*) \leq \text{poly}\left(S, A, H, \frac{1}{\varepsilon}, \frac{1}{\Delta_{\min}^*}, \log \frac{1}{P_{\min}}\right).$$

This implies that for any tabular MDP in \mathcal{M} , the AE_*^2 algorithm has regret bounded as:

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq (1 + \varepsilon) \cdot \mathbf{g}^* \log(T) + \text{poly}\left(S, A, H, \frac{1}{\Delta_{\min}^*}, \frac{1}{\varepsilon}, \log \frac{1}{P_{\min}}, \log \log T\right) \cdot \log^{1/2}(T).$$

To the best of our knowledge, this is the first regret bound in the setting of tabular reinforcement learning which is instance-optimal with lower-order terms scaling only polynomially in problem parameters, an exponential improvement over past work ([Ok et al., 2018](#); [Dong and Ma, 2022](#)). Furthermore, one can also show that $\sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}^{\text{opt}}(\bar{M}), \bar{M}) \geq \tilde{\Omega}\left(\frac{1}{\varepsilon^2} \cdot \frac{SA}{(\Delta_{\min}^*)^2}\right)$, so that our lower bound, [Theorem 2.3](#), implies that a burn-in time scaling polynomially in $S, A, \frac{1}{\varepsilon}$, and $\frac{1}{\Delta_{\min}^*}$ is necessary to learn an ε -optimal Graves-Lai allocation for every model in \mathcal{M} .

We remark that the prior work of [Dong and Ma \(2022\)](#) does not require that $P_h^M(s' | s, a) \geq P_{\min}$ as we do, yet their bound scales *polynomially* in the inverse probability of observing the trajectory that occurs with minimum non-zero probability (the work of [Ok et al. \(2018\)](#) only holds for ergodic MDPs, itself a very strong assumption). Our finite-time results are therefore, in general, significantly stronger, scaling only logarithmically in P_{\min} . Removing the P_{\min} assumption while still achieving reasonable lower-order terms is an interesting direction for future work. \triangleleft

Example 2.3 (Linear Bandits). Consider the setting of linear bandits in d dimensions with unit-variance Gaussian noise. Let \mathcal{M} denote the set of all linear bandit models defined with respect to some arm set $\mathcal{X} \subseteq \mathbb{R}^d$ and parameter set $\Theta \subseteq \mathbb{R}^d$. Concretely, each model $M \in \mathcal{M}$ takes the form

$$M(\pi) = \mathcal{N}(\langle \theta, x_\pi \rangle, 1),$$

for some $\theta \in \Theta$, where $x_\pi \in \mathcal{X}$ is an embedding of π . Let Δ_{\min}^* denote the minimum gap of M^* (which is unknown to the algorithm). Then it can be shown that

$$\text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}^*) \leq \text{poly}\left(d, \frac{1}{\varepsilon}, \frac{1}{\Delta_{\min}^*}\right)$$

which implies that the AE_*^2 algorithm has regret bounded as

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq (1 + \varepsilon) \cdot \mathbf{g}^* \log(T) + \text{poly}\left(d, \frac{1}{\Delta_{\min}^*}, \frac{1}{\varepsilon}, \log \log T\right) \cdot \log^{6/7}(T). \quad (14)$$

We remark that the scaling of [Eq. \(14\)](#) matches the state-of-the-art instance-optimal bounds for linear bandits (in that all have polynomial lower-order terms—our polynomial dependence on d is slightly worse as our upper bound on the AEC is somewhat coarse) ([Tirinzi et al., 2020](#); [Kirschner et al., 2021](#)). Notably, it is a simple corollary of a much more general result, while prior work relies on specialized algorithms tailored to linear bandits. \triangleleft

In [Appendices A.3, A.6, A.7](#) and [B.3](#) we formalize these examples and present additional examples, including structured bandits with bounded eluder dimension and finite-action contextual bandits. In all cases, we obtain lower-order terms scaling only polynomially with problem parameters, and in each setting either match the best-known existing bound, or are the first to provide any meaningful finite-time bounds.

Algorithm 1 Allocation Estimation via Adaptive Exploration (AE², Informal)

- 1: **input:** optimality tolerance ε , model class \mathcal{M} .
- 2: Initialize $s \leftarrow 1$ and $q \leftarrow$ class-dependent quantity.
- 3: Compute $\xi^1 \leftarrow \mathbf{Alg}_{\text{KL}}(\{\emptyset\})$ and $\widehat{M}^1 \leftarrow \mathbb{E}_{M \sim \xi^1}[M]$.
- 4: **for** $t = 1, 2, 3, \dots$ **do**
- 5: **if** $\exists \pi_{\widehat{M}^s} \in \pi_{\widehat{M}^s}$ s.t. $\forall M \in \mathcal{M}^{\text{alt}}(\pi_{\widehat{M}^s}), \sum_{i=1}^{s-1} \mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}_{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}_{M, \pi^i}(r^i, o^i)} \right] \geq \log(t \log t)$ **then**
- 6: Play $\pi_{\widehat{M}^s}$. // Exploit
- 7: **else** // Explore
- 8: Set $p^s \leftarrow q\lambda^s + (1 - q)\omega^s$ for

$$\lambda^s, \omega^s \leftarrow \arg \min_{\lambda, \omega \in \Delta_{\Pi}} \sup_{M \in \mathcal{M} \setminus \mathcal{M}_{\varepsilon/6}^{\text{gl}}(\lambda)} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\widehat{M}(\pi) \| M(\pi))]]}. \quad (15)$$

- 9: Draw $\pi^s \sim p^s$ and observe reward r^s and observation o^s .
 - 10: Compute $\widehat{M}^{s+1} = \mathbb{E}_{\widehat{M} \sim \xi^{s+1}}[\widehat{M}]$ for $\xi^{s+1} \leftarrow \mathbf{Alg}_{\text{KL}}(\{(\pi^i, r^i, o^i)\}_{i=1}^s)$, $s \leftarrow s + 1$.
-

3. Algorithm Overview

Finally, we present our algorithm, AE², in [Algorithm 1](#). AE² relies on an *online estimation oracle*, denoted by \mathbf{Alg}_{KL} , which at every step s , given access to data $\{(\pi^i, r^i, o^i)\}_{i=1}^{s-1}$ with $\pi^i \sim p^i$ and $(r^i, o^i) \sim M^*(\pi^i)$ returns a randomized estimate $\xi^s = \mathbf{Alg}_{\text{KL}}(\{(\pi^i, r^i, o^i)\}_{i=1}^{s-1}) \in \Delta_{\mathcal{M}}$ with the goal of approximating M^* ([Foster and Rakhlin, 2020](#); [Foster et al., 2021](#)). The estimates produced by \mathbf{Alg}_{KL} must ensure the total KL estimation error is bounded:

$$\mathbf{Est}_{\text{KL}}(s) := \sum_{i=1}^s \mathbb{E}_{\widehat{M} \sim \xi^i} [\mathbb{E}_{\pi \sim p^i} [D_{\text{KL}}(\widehat{M}(\pi) \| M^*(\pi))]] \lesssim O(\log s). \quad (16)$$

We show that, under the regularity conditions required by [Theorem 2.1](#), such a guarantee can be achieved, with $O(\cdot)$ hiding the log-covering number of \mathcal{M} .

AE² alternates between *exploit* steps and *explore* steps, tracking the number of explore steps that have been performed with a counter $s \in \mathbb{N}$. For each step $t \in \mathbb{N}$, the algorithm makes use of an estimator $\widehat{M}^s = \mathbb{E}_{\widehat{M} \sim \xi^s}[\widehat{M}]$, where $\xi^s = \mathbf{Alg}_{\text{KL}}(\{(\pi^i, r^i, o^i)\}_{i=1}^{s-1})$ is computed by calling the estimation oracle with data gathered at previous explore steps. Given the estimator, the algorithm first checks whether it has enough information to guarantee that the greedy decision is optimal, in which case it exploits ([Line 6](#)); otherwise it explores. In explore steps, the key component is the choice of the exploration distributions λ^s and ω^s in [Eq. \(15\)](#), which mimics the Allocation-Estimation Coefficient program. To understand the role of the Allocation-Estimation Coefficient here, we consider two cases. In the first case, if λ^s is an ε -optimal Graves-Lai allocation for M^* (that is, $M^* \in \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda^s)$), then playing λ^s will optimize the tradeoff between minimizing regret and collecting information, and will therefore match the optimal performance prescribed by the Graves-Lai Coefficient.

In the second case, if λ^s is not an ε -optimal Graves-Lai allocation for M^* , we have $M^* \notin \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda^s)$, so ω^s will place mass on actions that ensure $\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\widehat{M}(\pi) \| M^*(\pi))]]$ is large. Since p^s plays ω^s with constant probability, the quantity $\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M^*(\pi))]]$ will also be large. If our estimator is consistent and [Eq. \(16\)](#) holds, this can only happen a small number of times without violating [Eq. \(16\)](#); at most logarithmic in the number of exploration rounds. As such,

we can show that λ^s must be a near-optimal Graves-Lai allocation for M^* on all but a logarithmic number of exploration rounds, and that AE^2 achieves the optimal rate on such rounds, yielding the optimal performance rate of [Theorem 2.1](#). Critically, rather than exploring in a naive fashion (e.g., by sampling decisions uniformly), AE^2 explores specifically with the goal of learning a Graves-Lai optimal allocation for M^* and adapts to the structure of \mathcal{M} to perform this exploration efficiently.

We emphasize the simplicity of AE^2 . While most existing instance-optimal algorithms are quite complicated even for basic settings, AE^2 relies on a few simple components yet is far more general than existing approaches and performs comparably or better. See [Appendix A.2](#) for a full description.

4. Discussion

Our work initiates the systematic study of non-asymptotic instance-optimality in interactive decision making. We close by highlighting a number of interesting open problems and future directions raised by our work. On the technical side:

- Our upper bounds depend on a number of different problem-dependent parameters, such as the minimum gap in the lower-order terms. Can we improve the dependence on these parameters, or understand to what extent they are necessary?
- Our lower bounds concern the problem of learning a near-optimal allocation. Like the upper bounds, these results are likely loose in terms of dependence on various problem parameters, and new techniques will be required to tighten them. Furthermore, it remains to develop a complete understanding of the connections between this problem and the problem of minimizing regret. While our results show that “well-behaved” algorithms which achieve instance-optimal regret must pay a burn-in proportional to the cost of learning a near-optimal allocation, it remains unclear if this is truly necessary for algorithms which only have optimal expected regret (but, for example, could exhibit heavy-tailed behavior).
- While our algorithm achieves the instance-optimal rate, its regret could scale linearly over shorter time horizons, until it has learned a near-optimal allocation. Can we develop “best-of-both-worlds” algorithms that achieve the same instance-optimal guarantees of AE^2 , yet also achieves the minimax-optimal rate (for example, a $\mathcal{O}(\sqrt{T})$ -style guarantee) over shorter time horizons?
- The focus of this work is primarily on regret minimization, yet the challenge of learning the optimal allocation also arises in the PAC setting. Does the AEC extend to the PAC setting, and can algorithms be developed in the PAC setting which achieve the instance-optimal rate in the leading-order term, while scaling with an AEC-like quantity in the lower-order term?

More broadly, it will be interesting to explore whether our framework and algorithm design ideas can be used to develop practical and computationally efficient algorithms.

Acknowledgements

The authors would like to thank Johannes Kirschner for helpful discussions. The work of AW was supported in part by NSF TRIPODS 62-2945 and NSF HDR 62-0221. A portion of this work was completed while AW was an intern at Microsoft Research, and while visiting the Simons Institute for the Theory of Computing.

References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, 2011.
- Rajeev Agrawal, Demosthenis Teneketzis, and Venkatachalam Anantharam. Asymptotically efficient adaptive allocation schemes for controlled markov chains: Finite parameter space. Technical report, MICHIGAN UNIV ANN ARBOR COMMUNICATIONS AND SIGNAL PROCESSING LAB, 1988.
- Aymen Al Marjani and Alexandre Proutiere. Adaptive sampling for best policy identification in markov decision processes. In *International Conference on Machine Learning*, pages 7459–7468. PMLR, 2021.
- Aymen Al Marjani, Aurélien Garivier, and Alexandre Proutiere. Navigating to the best policy in markov decision processes. *Advances in Neural Information Processing Systems*, 34:25852–25864, 2021.
- Chamy Allenberg, Peter Auer, László Györfi, and György Ottucsák. *Hannan Consistency in On-Line Learning in Case of Unbounded Losses Under Partial Monitoring*, pages 229–243. Springer Berlin Heidelberg, Berlin, Heidelberg, 2006. ISBN 978-3-540-46650-5. doi: 10.1007/11894841_20. URL http://dx.doi.org/10.1007/11894841_20.
- Jean-Yves Audibert and Sébastien Bubeck. Minimax policies for adversarial and stochastic bandits. In *COLT*, volume 7, pages 1–122, 2009.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- Sébastien Bubeck, Michael Cohen, and Yuanzhi Li. Sparsity, variance and curvature in multi-armed bandits. In *Algorithmic Learning Theory*, pages 111–127. PMLR, 2018.
- Apostolos N Burnetas and Michael N Katehakis. Optimal adaptive policies for sequential allocation problems. *Advances in Applied Mathematics*, 17(2):122–142, 1996.
- Fan Chen, Song Mei, and Yu Bai. Unified algorithms for RL with decision-estimation coefficients: No-regret, PAC, and reward-free learning. *arXiv preprint arXiv:2209.11745*, 2022.
- Lijie Chen, Jian Li, and Mingda Qiao. Nearly instance optimal sample complexity bounds for top-k arm selection. In *Artificial Intelligence and Statistics*, pages 101–110. PMLR, 2017.
- Richard Combes, Stefan Magureanu, and Alexandre Proutiere. Minimal exploration in structured stochastic bandits. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 1761–1769, 2017.
- Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit feedback. In *Conference on Learning Theory (COLT)*, 2008.
- Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying pac and regret: Uniform pac bounds for episodic reinforcement learning. *Advances in Neural Information Processing Systems*, 30, 2017.

- Christoph Dann, Lihong Li, Wei Wei, and Emma Brunskill. Policy certificates: Towards accountable reinforcement learning. In *International Conference on Machine Learning*, pages 1507–1516. PMLR, 2019.
- Christoph Dann, Teodor Vanislavov Marinov, Mehryar Mohri, and Julian Zimmert. Beyond value-function gaps: Improved instance-dependent regret bounds for episodic reinforcement learning. *Advances in Neural Information Processing Systems*, 34:1–12, 2021.
- Rémy Degenne and Wouter M Koolen. Pure exploration with multiple correct answers. *Advances in Neural Information Processing Systems*, 32, 2019.
- Rémy Degenne, Wouter M Koolen, and Pierre Ménard. Non-asymptotic pure exploration by solving games. *Advances in Neural Information Processing Systems*, 32, 2019.
- Rémy Degenne, Pierre Ménard, Xuedong Shang, and Michal Valko. Gamification of pure exploration for linear bandits. In *International Conference on Machine Learning*, pages 2432–2442. PMLR, 2020a.
- Rémy Degenne, Han Shao, and Wouter Koolen. Structure adaptive algorithms for stochastic bandits. In *International Conference on Machine Learning*, pages 2443–2452. PMLR, 2020b.
- Omar Darwiche Domingues, Pierre Ménard, Emilie Kaufmann, and Michal Valko. Episodic reinforcement learning in finite mdps: Minimax lower bounds revisited. In *Algorithmic Learning Theory*, pages 578–598. PMLR, 2021.
- Kefan Dong and Tengyu Ma. Asymptotic instance-optimal algorithms for interactive decision making. *arXiv preprint arXiv:2206.02326*, 2022.
- Simon S Du, Sham M Kakade, Jason D Lee, Shachar Lovett, Gaurav Mahajan, Wen Sun, and Ruosong Wang. Bilinear classes: A structural framework for provable generalization in RL. *International Conference on Machine Learning*, 2021.
- Tanner Fiez, Lalit Jain, Kevin G Jamieson, and Lillian Ratliff. Sequential experimental design for transductive linear bandits. *Advances in neural information processing systems*, 32, 2019.
- Dylan J Foster and Alexander Rakhlin. Beyond UCB: Optimal and efficient contextual bandits with regression oracles. *International Conference on Machine Learning (ICML)*, 2020.
- Dylan J Foster, Zhiyuan Li, Thodoris Lykouris, Karthik Sridharan, and Eva Tardos. Learning in games: Robustness of fast convergence. *Advances in Neural Information Processing Systems*, 29, 2016.
- Dylan J Foster, Alexander Rakhlin, David Simchi-Levi, and Yunzong Xu. Instance-dependent complexity of contextual bandits and reinforcement learning: A disagreement-based perspective. *Conference on Learning Theory (COLT)*, 2020.
- Dylan J Foster, Sham M Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity of interactive decision making. *arXiv preprint arXiv:2112.13487*, 2021.

- Dylan J Foster, Noah Golowich, Jian Qian, Alexander Rakhlin, and Ayush Sekhari. A note on model-free reinforcement learning with the decision-estimation coefficient. *arXiv preprint arXiv:2211.14250*, 2022a.
- Dylan J Foster, Alexander Rakhlin, Ayush Sekhari, and Karthik Sridharan. On the complexity of adversarial decision making. *arXiv preprint arXiv:2206.13063*, 2022b.
- Dylan J. Foster, Noah Golowich, and Yanjun Han. Tight guarantees for interactive decision making with the decision-estimation coefficient. *arXiv preprint arXiv:2301.08215*, 2023.
- Aurélien Garivier and Emilie Kaufmann. Optimal best arm identification with fixed confidence. In *Conference on Learning Theory*, pages 998–1027, 2016.
- Aurélien Garivier, Tor Lattimore, and Emilie Kaufmann. On explore-then-commit strategies. In *Advances in Neural Information Processing Systems*, pages 784–792, 2016.
- Aurélien Garivier, Pierre Ménard, and Gilles Stoltz. Explore first, exploit next: The true shape of regret in bandit problems. *Mathematics of Operations Research*, 44(2):377–399, 2019.
- Todd L Graves and Tze Leung Lai. Asymptotically efficient adaptive choice of control laws in controlled Markov chains. *SIAM journal on control and optimization*, 35(3):715–743, 1997.
- Botao Hao, Tor Lattimore, and Csaba Szepesvari. Adaptive exploration in linear contextual bandit. *arXiv preprint arXiv:1910.06996*, 2019.
- Botao Hao, Tor Lattimore, and Csaba Szepesvari. Adaptive exploration in linear contextual bandit. In *International Conference on Artificial Intelligence and Statistics*, pages 3536–3545. PMLR, 2020.
- Elad Hazan and Satyen Kale. Better algorithms for benign bandits. *Journal of Machine Learning Research*, 12(4), 2011.
- Kevin Jamieson, Matthew Malloy, Robert Nowak, and Sébastien Bubeck. lil’ucb: An optimal exploration algorithm for multi-armed bandits. In *Conference on Learning Theory*, pages 423–439. PMLR, 2014.
- Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Contextual decision processes with low Bellman rank are PAC-learnable. In *International Conference on Machine Learning*, pages 1704–1713, 2017.
- Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman eluder dimension: New rich classes of RL problems, and sample-efficient algorithms. *Neural Information Processing Systems*, 2021.
- Kwang-Sung Jun and Chicheng Zhang. Crush optimism with pessimism: Structured bandits beyond asymptotic optimality. *Advances in Neural Information Processing Systems*, 33:6366–6376, 2020.
- Sham Machandranath Kakade. *On the sample complexity of reinforcement learning*. University of London, University College London (United Kingdom), 2003.

- Julian Katz-Samuels, Lalit Jain, Kevin G Jamieson, et al. An empirical process approach to the union bound: Practical algorithms for combinatorial and linear bandits. *Advances in Neural Information Processing Systems*, 33:10371–10382, 2020.
- Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best-arm identification in multi-armed bandit models. *The Journal of Machine Learning Research*, 17(1):1–42, 2016.
- Johannes Kirschner, Tor Lattimore, Claire Vernade, and Csaba Szepesvári. Asymptotically optimal information-directed sampling. In *Conference on Learning Theory*, pages 2777–2821. PMLR, 2021.
- Junpei Komiyama, Junya Honda, and Hiroshi Nakagawa. Regret lower bound and optimal algorithm in finite stochastic partial monitoring. *Advances in Neural Information Processing Systems*, 28, 2015.
- Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1):4–22, 1985.
- Tor Lattimore. Refining the confidence level for optimistic bandit strategies. *The Journal of Machine Learning Research*, 19(1):765–796, 2018.
- Tor Lattimore and Csaba Szepesvari. The end of optimism? an asymptotic analysis of finite-armed linear bandits. In *Artificial Intelligence and Statistics*, pages 728–737. PMLR, 2017.
- Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Gábor Lugosi and Shahar Mendelson. Mean estimation and regression under heavy-tailed distributions: A survey. *Foundations of Computational Mathematics*, 19(5):1145–1190, 2019.
- Stefan Magureanu, Richard Combes, and Alexandre Proutiere. Lipschitz bandits: Regret lower bound and optimal algorithms. In *Conference on Learning Theory*, pages 975–999. PMLR, 2014.
- Teodor V Marinov, Mehryar Mohri, and Julian Zimmert. Stochastic online learning with feedback graphs: Finite-time and asymptotic optimality. *arXiv preprint arXiv:2206.10022*, 2022a.
- Teodor Vanislavov Marinov, Mehryar Mohri, and Julian Zimmert. Open problem: Finite-time instance dependent optimality for stochastic online learning with feedback graphs. In *Conference on Learning Theory*, pages 5644–5649. PMLR, 2022b.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- Jungseul Ok, Alexandre Proutiere, and Damianos Tranos. Exploration in structured reinforcement learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- Ian Osband and Benjamin Van Roy. On lower bounds for regret in reinforcement learning. *arXiv preprint arXiv:1608.02732*, 2016.

- Daniel Russo. Simple bayesian algorithms for best arm identification. In *Conference on Learning Theory*, pages 1417–1418, 2016.
- Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic exploration. In *Advances in Neural Information Processing Systems*, pages 2256–2264, 2013.
- Daniel Russo and Benjamin Van Roy. Learning to optimize via information-directed sampling. *Operations Research*, 66(1):230–252, 2018.
- Ohad Shamir. A variant of azuma’s inequality for martingales with subgaussian tails. *arXiv preprint arXiv:1110.2392*, 2011.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484, 2016.
- Max Simchowitz and Kevin G Jamieson. Non-asymptotic gap-dependent regret bounds for tabular mdps. *Advances in Neural Information Processing Systems*, 32, 2019.
- Max Simchowitz, Kevin Jamieson, and Benjamin Recht. The simulator: Understanding adaptive sampling in the moderate-confidence regime. In *Conference on Learning Theory*, pages 1794–1834. PMLR, 2017.
- Wen Sun, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Model-based RL in contextual decision processes: PAC bounds and exponential improvements over model-free approaches. In *Conference on learning theory*, pages 2898–2933. PMLR, 2019.
- Andrea Tirinzoni, Matteo Pirota, Marcello Restelli, and Alessandro Lazaric. An asymptotically optimal primal-dual incremental algorithm for contextual linear bandits. *Advances in Neural Information Processing Systems*, 33:1417–1427, 2020.
- Andrea Tirinzoni, Matteo Pirota, and Alessandro Lazaric. A fully problem-dependent regret lower bound for finite-horizon mdps. *arXiv preprint arXiv:2106.13013*, 2021.
- Andrea Tirinzoni, Aymen Al-Marjani, and Emilie Kaufmann. Near instance-optimal pac reinforcement learning for deterministic mdps. *arXiv preprint arXiv:2203.09251*, 2022.
- Bart PG Van Parys and Negin Golrezaei. Optimal learning for structured bandits. *arXiv preprint arXiv:2007.07302*, 2020.
- Andrew Wagenmaker and Kevin Jamieson. Instance-dependent near-optimal policy identification in linear mdps via online experiment design. *arXiv preprint arXiv:2207.02575*, 2022.
- Andrew Wagenmaker and Aldo Pacchiano. Leveraging offline data in online reinforcement learning. *arXiv preprint arXiv:2211.04974*, 2022.
- Andrew J Wagenmaker, Max Simchowitz, and Kevin Jamieson. Task-optimal exploration in linear dynamical systems. In *International Conference on Machine Learning*, pages 10641–10652. PMLR, 2021.

- Andrew J Wagenmaker, Yifang Chen, Max Simchowitz, Simon Du, and Kevin Jamieson. Reward-free rl is no harder than reward-aware rl in linear markov decision processes. In *International Conference on Machine Learning*, pages 22430–22456. PMLR, 2022a.
- Andrew J Wagenmaker, Max Simchowitz, and Kevin Jamieson. Beyond no regret: Instance-dependent pac reinforcement learning. In *Conference on Learning Theory*, pages 358–418. PMLR, 2022b.
- Ruosong Wang, Russ R Salakhutdinov, and Lin Yang. Reinforcement learning with general value function approximation: Provably efficient approach via bounded eluder dimension. *Advances in Neural Information Processing Systems*, 33, 2020.
- Chen-Yu Wei and Haipeng Luo. More adaptive algorithms for adversarial bandits. In *Conference On Learning Theory*, pages 1263–1291. PMLR, 2018.
- Yuhong Yang and Andrew R Barron. An asymptotic property of model selection criteria. *IEEE Transactions on Information Theory*, 44(1):95–116, 1998.
- Andrea Zanette and Emma Brunskill. Tighter problem-dependent regret bounds in reinforcement learning without domain knowledge using value function bounds. In *International Conference on Machine Learning*, pages 7304–7312. PMLR, 2019.
- Zihan Zhang, Xiangyang Ji, and Simon Du. Is reinforcement learning more difficult than bandits? a near-optimal algorithm escaping the curse of horizon. In *Conference on Learning Theory*, pages 4528–4531. PMLR, 2021.
- Dongruo Zhou, Quanquan Gu, and Csaba Szepesvari. Nearly minimax optimal reinforcement learning for linear mixture markov decision processes. In *Conference on Learning Theory*, pages 4532–4576. PMLR, 2021.

Contents

I	Main Results	23
A	The AE^2 Algorithm: Regret Bounds and Examples	23
A.1	Regularity Conditions	23
A.2	The AE^2 Algorithm	25
A.3	AE^2 Algorithm: Regret Bound for Uniformly Regular Classes	28
A.4	The AE^2_\star Algorithm	31
A.5	AE^2_\star Algorithm: Regret Bound without Uniform Regularity	33
A.6	Application: Structured and Contextual Bandits	34
A.7	Application: Tabular Reinforcement Learning	39
A.8	Overview of Analysis	41
B	Lower Bounds for Learning the Optimal Allocation	44
B.1	Learning the Optimal Allocation: Minimax Formulation	45
B.2	Main Result	46
B.3	Examples	48
B.4	Discussion and Interpretation	50
C	Additional Related Work	53
II	Proofs	54
D	Additional Notation	55
E	Technical Tools	59
E.1	Online Learning	59
E.2	Properties of Graves-Lai Program	61
F	Proofs from Appendix A	71
F.1	Regret Bound for Uniformly Regular Classes (Theorem A.1)	72
F.2	Regret Bound without Uniform Regularity (Theorem A.2)	83
F.3	Estimation Guarantees	92
F.4	Supporting Lemmas	96
G	Proofs for Examples	103
G.1	Preliminaries: Regular Models	103
G.2	Structured Bandits with Gaussian Noise	106
G.3	Contextual Bandits with Finitely Many Actions (Example A.6)	115
G.4	Informative Arms (Example A.1)	116
G.5	Tabular Reinforcement Learning (Appendix A.7)	119
H	Proofs and Additional Results from Appendix B	132
H.1	Technical Lemmas	132

H.2	Proof of Theorem B.1	135
H.3	Proof of Theorem B.2	137
H.4	Proofs for Lower Bound Examples	142
H.5	Lower Bound on Regret for Algorithms with Well-Behaved Tails	147

Part I

Main Results

Part I is organized as follows. [Appendix A](#) presents our algorithm and main upper bounds, as well as examples. [Appendix B](#) presents complementary lower bounds. In [Appendix C](#) we review additional related work. We conclude with discussion of open problems and future directions in [Section 4](#). Proofs are deferred to the appendix.

Additional notation. For an integer $n \in \mathbb{N}$, we let $[n]$ denote the set $\{1, \dots, n\}$. For a set \mathcal{Z} , we let $\Delta_{\mathcal{Z}}$ denote the set of all probability distributions over \mathcal{Z} . We adopt standard big-oh notation, and write $f = \tilde{O}(g)$ to denote that $f = O(g \cdot \max\{1, \text{polylog}(g)\})$. We use \lesssim only in informal statements to emphasize the most notable elements of an inequality. We will let $\text{lin}(\cdot)$ denote a function multi-linear and poly-logarithmic in its arguments. For a decision $\pi \in \Pi$, we use $\mathbb{I}_{\pi} \in \Delta_{\Pi}$ to denote the delta distribution which places probability mass 1 on π .

Define the Kullback-Leibler divergence by

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

Appendix A. The AE^2 Algorithm: Regret Bounds and Examples

This section presents our main algorithm and regret bounds. We begin by introducing the most basic variant of our algorithm, AE^2 , and using it to provide instance-optimal regret bounds for simple settings ([Appendices A.2](#) and [A.3](#)); with preliminaries in [Appendix A.1](#). We then give a refined variant of the algorithm, AE_*^2 , which adapts to the minimum gap Δ_{\min}^* and leads to regret bounds under relaxed regularity conditions ([Appendix A.4](#) and [Appendix A.5](#)). We use this variant to provide applications to structured and contextual bandits ([Appendix A.6](#)) and tabular reinforcement learning ([Appendix A.7](#)). We conclude with an overview of our analysis in [Appendix A.8](#). For all results in this section, we assume that [Assumptions 1.1](#) to [1.3](#) hold.

A.1. Regularity Conditions

To present our algorithm and results, we first introduce several regularity conditions for the model class \mathcal{M} .

Likelihood ratios. We next make two assumptions concerning smoothness of KL divergences and behavior of log-likelihood ratios.

Assumption A.1 (Smooth KL). *There exists $L_{\text{KL}} > 0$ such that for all $M, M', M'' \in \mathcal{M}$ and $\pi \in \Pi$,*

$$\left| D_{\text{KL}}(M(\pi) \parallel M''(\pi)) - D_{\text{KL}}(M'(\pi) \parallel M''(\pi)) \right| \leq L_{\text{KL}} \sqrt{D_{\text{KL}}(M(\pi) \parallel M'(\pi))}.$$

Assumption A.2 (Sub-Gaussian Log-Likelihood). *There exists $V_{\mathcal{M}} > 0$ such that for all $M, M', M'' \in \mathcal{M}$ and $\pi \in \Pi$,*

$$\mathbb{P}_{(r,o) \sim M(\pi)} \left[\left| \log \frac{\mathbb{P}^{M',\pi}(r,o)}{\mathbb{P}^{M'',\pi}(r,o)} - \mathbb{E}_{(r',o') \sim M(\pi)} \left[\log \frac{\mathbb{P}^{M',\pi}(r',o')}{\mathbb{P}^{M'',\pi}(r',o')} \right] \right| \geq x \right] \leq 2 \exp(-x^2/V_{\mathcal{M}}^2)$$

for all $x \geq 0$.

[Assumption A.1](#) and [Assumption A.2](#) facilitate finite-sample estimation guarantees with respect to the KL divergence. Both assumptions are met by standard problem classes, including general structured bandit problems with Gaussian noise. Existing works that consider general model classes make similar assumptions ([Dong and Ma, 2022](#)).

Estimation. To provide estimation guarantees that accommodate infinite classes \mathcal{M} , we assume certain covering properties. We will consider the following notion of a cover.

Definition A.1 ((ρ, μ) -Cover). *We say that a set $\mathcal{M}_{\text{cov}} \subseteq \mathcal{M}$ is a (ρ, μ) -cover of \mathcal{M} if there exists some event \mathcal{E} such that:*⁴

1. $\sup_{M \in \mathcal{M}} \sup_{\pi \in \Pi} \mathbb{P}^{M, \pi}(\mathcal{E}^c) \leq \mu$.
2. For each $M \in \mathcal{M}$, there exists some $M' \in \mathcal{M}_{\text{cov}}$ such that

$$\left| \log \mathbb{P}^{M, \pi}(r, o) - \log \mathbb{P}^{M', \pi}(r, o) \right| \leq \rho$$

for all $(r, o) \in \mathcal{R} \times \mathcal{O}$ with $\sup_{M'' \in \mathcal{M}} \mathbb{P}^{M'', \pi}(r, o \mid \mathcal{E}) > 0$.

We denote the size of the smallest such cover by $N_{\text{cov}}(\mathcal{M}, \rho, \mu)$.

[Definition A.1](#) states that the log-likelihoods are “covered” under some good event \mathcal{E} which occurs with high probability: for any model in the class, we can find some model in the cover with log-likelihoods that are “close” on \mathcal{E} . We assume that the covering number for the model class \mathcal{M} is bounded, and has reasonable (“parametric”) growth.

Assumption A.3 (Bounded Covering Number). *For some parameters $d_{\text{cov}} \geq 1, C_{\text{cov}} \geq 1$, we have*

$$\log N_{\text{cov}}(\mathcal{M}, \rho, \mu) \leq d_{\text{cov}} \cdot \log \left(\frac{C_{\text{cov}}}{\rho \mu} \right).$$

Note that the rate of growth of the covering number required by [Assumption A.3](#) is the standard rate of growth for parametric (e.g., linear) classes. Our results easily extend to accommodate general growth rates, but we adopt [Assumption A.3](#) because it suffices for all of the examples we will consider, and simplifies presentation.

Information content of optimal decisions. As noted in the introduction, the Graves-Lai Coefficient $g^* = \text{glc}(\mathcal{M}, M^*)$ can be thought of as the minimal regret needed to distinguish M^* from all possible models with different optimal decisions. As playing the optimal decision, π_* , incurs no regret, any allocation η which is optimal for the Graves-Lai program, [Eq. \(3\)](#), will still be optimal if we increase the number of plays of π_* arbitrarily. As we are interested in finite-time behavior in this work, it is undesirable to consider allocations for which the number of pulls of optimal decisions are arbitrarily large. Instead, we would like to consider allocations which play optimal decisions only as long as they still provides useful information about models in the alternate set. The following definition gives a formal quantification of this.

4. Note that we require that $\mathcal{M}_{\text{cov}} \subseteq \mathcal{M}$, i.e. that \mathcal{M}_{cov} is a *proper* cover.

Definition A.2 (Information Content of Optimal Decision). *Fix $\varepsilon \in (0, 1/2]$. For a model $M \in \mathcal{M}$, we define $n_\varepsilon^M > 0$ as the minimum value such that, for any allocation $\eta \in \mathbb{R}_+^\Pi$ satisfying*

$$(1 + \varepsilon)g^M \geq \sum_{\pi \in \Pi} \eta(\pi) \Delta^M(\pi) \quad \text{and} \quad \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi \in \Pi} \eta(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) \geq 1 - \varepsilon,$$

we have

$$\inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi \in \Pi, \pi \notin \pi_M} \eta(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) + \sum_{\pi \in \pi_M} n_\varepsilon^M D_{\text{KL}}(M(\pi) \| M'(\pi)) \geq 1 - 2\varepsilon.$$

We denote $n_\varepsilon^{\mathcal{M}} := \sup_{M \in \mathcal{M}} n_\varepsilon^M$.

Intuitively, any allocation which is ε -optimal for the Graves-Lai program [Eq. \(3\)](#) need not play any optimal decision $\pi_M \in \pi_M$ more than n_ε^M times. Therefore, for model M , n_ε^M can be thought of as a quantification of the extent to which playing optimal decisions provides useful information—no additional useful information can be acquired on models in the alternate set $\mathcal{M}^{\text{alt}}(M)$ by playing optimal decisions more than n_ε^M times. As we will see, n_ε^M is bounded polynomially in problem parameters for many classes of interest.

Uniformly regular classes. We refer to a class as *uniformly regular* if $n_\varepsilon^{\mathcal{M}} < \infty$, and the following assumption on the minimum gaps holds.

Assumption A.4 (Lower-Bounded Minimum Gap). *We have $\inf_{M \in \mathcal{M}} \Delta_{\min}^M > 0$. We denote by $\Delta_{\min} > 0$ a (known) lower bound on $\inf_{M \in \mathcal{M}} \Delta_{\min}^M$.*

Note that [Assumption A.4](#) implies that for all $M \in \mathcal{M}$, π_M is unique. For the results concerning the most basic version of our algorithm, AE^2 ([Appendices A.2 and A.3](#)), we assume for expositional purposes that the class \mathcal{M} is uniformly regular. Our more general algorithm, AE_*^2 ([Appendix A.5](#)), achieves guarantees similar to those of AE^2 , but without uniform regularity. In particular, AE_*^2 replaces dependence on Δ_{\min} with the minimum gap $\Delta_{\min}^* := \Delta_{\min}^{M^*}$ for the true model, and replaces dependence on $n_\varepsilon^{\mathcal{M}}$ with $n_\varepsilon^* := n_\varepsilon^{M^*}$. We note, however, that in cases where a lower bound Δ on the minimum gap Δ_{\min}^* of the true model is known a-priori, [Assumption A.4](#) can be satisfied by restricting the model class to models with minimum gap at least Δ .

A.2. The AE^2 Algorithm

We now present the most basic variant of our main algorithm, AE^2 ([Algorithm 2](#)). This will serve as the starting point for the most general version of our algorithm, AE_*^2 ([Appendix A.5](#)). To describe the algorithm, we first introduce the primitive of an *online estimation oracle* ([Foster and Rakhlin, 2020; Foster et al., 2021](#)).

Estimation oracles. [Algorithm 2](#) makes use of an *online estimation oracle*, denoted by Alg_{KL} , which is an algorithm that, given knowledge of the class \mathcal{M} , estimates the underlying model $M^* \in \mathcal{M}$ from data in a sequential fashion. When invoked at step $s \in \mathbb{N}$ with the data $(\pi^1, r^1, o^1), \dots, (\pi^{s-1}, r^{s-1}, o^{s-1})$ observed so far, the estimation oracle builds an estimate

$$\widehat{M}^s = \text{Alg}_{\text{KL}} \left(\{(\pi^i, r^i, o^i)\}_{i=1}^{s-1} \right)$$

Algorithm 2 Allocation Estimation via Adaptive Exploration (AE²)

- 1: **input:** optimality tolerance ε , model class \mathcal{M} .
 - 2: Initialize $s \leftarrow 1$, $n_{\max} \leftarrow n_{\max}(\mathcal{M}, \varepsilon/6)$, and $q \leftarrow \frac{4n_{\max} + \varepsilon \mathbf{g}^{\mathcal{M}}}{4n_{\max} + 2\varepsilon \mathbf{g}^{\mathcal{M}}}$ for $\underline{\mathbf{g}}^{\mathcal{M}} := \inf_{M \in \mathcal{M}: \mathbf{g}^M > 0} \mathbf{g}^M$.
 - 3: Compute $\xi^1 \leftarrow \mathbf{Alg}_{\text{KL}}(\{\emptyset\})$ and $\widehat{M}^1 \leftarrow \mathbb{E}_{M \sim \xi^1}[M]$.
 - 4: **for** $t = 1, 2, 3, \dots$ **do**
 - 5: **if** $\exists \pi_{\widehat{M}^s} \in \pi_{\widehat{M}^s}$ s.t. $\forall M \in \mathcal{M}^{\text{alt}}(\pi_{\widehat{M}^s})$, $\sum_{i=1}^{s-1} \mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}^{M, \pi^i}(r^i, o^i)} \right] \geq \log(t \log t)$ **then**
 - 6: Play $\pi_{\widehat{M}^s}$. // Exploit
 - 7: **else** // Explore
 - 8: Set $p^s \leftarrow q\lambda^s + (1-q)\omega^s$ for

$$\lambda^s, \omega^s \leftarrow \arg \min_{\lambda, \omega \in \Delta_{\Pi}} \sup_{M \in \mathcal{M} \setminus \mathcal{M}_{\varepsilon/6}^{\text{gl}}(\lambda; n_{\max})} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\widehat{M}(\pi) \| M(\pi))]]}. \quad (18)$$
 - 9: Draw $\pi^s \sim p^s$ and observe reward r^s and observation o^s .
 - 10: Compute estimate $\xi^{s+1} \leftarrow \mathbf{Alg}_{\text{KL}}(\{(\pi^i, r^i, o^i)\}_{i=1}^s)$ and $\widehat{M}^{s+1} = \mathbb{E}_{\widehat{M} \sim \xi^{s+1}}[\widehat{M}]$.
 - 11: $s \leftarrow s + 1$.
-

which aims to approximate the true model M^* . Following (Foster et al., 2021; Chen et al., 2022; Foster et al., 2022a), we make use of *randomized* estimation oracles that, at each step produces $\xi^s = \mathbf{Alg}_{\text{KL}}(\{(\pi^i, r^i, o^i)\}_{i=1}^{s-1})$, where $\xi^s \in \Delta_{\mathcal{M}}$ is a randomization distribution, and draw $\widehat{M} \sim \xi^s$. We measure the oracle’s performance in terms of cumulative estimation error, defined as follows.

Definition A.3 (Cumulative Estimation Error). *Consider the process where, for each round $i \in \mathbb{N}$, given $(\pi^1, r^1, o^1), \dots, (\pi^{i-1}, r^{i-1}, o^{i-1})$ with $\pi^i \sim p^i$ and $(r^i, o^i) \sim M^*(\pi^i)$, the estimation oracle returns $\xi^i = \mathbf{Alg}_{\text{KL}}((\pi^1, r^1, o^1), \dots, (\pi^{i-1}, r^{i-1}, o^{i-1}))$. For any $s \in \mathbb{N}$, we define the oracle’s cumulative KL estimation error under this process as:*

$$\mathbf{Est}_{\text{KL}}(s) := \sum_{i=1}^s \mathbb{E}_{\widehat{M} \sim \xi^i} \left[\mathbb{E}_{\pi \sim p^i} \left[D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi)) \right] \right].$$

Algorithm 2 can be invoked with any off-the-shelf algorithm for estimation, but our main results make use of the fact that under Assumption A.2 and Assumption A.3, there exists an estimation oracle \mathbf{Alg}_{KL} (Algorithm 7 in Appendix F.3) which ensures that with probability at least $1 - \delta$, for all $s \in \mathbb{N}$:

$$\mathbf{Est}_{\text{KL}}(s) \lesssim V_{\mathcal{M}} \cdot d_{\text{cov}} \cdot \log^{3/2} \left(\frac{C_{\text{cov}} \cdot s}{\delta} \right). \quad (17)$$

That is, the estimation oracle ensures that the KL divergence between the true model M^* and the estimates returned scales at most poly-logarithmically in the exploration horizon. Note that on its own, this guarantee does not necessarily imply that $\widehat{M}^s = \mathbb{E}_{M \sim \xi^s}[M] \rightarrow M^*$ —low online estimation error only requires that \widehat{M}^s is on average close to M^* on the decisions we have actually played.

Algorithm overview. We now present a formal overview of our algorithm AE^2 . We restate it here for convenience in [Algorithm 2](#). The algorithm alternates between *exploit* steps and *explore* steps, tracking the number of explore steps that have been performed with a counter $s \in \mathbb{N}$. For each step $t \in \mathbb{N}$, the algorithm makes use of an estimator $\widehat{M}^s = \mathbb{E}_{\widehat{M} \sim \xi^s}[\widehat{M}]$, where $\xi^s = \text{Alg}_{\text{KL}}(\{(\pi^i, r^i, o^i)\}_{i=1}^{s-1})$ is computed by calling the estimation oracle with data gathered at previous explore steps. Given the estimator, AE^2 performs a test based on likelihood ratios ([Line 5](#)) to check whether it has collected enough information to rule out all models for which $\pi_{\widehat{M}^s} \in \pi_{\widehat{M}^s}$ is not an optimal decision. If so, it *exploits*, and plays $\pi_{\widehat{M}^s}$ ([Line 6](#)), as in this case $\pi_{\widehat{M}^s} = \pi_*$ with high probability. If the test fails, the algorithm must gather more information to eliminate alternatives, and it *explores* ([Line 8](#)). The key component of the explore phase is the choice of the exploration distribution in [Eq. \(18\)](#), which is based on the Allocation-Estimation Coefficient program, but incorporates some small modifications: 1) First, \widehat{M} is randomized according to the distribution ξ^s , 2) Second, the set $M^* \in \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)$ is replaced with a smaller set $M^* \in \mathcal{M}_\varepsilon^{\text{gl}}(\lambda; n_{\max})$, which requires that λ obeys certain normalization constraints; this is detailed below. Using the distributions λ^s (representing a normalized allocation) and ω^s (representing an exploration distribution) returned in [Eq. \(18\)](#), the algorithm computes a mixture $p^s = q\lambda^s + (1-q)\omega^s$, where $q \in (0, 1)$ is a carefully chosen parameter, and plays $\pi^s \sim p^s$ ([Line 9](#)). The reward and observation (r^s, o^s) that result from playing π^s are then used to update the estimation oracle for subsequent rounds ([Line 10](#)).

To understand the intuition behind the explore phase and why the Allocation-Estimation Coefficient plays a useful role here, we can consider two cases. In the first case, if λ^s is an ε -optimal Graves-Lai allocation for M^* (that is, $M^* \in \mathcal{M}_{\varepsilon/6}^{\text{gl}}(\lambda^s; n_{\max})$), then playing λ^s will optimize the tradeoff between minimizing regret on M^* and collecting information that allows one to distinguish M^* from $M \in \mathcal{M}^{\text{alt}}(M^*)$, and will therefore match the optimal performance prescribed by the Graves-Lai Coefficient, incurring regret scaling as g^* .

In the second case, if λ^s is not an ε -optimal Graves-Lai allocation for M^* , we have $M^* \notin \mathcal{M}_{\varepsilon/6}^{\text{gl}}(\lambda^s; n_{\max})$, so by the definition of the AEC ([Eq. \(18\)](#)), ω^s will place mass on actions that ensure $\mathbb{E}_{\widehat{M} \sim \xi^s}[\mathbb{E}_{\pi \sim \omega}[D_{\text{KL}}(\widehat{M}(\pi) \| M^*(\pi))]]$ is large; exactly how large this quantity is will be quantified by the value of the AEC. Since p^s plays ω^s with constant probability, the quantity

$$\mathbb{E}_{M \sim \xi^s}[\mathbb{E}_{\pi \sim p^s}[D_{\text{KL}}(M^*(\pi) \| M(\pi))]]$$

will also be large, but if the estimation oracle is consistent in the sense of [Definition A.3](#), this can only happen a small number of times. In particular, if [Eq. \(17\)](#) holds, the number of times in which we encounter this second case is at most *logarithmic* in the number of exploration rounds. As such, we can show that λ^s must be a near-optimal Graves-Lai allocation for M^* on all but a logarithmic number of exploration rounds, and that AE^2 achieves the optimal rate on such rounds.

Critically, rather than exploring in a naive fashion (e.g., by sampling decisions uniformly), AE^2 explores only to the extent necessary to learn a Graves-Lai allocation for M^* . There may exist instances $M \neq M^*$ which differ significantly from M^* but have a similar Graves-Lai allocations— AE^2 will make no effort to distinguish such instances since, as long as it knows that one of these instances is correct, it can simply play their shared Graves-Lai allocation. This notion of exploration, which is targeted toward distinguishing instances that have different Graves-Lai allocations, is precisely the notion captured by the Allocation-Estimation Coefficient.

Normalization factor for allocations. As noted in [Appendix A.1](#), while an optimal Graves-Lai allocation may place an arbitrarily large number of pulls on an optimal decision, for finite-time

guarantees it is useful to restrict to allocations which place only finite mass on optimal decisions. To this end, AE^2 restricts the optimization problem based on the AEC in Eq. (18) to only consider normalized allocations λ for which the *normalization factor* is at most

$$n_{\max}(\mathcal{M}, \varepsilon) := \frac{64}{\Delta_{\min}^2} \cdot \left(\frac{1}{\varepsilon} + V_{\mathcal{M}} n_{\varepsilon}^{\mathcal{M}} \right) \cdot \max_{M \in \mathcal{M}} g^M. \quad (19)$$

where the normalization factor refers to the value n in the definition of $\Lambda(M; \varepsilon)$ (see Eq. (8)). In particular, to enforce this restriction, the optimization problem in Eq. (18) restricts the max-player to $M \in \mathcal{M} \setminus \mathcal{M}_{\varepsilon/6}^{\text{gl}}(\lambda; n_{\max})$, where $n_{\max} := n_{\max}(\mathcal{M}, \varepsilon/6)$ and $\mathcal{M}_{\varepsilon/6}^{\text{gl}}(\lambda; n_{\max})$ is defined identically to $\mathcal{M}_{\varepsilon/6}^{\text{gl}}(\lambda)$ in Eq. (9), but with n restricted to $n \leq n_{\max}$. As we show in Lemma E.4, for $n_{\max}(\mathcal{M}, \varepsilon)$ defined as in Eq. (19), the optimal value in Eq. (18) can be bounded by the AEC.

Computational efficiency. The primary computational burden in AE^2 lies in solving the optimization problem (18) to compute the exploration distributions. In general there is little hope of solving this efficiently (i.e., in time sublinear in $|\Pi|$ and $|\mathcal{M}|$)—indeed, in some cases it may be that to even determine whether $M \in \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda)$ will require enumerating the model class \mathcal{M} . However, for nicely structured problems, we anticipate that this program can be solved, or at least approximated, efficiently. As the focus of our work is primarily statistical, we leave further exploration as to when the algorithm can be implemented efficiently to future work.

Simplicity. We emphasize the simplicity of AE^2 . Most existing algorithms which achieve instance-optimality are quite complex, even in specialized settings such as linear bandits. In contrast, AE^2 is very simple and intuitive, and relies only on three basic components: an explore-exploit test, an estimation oracle, and a single optimization to compute the exploration distributions. Despite its simplicity, as we show, AE^2 obtains comparable or better performance over existing approaches.

Relation to existing approaches. At a very high level, AE^2 bears some similarity to the E2D algorithm of Foster et al. (2021), which achieves the *minimax* optimal rate for general classes \mathcal{M} in the DMSO framework. Both algorithms rely on online estimation algorithms, and both solve min-max programs based on the output of the estimator to determine which allocations to play. However, the algorithm design and analysis principles for the two algorithms, and in particular the motivation for the min-max programs they solve, differ significantly.

A.3. AE^2 Algorithm: Regret Bound for Uniformly Regular Classes

We present upper bounds for AE^2 in the setting where our class \mathcal{M} is uniformly regular: Assumption A.4 holds and $n_{\varepsilon}^{\mathcal{M}} < \infty$; these assumptions are relaxed by the AE_*^2 algorithm in the sequel. To state the regret bound for AE^2 in the tightest form possible, we introduce the following variant of the Allocation-Estimation Coefficient, which incorporates randomized estimators $\xi \in \Delta_{\mathcal{M}}$:

$$\overline{\text{aec}}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}_0, \xi) := \inf_{\lambda, \omega \in \Delta_{\Pi}} \sup_{M \in \mathcal{M} \setminus \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda)} \frac{1}{\mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \| M(\pi))]]}, \quad (20)$$

with $\overline{\text{aec}}_{\varepsilon}(\mathcal{M}, \xi) := \overline{\text{aec}}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}, \xi)$ and $\overline{\text{aec}}_{\varepsilon}(\mathcal{M}) := \sup_{\xi \in \Delta_{\mathcal{M}}} \overline{\text{aec}}_{\varepsilon}(\mathcal{M}, \xi)$. Note that one can always bound $\overline{\text{aec}}_{\varepsilon}(\mathcal{M}) \leq \text{aec}_{\varepsilon}(\mathcal{M})$ due to the convexity of the KL divergence. In fact, these definitions are equivalent up to dependence on problem-dependent parameters in Appendix A.1 (indeed, our

lower bounds in [Appendix B](#) scale with the latter quantity), but the former can be simpler to bound for some of the examples we consider.

Our main theorem concerning the performance of AE^2 is as follows.

Theorem A.1 (Regret Bound for AE^2). *For any $\varepsilon \in (0, 1/2]$, there exists a choice for the estimation oracle Alg_{KL} such that for all $T \in \mathbb{N}$, under [Assumptions A.1](#) to [A.4](#) and if $\mathbf{g}^* > 0$, the expected regret of AE^2 is bounded by*

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq (1 + \varepsilon)\mathbf{g}^* \cdot \log(T) + \overline{\text{aec}}_{\varepsilon/12}(\mathcal{M}) \cdot C_{\text{aec}} \cdot \log^{3/2}(\log T) + C_{\text{low}} \cdot \log^{1/2}(T), \quad (21)$$

where

$$C_{\text{aec}} := c \cdot \frac{V_{\mathcal{M}}^2 d_{\text{cov}} \log(C_{\text{cov}}) \cdot \max_{M \in \mathcal{M}} \mathbf{g}^M}{\varepsilon \Delta_{\min}^3} \cdot \left(\varepsilon^{-1} + V_{\mathcal{M}} n_{\varepsilon/6}^M \right) \cdot \log(C_{\text{low}}),$$

for a universal numerical constant $c > 0$, and C_{low} is a lower-order constant given by

$$C_{\text{low}} := \text{lin} \left(\max_{M \in \mathcal{M}} \mathbf{g}^M, \text{aec}_{\varepsilon/12}^{1/2}(\mathcal{M}), \frac{1}{\varepsilon^2}, \frac{1}{\Delta_{\min}^3}, n_{\varepsilon/6}^M, L_{\text{KL}}^2, V_{\mathcal{M}}^{13/2}, d_{\text{cov}}, \log(C_{\text{cov}}), \log \log T \right),$$

where $\text{lin}(\cdot)$ denotes a function multi-linear and poly-logarithmic in its arguments.

We prove [Theorem A.1](#) in [Appendix F.1](#), and give a proof sketch in [Appendix A.8](#). [Theorem A.1](#) shows that AE^2 achieves the asymptotically optimal Graves-Lai rate for M^* , as given in [Proposition 1.1](#), up to a $(1 + \varepsilon)$ approximation factor. In more detail, if we label the terms in [Eq. \(21\)](#) as

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq \underbrace{(1 + \varepsilon)\mathbf{g}^* \cdot \log(T)}_{\text{(I)}} + \underbrace{\overline{\text{aec}}_{\varepsilon/12}(\mathcal{M}) \cdot C_{\text{aec}} \cdot \log^{3/2}(\log T)}_{\text{(II)}} + \underbrace{C_{\text{low}} \cdot \log^{1/2}(T)}_{\text{(III)}},$$

the regret bound can be seen to consist of:

- The *leading-order term* (I) = $(1 + \varepsilon)\mathbf{g}^* \cdot \log(T)$. This is the only term that scales linearly with $\log(T)$, as a consequence we have $\lim_{T \rightarrow \infty} \frac{\mathbb{E}^{M^*}[\mathbf{Reg}(T)]}{\log(T)} \leq (1 + \varepsilon)\mathbf{g}^*$, which matches the instance-optimal rate given in [Proposition 1.1](#) up to a factor of $(1 + \varepsilon)$.
- A lower-order term (II) = $\overline{\text{aec}}_{\varepsilon/12}(\mathcal{M}) \cdot C_{\text{aec}} \cdot \log^{3/2}(\log T)$, which is polylogarithmic in $\log(T)$, and scales with $\overline{\text{aec}}_{\varepsilon/12}(\mathcal{M})$, as well as regularity parameters from [Appendix A.1](#).
- A second lower-order term (III) = $C_{\text{low}} \cdot \log^{1/2}(T)$. This term scales with $\log^{1/2}(T) = o(\log(T))$ and, like the term (II), scales with the AEC and regularity parameters from [Appendix A.1](#). Compared to (II), this term has worse dependence on $\log(T)$, but enjoys sublinear $\overline{\text{aec}}_{\varepsilon/12}^{1/2}(\mathcal{M})$ scaling with the AEC.

Critically, both of the $o(\log T)$ lower-order terms above do not scale with (often exponentially large) terms such as $|\text{II}|$ or $|\mathcal{M}|$ found in prior work, and instead scale principally with $\overline{\text{aec}}_{\varepsilon}(\mathcal{M})$, which, as

we will show in [Appendix B](#), is unavoidable in a certain sense. In particular, note that once T is large enough that

$$\log(T) \geq \tilde{\Omega}^+(\text{aec}_{\varepsilon/12}(\mathcal{M})),$$

the leading-order term $(I) = (1 + \varepsilon)g^* \cdot \log(T)$ term in [Theorem A.1](#) will dominate the regret. This is precisely the time horizon given by the lower bound in [Theorem 2.3](#), which is necessary for an algorithm to learn a near-optimal allocation for the Graves-Lai program. We offer a more thorough comparison of [Theorem A.1](#) with our lower bounds in [Appendix B.4](#). Below, we discuss the lower-order terms and asymptotic performance in greater detail.

Remark A.1 (Additional Lower-Order Terms). *The lower-order terms in [Theorem A.1](#) depend on the model class \mathcal{M} through the regularity, covering, and smoothness assumptions, as well as the minimum gap ([Appendix A.1](#)). For many of the examples we consider, the Allocation-Estimation Coefficient will dominate these other terms, yet there may exist classes where this is not the case. Resolving the optimal dependence on these problem-dependent parameters in the lower-order terms, as well as understanding when these parameters are necessary, remains an interesting direction for future work.*

In addition, let us mention that while both lower-order terms scale with $o(\log T)$ (note that the scaling is no larger than $O(\sqrt{\log T} \cdot \text{polylog } \log T)$), it is not clear what the optimal dependence on T should be for the lower-order terms. For example, one might hope to replace the dependence on $\log^{1/2}(T)$ with $\log^a(T)$ for some constant $a < 1/2$, or even with $\text{polylog}(\log(T))$. Precisely characterizing the optimal $\log(T)$ scaling for lower-order terms remains an interesting open question. To this end, we remark that [Jun and Zhang \(2020\)](#) show that in some cases, an $\Omega(\log \log T)$ term is indeed necessary.

Remark A.2 (Asymptotic Performance). *Asymptotically, as $T \rightarrow \infty$, the regret of AE^2 scales with $(1 + \varepsilon)g^* \cdot \log(T)$, which is a factor of $(1 + \varepsilon)$ off from the asymptotic lower bound in [Proposition 1.1](#). For any fixed $T \in \mathbb{N}$ of interest, as long as $\text{aec}_{\varepsilon}(\mathcal{M}) = \text{poly}(\varepsilon^{-1})$, one can obtain an asymptotic constant of 1 by choosing $\varepsilon = 1/\log^a(T)$ for a sufficiently small constant $a > 0$. For example, when $|\Pi| < \infty$, it is always possible to bound $\text{aec}_{\varepsilon}(\mathcal{M}) \lesssim \text{poly}(|\Pi|)/\varepsilon^4$ (see [Proposition A.1](#) below), so choosing ε as above ensures that the lower-order terms scale $o(\log T)$, while the leading-order term scales as $g^* \cdot \log(T)$ asymptotically.*

A.3.1. EXAMPLE: SEARCHING FOR AN INFORMATIVE ARM

We next provide an example of a uniformly regular class in order to illustrate a case where [Theorem A.1](#) holds. In particular, we revisit the *informative arm* setting described in the introduction ([Example 1.1](#)). Recall that we exhibited a model class for which the complexity of learning the Graves-Lai allocation is not governed by existing complexity measures, and can be larger than the minimax optimal rate for learning with \mathcal{M} . In what follows, we show that on this example the Allocation-Estimation Coefficient correctly adapts to the complexity of this model class. We emphasize that the main applications of our results, which take advantage of the more general AE_{\star}^2 algorithm, will be given in [Appendix A.6](#) and [Appendix A.7](#), and we also present additional examples of uniformly regular classes in [Appendix G.2.1](#).

Example A.1 (Searching for an Informative Arm (revisited)). Let \mathcal{M} denote the model class constructed in [Example 1.1](#), with parameters $A, N \geq 5$ and $\beta \in [4/A, 9/10]$. We additionally discretize

the space so that, for each $M \in \mathcal{M}$ and $\pi \in [A]$, we have $f^M(\pi) \in \{0, \Delta_{\min}, 2\Delta_{\min}, \dots, \lfloor \frac{1}{\Delta_{\min}} \rfloor \Delta_{\min}\}$,⁵ and furthermore restrict \mathcal{M} so that it does not include instances with multiple optimal arms. For this class, one can show that [Assumptions A.1 to A.4](#) hold with $L_{\text{KL}}, V_{\mathcal{M}} \leq O(\log A)$ and $d_{\text{cov}} = O(A), C_{\text{cov}} = O(N)$ (see [Appendix G.4](#)). Furthermore, we can bound $n_{\varepsilon}^M \leq \frac{2}{\Delta_{\min}^2}$, and

$$\text{aec}_{\varepsilon}(\mathcal{M}) \leq \frac{64N}{\beta^2} + \frac{16A}{\Delta_{\min}^2}.$$

As a result, for this class, AE^2 has expected regret bounded as

$$\mathbb{E}^{M^*} [\mathbf{Reg}(T)] \leq (1 + \varepsilon) \mathbf{g}^* \cdot \log(T) + N \cdot \text{poly}\left(A, \frac{1}{\varepsilon}, \frac{1}{\Delta_{\min}}, \log N, \log \log T\right) \cdot \log^{1/2}(T).$$

Note that here the only term that scales linearly with N is the Allocation-Estimation Coefficient—every other class-dependent term appearing in the regret bound scales at most logarithmically in N . We are particularly interested in situations where the cost of finding the correct informative arm is much larger than any existing complexity measures for the problem: that is, when β is constant and $N \gg A, \frac{1}{\Delta_{\min}}$. In this case, we have $\text{aec}_{\varepsilon}(\mathcal{M}) \leq O(N)$, and the Allocation-Estimation Coefficient correctly captures the intuitive complexity of learning the optimal allocation. In particular, the dependence on N reflects the fact that we need to test each informative arm at least once. Furthermore, as we show in [Example B.1](#), we can lower bound $\text{aec}_{\varepsilon}(\mathcal{M}) \geq \Omega(N)$ as well, so in the regime where $N \gg A, \frac{1}{\Delta_{\min}}$, the AEC is the dominant lower-order term. \triangleleft

See [Appendix G.2](#) for the proof of this example.

A.4. The AE_{*}^2 Algorithm

While it may be reasonable to assume that the minimum gap of M^*, Δ_{\min}^* , is bounded away from 0, and that the amount of useful information playing π_* provides is also bounded on M^* , assuming that this is true for every model in the model class (as in the prequel) is a significantly stronger assumption. For example, if we let \mathcal{M} denote the space of all multi-armed bandits with means in $[0, 1]$, the only possible value of Δ_{\min} is 0, as we can always find some instance with minimum gap arbitrarily close to 0. In this section, we dispense with the uniform regularity assumption: we relax [Assumption A.4](#), and additionally prove that it suffices if only $n_{\varepsilon}^* := n_{\varepsilon}^{M^*}$ (as opposed to n_{ε}^M) is bounded.

Our main algorithm for this section, AE_{*}^2 , is given in [Algorithm 3](#). It is very similar to AE^2 but to remove the requirement of uniform regularity, the algorithm avoids solving [Eq. \(18\)](#) over the entire model class \mathcal{M} , and instead solves it over a carefully restricted model class. For $x, y > 0$, define

$$\mathcal{M}_{x,y} := \{M \in \mathcal{M} : \Delta_{\min}^M \geq x, n_{\varepsilon}^M \leq y\}. \quad (23)$$

AE_{*}^2 breaks its explore rounds into doubling epochs. For each epoch ℓ , [Eq. \(22\)](#) in [Algorithm 3](#) solves an AEC-like optimization problem over a restricted class $\mathcal{M}^{\ell} \subseteq \mathcal{M}_{\Delta^{\ell}, \frac{1}{\Delta^{\ell}}}$, which is chosen in [Line 9](#) to explicitly ensure that the value of the optimization problem in [Eq. \(22\)](#), is bounded; this is guaranteed by the definition of Δ^{ℓ} in [Line 8](#). Similar to AE^2 , the value of the optimization

5. The discretization is required to satisfy the uniform regularity assumption—we include it here to provide a concrete example of a uniformly regular class. However, it can be shown that, without this discretization assumption, [Theorem A.2](#) applies to the original formulation given in [Example 1.1](#) with the AEC again scaling as $O(N)$.

Algorithm 3 Adaptive Exploration for Allocation Estimation for classes without uniform regularity (AE_*^2)

1: **input:** Optimality tolerance ε , estimation oracle \mathbf{Alg}_{KL} , growth parameters $\alpha_q, \alpha_n, \alpha_{\mathcal{M}} \geq 0$.
 2: $s \leftarrow 1, \ell \leftarrow 1, \delta \leftarrow \frac{\varepsilon}{4+2\varepsilon}, q^s \leftarrow 1 - s^{-\alpha_q}, n^s \leftarrow s^{\alpha_n}$.
 3: Compute $\xi^1 \leftarrow \mathbf{Alg}_{\text{KL}}(\{\emptyset\})$ and $\widehat{M}^1 \leftarrow \mathbb{E}_{M \sim \xi^1}[M]$.
 4: **for** $t = 1, 2, 3, \dots$ **do**
 5: **if** $s \geq 2^\ell$ **then** // Form active set and cover
 6: $\ell \leftarrow \ell + 1$.
 7: $\Delta^\ell \leftarrow \arg \min_{\Delta \geq 0} \Delta \quad \text{s.t.} \quad \text{aec}_{\delta/2}^{\mathcal{M}}(\mathcal{M}_{\Delta, \frac{1}{\Delta}}) \leq s^{\alpha_{\mathcal{M}}}$.
 8: $\mathcal{M}^\ell \leftarrow \mathcal{M}_{\Delta^\ell, \frac{1}{\Delta^\ell}} \cap \{M \in \mathcal{M} : n_\varepsilon^M + \frac{1}{\Delta_{\min}^M} + \frac{4g^M}{\Delta_{\min}^M} + \frac{2n_\varepsilon^M}{g^M} + \frac{4}{\Delta_{\min}^M \delta} \leq \sqrt{n^s}\}$.
 9: $\mathcal{M}_{\text{cov}}^\ell \leftarrow (\rho_\ell, \mu_\ell)$ -cover of \mathcal{M}^ℓ for $\rho_\ell \leftarrow 2^{-\ell}, \mu_\ell \leftarrow 2^{-5\ell}, \mathfrak{D}^\ell \leftarrow \emptyset$.
 10: **if** $\exists \pi_{\widehat{M}^s} \in \pi_{\widehat{M}^s}$ s.t. $\forall M \in \mathcal{M}^{\text{alt}}(\pi_{\widehat{M}^s}), \sum_{i=1}^{s-1} \mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}_{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}_{M, \pi^i}(r^i, o^i)} \right] \geq \log(t \log t)$ **then**
 11: Play $\pi_{\widehat{M}^s}$. // Exploit
 12: **else** // Explore
 13: Set $p^s \leftarrow q^s \lambda^s + (1 - q^s) \omega^s$ for

$$\lambda^s, \omega^s \leftarrow \arg \min_{\lambda, \omega \in \Delta_\Pi} \sup_{M \in \mathcal{M}^\ell \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda; n^s)} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\widehat{M}(\pi) \| M(\pi))]]}. \quad (22)$$

14: Draw $\pi^s \sim p^s$, observe r^s, o^s , set $\mathfrak{D}^\ell \leftarrow \mathfrak{D}^\ell \cup \{(\pi^s, r^s, o^s)\}$.
 15: Compute estimate $\xi^{s+1} \leftarrow \mathbf{Alg}_{\text{KL}}(\mathfrak{D}^\ell, \mathcal{M}_{\text{cov}}^\ell)$ and $\widehat{M}^{s+1} = \mathbb{E}_{M \sim \xi^{s+1}}[M]$.
 16: $s \leftarrow s + 1$.

problem in Eq. (22) quantifies how much information we are gaining about the Graves-Lai allocation of M^* , and the regret of the explore phase can be bounded in terms of the value of this optimization. By restricting \mathcal{M}^ℓ so that the value of Eq. (22) is always bounded, we can therefore ensure that the regret during the exploration phase is bounded.

Intuitively, this restriction of \mathcal{M}^ℓ reduces the space of models we must distinguish M^* from in order to identify its Graves-Lai allocation: rather than distinguishing M^* from all models in \mathcal{M} , we must only distinguish it from models in \mathcal{M}^ℓ , which could be significantly easier. The caveat is that, since we do not know the value of n_ε^* or Δ_{\min}^* , M^* may not always be in \mathcal{M}^ℓ . In such cases, little can be said about the exploration phase—we are not able to provide any meaningful guarantees on how much information λ^s and ω^s acquire about M^* . To mitigate this, as s increases we gradually relax the criteria for inclusion in \mathcal{M}^ℓ , ensuring that for large enough s , M^* will be in \mathcal{M}^ℓ . In particular, one can show that the number of exploration rounds needed to guarantee $M^* \in \mathcal{M}^\ell$ scales with $\text{aec}_\varepsilon^{\mathcal{M}}(M^*)$, for

$$\mathcal{M}^* := \{M \in \mathcal{M} : \Delta_{\min}^M \geq \Delta_*, n_{\varepsilon/36}^M \leq 1/\Delta_*\} \quad \text{for} \quad \Delta_* := \min\{\Delta_{\min}^*, 1/n_{\varepsilon/36}^*\}. \quad (24)$$

That is, \mathcal{M}^* is the restriction of \mathcal{M} to models with gap at least $\min\{\Delta_{\min}^*, 1/n_{\varepsilon/36}^*\}$ (implying all models in \mathcal{M}^* have a unique optimal decision), and for which the information content of the optimal decision is at most $\max\{1/\Delta_{\min}^*, n_{\varepsilon/36}^*\}$.

Estimation oracle. While AE^2 simply requires that the estimation oracle Alg_{KL} returns randomized estimators supported on $\Delta_{\mathcal{M}}$, for AE_*^2 , we wish to ensure that the estimators produced are instead only supported on \mathcal{M}^ℓ . To this end, we restrict the estimator to $\mathcal{M}_{\text{cov}}^\ell$, the (ρ_ℓ, μ_ℓ) -cover of \mathcal{M}^ℓ . We denote the resulting estimation oracle with $\text{Alg}_{\text{KL}}(\mathcal{D}^\ell, \mathcal{M}_{\text{cov}}^\ell)$, where the first argument represents the set of available observations, and the second argument the set over which the estimation oracle must return an estimate.

Computational efficiency of AE_*^2 . Similar to AE^2 , it is not clear how to solve the main optimization required by AE_*^2 , Eq. (22), in general. In addition, unlike AE^2 , AE_*^2 maintains a version space of models, \mathcal{M}^ℓ , which could increase the computational burden further. We emphasize that the focus of this work is primarily statistical, and leave addressing the computational challenges for future work.

A.5. AE_*^2 Algorithm: Regret Bound without Uniform Regularity

The following theorem provides the main guarantee for AE_*^2 .

Theorem A.2 (Regret Bound for AE_*^2). *For any $\varepsilon \in (0, 1/2]$, if Assumptions A.1 to A.3 hold and $g^* > 0$, AE_*^2 (Algorithm 3) ensures that for all $T \in \mathbb{N}$, the expected regret is bounded as*

$$\mathbb{E}^{M^*}[\text{Reg}(T)] \leq (1 + \varepsilon)g^* \cdot \log(T) + \left(\overline{\text{aec}}_{\varepsilon/12}^{\mathcal{M}}(\mathcal{M}^*)\right)^3 \cdot C_{\text{aec}} \cdot \log^{3/2}(\log T) + C_{\text{low}} \cdot \log^{6/7}(T),$$

where

$$C_{\text{aec}} := \tilde{O}\left(\frac{V_{\mathcal{M}}^3(V_{\mathcal{M}} + L_{\text{KL}}) \cdot d_{\text{cov}} \log(C_{\text{cov}})}{\varepsilon \Delta_{\min}^*}\right),$$

and C_{low} is a lower-order constant given by

$$C_{\text{low}} := \text{poly}\left(g^*, \frac{1}{\Delta_{\min}^*}, n_{\varepsilon/6}^*, \frac{1}{\varepsilon}, V_{\mathcal{M}}, L_{\text{KL}}, d_{\text{cov}}, \log C_{\text{cov}}, \log \log T\right).$$

The proof of Theorem A.2 is given in Appendix F.2. As Theorem A.2 illustrates, at the expense of a slightly larger polynomial dependence on the Allocation-Estimation Coefficient, and slightly larger lower-order terms, we can obtain near instance-optimal regret—matching the instance-optimal lower bound given in Proposition 1.1 up to a factor of $(1 + \varepsilon)$ —without requiring any assumption on the minimum gap, or boundedness of $n_\varepsilon^{\mathcal{M}}$. Rather than scaling with the minimum gap for the entire class, Δ_{\min} , Theorem A.2 scales only with the minimum gap of the ground truth model, Δ_{\min}^* , which could be substantially larger than Δ_{\min} . An additional advantage of Theorem A.2 is that it scales with $\text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}^*)$ as opposed to $\text{aec}_\varepsilon(\mathcal{M})$; for the examples we consider in the sequel, the former quantity enjoys better dependence on problem-dependent parameters. For example, we show in Appendix B that for standard classes, $\text{aec}_\varepsilon(\mathcal{M})$ can scale with the minimum gap amongst all models in the class. On the other hand $\text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}^*)$ typically scales with Δ_{\min}^* .

We emphasize that AE_*^2 requires no prior knowledge of Δ_{\min}^* or n_ε^* —it is able to adapt to the minimum gap and regularity of the underlying model.

Remark A.3 (Dependence on n_ε^*). *As we show in the following examples, n_ε^* is typically bounded polynomially in standard problem parameters, though in practice this needs to be verified for each problem instance. We remark that some scaling in terms of n_ε^* seems unavoidable—if there is a significant amount of information to be gained playing the optimal decision, any algorithm which is nearly instance-optimal will play the optimal decision at least n_ε^* times, and therefore the “effective horizon” to eliminate alternate instances scales with n_ε^* . As we are the first to formalize this notion of how informative the optimal decision is, we believe more research in this direction is required.*

A.6. Application: Structured and Contextual Bandits

We now instantiate [Theorem A.2](#) to give regret bounds for AE_*^2 in standard settings of interest, bounding the Allocation-Estimation Coefficient for each setting. We begin by focusing on structured bandit settings and contextual bandits, then turn to tabular reinforcement learning in the sequel. We recall that, to map bandit problems to the DMSO framework, we take the decision space Π to be the set of “arms”, the observation space $\mathcal{O} = \{\emptyset\}$, and the reward space \mathcal{R} to be the rewards from the bandit (while we do not explicitly include the rewards in the observation space, we assume they are observed). We defer proofs for all examples to [Appendix G](#).

A.6.1. THE UNIFORM EXPLORATION COEFFICIENT

For the main examples we consider, we proceed by first bounding the Allocation-Estimation Coefficient in terms of another, somewhat simpler parameter we refer to as the *uniform exploration coefficient*.

Definition A.4 (Uniform Exploration Coefficient). *For a randomized estimator $\xi \in \Delta_{\mathcal{M}}$, we define the uniform exploration coefficient with respect to ξ at scale $\varepsilon > 0$ as the value of the following program:*

$$C_{\text{exp}}^{\xi}(\varepsilon) := \min_{C \in \mathbb{R}_+, p \in \Delta_{\Pi}} \left\{ C \mid \forall M, M' \in \mathcal{M} : \begin{array}{l} \max_{M'' \in \{M, M'\}} \mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{\pi \sim p} [D_{\text{KL}}(\bar{M}(\pi) \parallel M''(\pi))]] \leq 1/C \\ \implies \max_{p' \in \Delta_{\Pi}} \mathbb{E}_{\pi \sim p'} [D_{\text{KL}}(M(\pi) \parallel M'(\pi))] \leq \varepsilon \end{array} \right\}.$$

We define $p_{\text{exp}}^{\xi}(\varepsilon)$ as any minimizing distribution for this program, and let

$$C_{\text{exp}}(\mathcal{M}, \varepsilon) := \sup_{\xi \in \Delta_{\mathcal{M}}} C_{\text{exp}}^{\xi}(\varepsilon)$$

denote the uniform exploration constant for class \mathcal{M} .

Intuitively, the uniform exploration coefficient characterizes the extent to which it is possible to explore by uniformly covering the decision space. In particular, one can always choose p to be uniform over Π , which gives $C_{\text{exp}}(\mathcal{M}, \varepsilon) \lesssim |\Pi|/\varepsilon$, but in cases where information is shared between actions, the parameter is significantly smaller, as we will show for familiar examples below. For example, in the case of linear bandits with dimension d , we have $C_{\text{exp}}(\mathcal{M}, \varepsilon) \leq \tilde{O}(\frac{d \log 1/\varepsilon}{\varepsilon})$.

The following result shows that the Allocation-Estimation Coefficient can be bounded in terms of the uniform exploration coefficient.

Proposition A.1 (Informal). *For $\mathcal{M}_0 \subseteq \mathcal{M}$, we can bound $\overline{\text{aec}}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}_0) \leq C_{\text{exp}}(\mathcal{M}_0, \delta)$ for any*

$$\sqrt{\delta} \leq \min_{M \in \mathcal{M}_0} \min \left\{ \min \left\{ \frac{1}{81L_{\text{KL}}}, \frac{\Delta_{\min}^M}{34V_{\mathcal{M}}} \right\} \cdot \frac{\varepsilon}{2g^M/\Delta_{\min}^M + n_{\varepsilon/36}^M}, \frac{\Delta_{\min}^M}{3} \right\}.$$

The full statement of [Proposition A.1](#) is given in [Lemma E.6](#). Using [Proposition A.1](#), we obtain guarantees for AE_*^2 on several familiar classes, beginning with several bandit settings. We remark that [Proposition A.1](#) is not in general tight—it simply shows that the Allocation-Estimation Coefficient is bounded by a simple, general, and interpretable notion of how easily a class can be explored. As such the bounds on the Allocation-Estimation Coefficient in the following examples can almost certainly be improved using more specialized tools.

A.6.2. FINITE-ARMED BANDITS

We first consider the simplest bandit setting: multi-armed bandits with finite arms.

Example A.2 (Finite-Armed Bandit). Fix $A > 0$, and consider the class of finite-armed bandits with A arms and unit-variance Gaussian noise:

$$\mathcal{M} = \{M(\pi) = \mathcal{N}(f^M(\pi), 1) \mid f^M \in [0, 1]^A\}.$$

It is straightforward to verify that [Assumptions A.1](#) to [A.3](#) hold with $L_{\text{KL}}, V_{\mathcal{M}} \leq 4$ and $d_{\text{cov}} = O(A)$, $C_{\text{cov}} = O(1)$, and it can also be shown that, as long as $f^{M^*}(\pi_*) < 1$, we can bound $n_\varepsilon^* \leq c \cdot \frac{A^2}{\varepsilon(\Delta_{\min}^*)^4}$. In addition, we can bound $C_{\text{exp}}(\mathcal{M}^*, \varepsilon) \leq 4A/\varepsilon$, so [Proposition A.1](#) gives the following result.

Proposition A.2. *For the finite-armed bandit problem with A actions, there exists a universal constant $c > 0$ such that*

$$\overline{\text{aEC}}_\varepsilon^{\mathcal{M}}(\mathcal{M}^*) \leq c \cdot \frac{A^{15}}{\varepsilon^8(\Delta_{\min}^*)^{24}}. \quad (25)$$

We immediately obtain the following corollary to [Theorem A.2](#).

Corollary A.1. *For finite-armed bandits with Gaussian noise, as long as $f^{M^*}(\pi_*) < 1$, AE_*^2 has regret bounded by*

$$\mathbb{E}^{M^*}[\text{Reg}(T)] \leq (1 + \varepsilon)\mathbf{g}^* \cdot \log(T) + \text{poly}\left(A, \frac{1}{\varepsilon}, \frac{1}{\Delta_{\min}^*}, \log \log T\right) \cdot \log^{6/7}(T).$$

With more refined analyses, various works have achieved instance-optimal regret bounds for finite-armed bandits with tighter lower-order terms than [Corollary A.1](#) ([Garivier et al., 2016](#); [Kaufmann et al., 2016](#); [Lattimore, 2018](#); [Garivier et al., 2019](#)). We emphasize that [Corollary A.1](#) is a special case of a much more general result. In particular, we proved the bound on the Allocation-Estimation Coefficient, [Eq. \(25\)](#), using tools which hold for general classes (e.g. [Proposition A.1](#)). An analysis of AE_*^2 specialized to finite-armed bandits would likely yield a tighter result. \triangleleft

A.6.3. STRUCTURED BANDITS

Many bandit problems exhibit richer structure than the multi-armed bandit setting, and the study of these settings has been the focus of much of the recent work on instance-optimal learning. We next consider one such setting, that of structured bandits with bounded *eluder dimension* ([Russo and Van Roy, 2013](#)).

Example A.3 (Structured Bandits with Bounded Eluder Dimension). Consider a bandit problem with unit-variance Gaussian noise but where the means are now given by a *general function class* \mathcal{F} mapping from Π to $[0, 1]$:

$$\mathcal{M} = \{M(\pi) = \mathcal{N}(f(\pi), 1) \mid f \in \mathcal{F}\}. \quad (26)$$

For such general settings, we might hope to capture the complexity of learning in terms of generalized notions of dimension for \mathcal{F} . We consider one such notion here: the eluder dimension.

Definition A.5 (Eluder Dimension (Russo and Van Roy, 2013)). Let $\check{d}_E(\mathcal{F}, \varepsilon')$ denote the length of the longest sequence of actions $\{\pi_1, \dots, \pi_d\}$ such that, for each $n \leq d$, there exist functions $f, f' \in \mathcal{F}$ with $\sqrt{\sum_{i=1}^{n-1} (f(\pi_i) - f'(\pi_i))^2} \leq \varepsilon'$ but $f(\pi_n) - f'(\pi_n) > \varepsilon'$. The eluder dimension of function class \mathcal{F} at scale ε is then defined as $d_E(\mathcal{F}, \varepsilon) = \sup_{\varepsilon' \geq \varepsilon} \check{d}_E(\mathcal{F}, \varepsilon') \vee 1$.

The eluder dimension can be thought of as quantifying how easily a function class can be “explored”: evaluating a pair of functions on $d_E(\mathcal{F}, \varepsilon)$ points allows one to determine whether they are nearly identical over the entire space. It is known to be bounded for many standard classes—for example, for linear function classes with dimension d , $d_E(\mathcal{F}, \varepsilon) \leq \tilde{O}(d \cdot \log 1/\varepsilon)$ —and is also closely related to the *disagreement coefficient* (Foster et al., 2020). Furthermore, it can be shown to be a sufficient condition for bounded (worst-case) regret in general bandit problems (Russo and Van Roy, 2013). The following result shows that the eluder dimension bounds the Allocation-Estimation Coefficient.

Proposition A.3. For the structured bandit class \mathcal{M} considered in (26), we have

$$C_{\text{exp}}(\mathcal{M}^*, \delta) \leq \frac{16d_E(\mathcal{F}, \sqrt{\delta}/2)}{\delta}.$$

This implies that

$$\overline{\text{aec}}_\varepsilon^{\mathcal{M}}(\mathcal{M}^*) \leq \frac{16d_E(\mathcal{F}, \sqrt{\delta}/2)}{\delta} \quad \text{for scale } \delta = c \cdot \frac{\varepsilon^2 \Delta_\star^8}{d_E(\mathcal{F}, \frac{1}{2} \Delta_\star)^2} \quad \text{and } \Delta_\star := \min \left\{ \Delta_{\min}^\star, 1/n_{\varepsilon/36}^\star \right\},$$

where $c > 0$ is a universal constant.

Proposition A.3 highlights the ability of the Allocation-Estimation Coefficient to adapt to the inherent complexity of “exploring” for the model class under consideration. We henceforth abbreviate $d_E := d_E(\mathcal{F}, \sqrt{\delta}/2)$.

It is straightforward to show that Assumptions A.1 and A.2 are met in this setting with $L_{\text{KL}}, V_{\mathcal{M}} \leq 4$ (see Appendix G.2). Furthermore, Assumption A.3 can be shown to hold with d_{cov} scaling with the covering number of \mathcal{F} in the distance $d(f, f') = \sup_{\pi \in \Pi} |f(\pi) - f'(\pi)|$ and $C_{\text{cov}} = O(1)$. In general, it must be shown that n_ε^\star is bounded for each \mathcal{M}^* and class of interest; as we show in the following examples, it is bounded for standard structured bandit settings such as linear bandits. We have the following corollary to Theorem A.2.

Corollary A.2. In the structured bandit setting with bounded eluder dimension considered above, AE_\star^2 has regret bounded as

$$\mathbb{E}^{M^\star} [\mathbf{Reg}(T)] \leq (1 + \varepsilon) \mathbf{g}^\star \cdot \log(T) + \text{poly}(d_E, d_{\text{cov}}, \frac{1}{\varepsilon}, \frac{1}{\Delta_{\min}^\star}, n_{\varepsilon/36}^\star, \log \log T) \cdot \log^{6/7}(T).$$

To the best of our knowledge, Corollary A.2 is the first result to show that it is possible to obtain the instance-optimal rate in classes with bounded eluder dimension, with lower-order terms scaling only polynomially in the eluder dimension. More generally, Corollary A.2 illustrates that AE_\star^2 can adapt to the structural properties of the given model class, and achieve regret scaling with existing notions of intrinsic dimension. \triangleleft

We next consider two examples of structured bandits where it is known that the eluder dimension is bounded: linear bandits and generalized linear models. While these results are immediate given [Corollary A.2](#), the additional structure present in these settings allows us to obtain somewhat more explicit results.

Example A.4 (Linear Bandits). Consider the class of linear bandits with unit-variance Gaussian noise defined as

$$\mathcal{M} = \{M(\pi) = \mathcal{N}(\langle \theta, x_\pi \rangle, 1) \mid \theta \in \Theta\},$$

where $\Theta \subseteq \mathbb{R}^d$ is some convex set with ℓ_2 diameter $O(1)$ and $\mathcal{X} := \{x_\pi : \pi \in \Pi\} \subseteq \mathbb{R}^d$ is the arm set, which we assume has $\|x_\pi\|_2 \leq 1$ for all $\pi \in \Pi$. As in [Example A.3](#), [Assumptions A.1](#) and [A.2](#) are met in this setting with $L_{\text{KL}}, V_{\mathcal{M}} \leq 4$; furthermore, [Assumption A.3](#) is also met with $d_{\text{cov}} = O(d)$ and $C_{\text{cov}} = O(1)$. We then have the following bound on the AEC.

Proposition A.4. *For the linear bandit class \mathcal{M} defined above, we have*

$$\overline{\text{aec}}_\varepsilon^{\mathcal{M}}(\mathcal{M}^*) \leq c \cdot \frac{d^3}{\varepsilon^2} \cdot \left(\frac{1}{\Delta_{\min}^*} + n_{\varepsilon/36}^* \right)^8 \cdot \text{polylog} \left(d, \frac{1}{\varepsilon}, \frac{1}{\Delta_{\min}^*}, n_{\varepsilon/36}^* \right)$$

for a universal constant $c > 0$.

Using [Proposition A.4](#), we obtain the following corollary to [Theorem A.2](#).

Corollary A.3. *In the linear bandit setting defined above, AE_\star^2 has regret bounded as*

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq (1 + \varepsilon)g^* \cdot \log(T) + \text{poly} \left(d, \frac{1}{\varepsilon}, \frac{1}{\Delta_{\min}^*}, n_{\varepsilon/36}^*, \log \log T \right) \cdot \log^{6/7}(T).$$

[Corollary A.4](#) has lower-order terms scaling similarly to the best known lower-order terms in the linear bandit setting ([Tirinzoni et al., 2020](#); [Kirschner et al., 2021](#)). These works, however, develop algorithms which are specialized to the linear bandit setting, while [Corollary A.4](#) is the instantiation of a more general result designed for arbitrary general decision-making settings.

The following result shows that under general conditions, we can bound the parameter n_ε^* for linear bandits.

Proposition A.5 (Informal). *For linear bandits satisfying certain regularity conditions, n_ε^* is bounded by a polynomial function of problem parameters and a geometry-dependent term scaling with the structure of \mathcal{X} and Θ .*

The full statement of [Proposition A.5](#) is given in [Proposition G.2](#). The regularity condition for [Proposition A.5](#) requires primarily that θ^* lies sufficiently far within the interior of Θ (see [Appendix G.2.4](#) for further details). We remark that the guarantees given in both [Tirinzoni et al. \(2020\)](#) and [Kirschner et al. \(2021\)](#) scale with geometric parameters very similar to n_ε^* .

◁

Example A.5 (Generalized Linear Models). In the generalized linear model setting, we take the model class to be

$$\mathcal{M} = \{M(\pi) = \mathcal{N}(g(\langle \theta, x_\pi \rangle), 1) \mid \theta \in \Theta\},$$

where Θ and \mathcal{X} are as in [Example A.4](#), and $g(\cdot)$ is a known link function which is increasing and Lipschitz, but potentially nonlinear. Let g_{\max} and g_{\min} denote upper and lower bounds on the derivative of g , respectively:

$$g_{\max} := \max_{\theta \in \Theta, x \in \text{co}(\mathcal{X})} g'(\langle \theta, x \rangle) \quad \text{and} \quad g_{\min} := \min_{\theta \in \Theta, x \in \text{co}(\mathcal{X})} g'(\langle \theta, x \rangle).$$

As in the linear bandit setting, we can show that [Assumptions A.1](#) and [A.2](#) are both met with $L_{\text{KL}}, V_{\mathcal{M}} \leq 4$, and that [Assumption A.3](#) is also met with $d_{\text{cov}} = O(d)$ and $C_{\text{cov}} = O(g_{\max})$. Furthermore, under the same conditions as for linear bandits, n_{ε}^* can be bounded for generalized linear models exactly as for linear bandits, but with an additional scaling of $(\frac{g_{\max}}{g_{\min}})^2$. We then have the following.

Proposition A.6. *For the generalized linear model class \mathcal{M} defined above, we have*

$$\overline{\text{aec}}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}^*) \leq c \cdot \frac{d^3 g_{\max}^3}{\varepsilon^2 g_{\min}^3} \cdot \left(\frac{1}{\Delta_{\min}^*} + n_{\varepsilon/36}^* \right)^8 \cdot \text{polylog} \left(d, \frac{1}{\varepsilon}, \frac{1}{\Delta_{\min}^*}, n_{\varepsilon/36}^* \right)$$

for a universal constant $c > 0$.

Using [Proposition A.4](#), we obtain the following corollary to [Theorem A.2](#).

Corollary A.4. *In the generalized linear model setting defined above, AE_{\star}^2 has regret bounded as*

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq (1 + \varepsilon)g^* \cdot \log(T) + \text{poly} \left(d, \frac{g_{\max}}{g_{\min}}, \frac{1}{\varepsilon}, \frac{1}{\Delta_{\min}^*}, n_{\varepsilon/36}^*, \log \log T \right) \cdot \log^{6/7}(T).$$

To the best of our knowledge, this is the first result to obtain finite-time instance-optimality for generalized linear models with lower-order terms polynomial in problem parameters. \triangleleft

A.6.4. CONTEXTUAL BANDITS

The previous examples illustrate that AE_{\star}^2 is able to learn efficiently in a variety of structured bandit settings. We now show that it leads to new guarantees for finite-action contextual bandits with general function approximation.

Example A.6 (Contextual Bandits with Finitely Many Arms). Consider the contextual bandit setting with context set \mathcal{X} (which could be arbitrarily large) and action set \mathcal{A} such that $A := |\mathcal{A}| < \infty$. Let $p_{\mathcal{X}}$ denote the context distribution, which we assume is known to the learner. The learning protocol is then, for step $t = 1, 2, 3, \dots$:

1. Environment samples context $x^t \sim p_{\mathcal{X}}$.
2. Learner chooses action $a^t \in \mathcal{A}$, receives reward r_t .

We assume that $r^t = f^*(x^t, a^t) + w^t$ for $w^t \sim \mathcal{N}(0, 1)$, for some $f^* : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$. We assume as well that the learner is given access to a set of functions \mathcal{F} such that $f^* \in \mathcal{F}$.

To view this setting as a special case of the DMSO framework, we take the decision space to be the set $\Pi = (\mathcal{X} \rightarrow \mathcal{A})$ of all policies mapping from \mathcal{X} to \mathcal{A} , and take $\mathcal{O} = \mathcal{X}$ as the observation space. The learner's decision at round t is a policy π^t , and they receive a reward-observation pair $(r^t, o^t) = (r^t, x^t)$ under the process $x^t \sim p_{\mathcal{X}}, r^t \sim \mathcal{N}(f^*(x^t, \pi^t(x^t)), 1)$. The model class \mathcal{M} is the set of all instances of this form for $f^* \in \mathcal{F}$.

The following result shows that the Allocation-Estimation Coefficient is bounded by the number of actions A , and is *independent* of the size of the context space. See [Appendix G.3](#) for a proof.

Proposition A.7. *For the contextual bandit setting, we can bound*

$$C_{\text{exp}}(\mathcal{M}^*, \varepsilon) \leq \frac{4A}{\varepsilon},$$

which implies that

$$\overline{\text{aEC}}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}^*) \leq c \cdot \frac{A^3}{\varepsilon^2} \cdot \left(\frac{1}{\Delta_{\min}^*} + n_{\varepsilon/36}^* \right)^8,$$

for a universal constant $c > 0$.

As in the cases of bandits with bounded eluder dimension, n_{ε}^* must be bounded for each M^* and class \mathcal{F} of interest. It is straightforward to show, however, that [Assumptions A.1](#) and [A.2](#) are met in this setting with $L_{\text{KL}}, V_{\mathcal{M}} \leq 4$, and, furthermore, that [Assumption A.3](#) is also met with d_{cov} scaling as the covering number of \mathcal{F} in the distance $d(f, f') = \sup_{x \in \mathcal{X}, a \in \mathcal{A}} |f(x, a) - f'(x, a)|$, and $C_{\text{cov}} = O(1)$. We then have the following corollary.

Corollary A.5. *In the finit-action contextual bandit setting considered above, AE_{\star}^2 has regret bounded as*

$$\mathbb{E}^{M^*} [\mathbf{Reg}(T)] \leq (1 + \varepsilon) \mathbf{g}^* \cdot \log(T) + \text{poly}\left(A, d_{\text{cov}}, \frac{1}{\varepsilon}, \frac{1}{\Delta_{\min}^*}, n_{\varepsilon/36}^*, \log \log T\right) \cdot \log^{6/7}(T).$$

To the best of our knowledge, [Corollary A.5](#) is the first instance-optimal guarantee in the contextual bandit setting with general function approximation that obtains lower-order term scaling polynomially in problem parameters. Notably, the lower-order term scales independently of the size of the context space, $|\mathcal{X}|$. We anticipate that extending this result to contextual bandit settings that have large action spaces, but which exhibit additional structure allowing for efficient exploration (e.g., linearity), will be straightforward. \triangleleft

A.7. Application: Tabular Reinforcement Learning

As a final application of our results, we turn to the setting of episodic tabular reinforcement learning.

Episodic Markov decision processes. Recall that episodic reinforcement learning is a special case of the DMSO framework in which each model $M \in \mathcal{M}$ is an episodic Markov Decision Process (MDP) given by the tuple $M = (\mathcal{S}, \mathcal{A}, H, \{P_h^M\}_{h=1}^H, \{R_h^M\}_{h=1}^H, s_1)$. Here \mathcal{S} is a set of states, \mathcal{A} a set of actions, H the horizon, $P_h^M : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$ the probability transition kernel at step h , $R_h^M : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathbb{R}}$ the reward distribution at step h , and s_1 a deterministic initial state, which we take to be fixed across models. We assume that $R_h^M(s, a)$ is unit-variance Gaussian, and that $\mathbb{E}_{r_h \sim R_h^M(s, a)}[r_h] \in [0, 1/H]$.⁶

The decision space Π consists of non-stationary policies $\pi = (\pi_1, \dots, \pi_H)$, where $\pi_h : \mathcal{S} \rightarrow \mathcal{A}$. For a fixed policy π , an episode proceeds in an MDP M proceeds as follows. First, beginning from the initial state s_1 , we take action $a_1 \sim \pi_1(s_1)$, receive reward $r_1 \sim R_1^M(s_1, a_1)$, and transition to $s_2 \sim P_1^M(\cdot | s_1, a_1)$. This continues for H steps at which point the episode terminates and the process repeats. We define $f^M(\pi) := \mathbb{E}^{M, \pi} [\sum_{h=1}^H r_h]$ as the expected reward achieved over the entire episode under this process.

6. This assumption only serves to ensure that $f^M(\pi) \in [0, 1]$, in line with the convention for the rest of the paper. Our results continue to hold up to $\text{poly}(H)$ factors if $\mathbb{E}_{r_h \sim R_h^M(s, a)}[r_h] \in [0, 1]$.

Each round $t \in [T]$ in the DMSO framework corresponds to an episode in the underlying MDP M^* . At each round, the learner selects a policy π^t , and receives reward $r^t = \sum_{h=1}^H r_h^t$ and $o^t = (s_1^t, a_1^t, r_1^t, \dots, s_H^t, a_H^t, r_H^t)$, where $(s_1^t, a_1^t, r_1^t, \dots, s_H^t, a_H^t, r_H^t)$ is the trajectory that results from executing π^t in M^* for a single episode.

Tabular model class. In the tabular RL setting, it is assumed that $S := |\mathcal{S}|$ and $A := |\mathcal{A}|$ are both finite, and we take Π to be the set of all deterministic policies. In addition to assuming that \mathcal{M} consists of tabular MDPs, we restrict to the following subclass:

$$\mathcal{M}_{\text{tab}}(P_{\min}) := \left\{ M = (\mathcal{S}, \mathcal{A}, H, \{P_h^M\}_{h=1}^H, \{R_h^M\}_{h=1}^H, s_1) : \min_{s,a,s',h} P_h^M(s' | s, a) \geq P_{\min} \right\}. \quad (27)$$

While the assumption that $\min_{s,a,s',h} P_h^M(s' | s, a) \geq P_{\min}$ may be seen as restrictive, the guarantees we provide scale only with $\log \frac{1}{P_{\min}}$, so P_{\min} can be taken to be extremely small without affecting the result significantly.

Note that when our results are specialized to reinforcement learning, Δ_{\min}^* denotes the gap between the performance of the optimal policy, and the next-best *deterministic* policy. This quantity can be lower bounded in terms of other standard quantities including gaps in the rewards at each state and the transition probabilities.

Toward instantiating [Theorem A.2](#) in this tabular RL setting, we first provide a bound on the Allocation-Estimation Coefficient, which we establish by first bounding the Uniform Exploration Coefficient.

Proposition A.8. *For $\mathcal{M} \leftarrow \mathcal{M}_{\text{tab}}(P_{\min})$, we can bound⁷*

$$C_{\text{exp}}^H(\mathcal{M}^*, \varepsilon) \leq c \cdot \frac{SAH^2 \cdot \log^2 H}{\varepsilon^2}$$

for a universal constant $c > 0$, which implies that

$$\overline{\text{aEC}}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}^*) \leq c \cdot \frac{S^5 A^5 H^{14} \cdot \log^{10} H}{\varepsilon^4} \cdot \left(\frac{1}{\Delta_{\min}^*} + n_{\varepsilon/36}^* \right)^{24} \cdot \log^4 \frac{1}{P_{\min}}$$

for a universal constant $c > 0$.

Next, it can be shown that [Assumptions A.1](#) to [A.3](#) hold for $\mathcal{M} \leftarrow \mathcal{M}_{\text{tab}}(P_{\min})$ with constants (see [Appendix G.5](#)):

$$L_{\text{KL}} = V_{\mathcal{M}} = O(H \cdot \log 1/P_{\min}), \quad d_{\text{cov}} = O(S^2 AH), \quad C_{\text{cov}} = O(H/P_{\min}).$$

We then obtain the following corollary to [Theorem A.2](#).

Corollary A.6. *For $\mathcal{M} \leftarrow \mathcal{M}_{\text{tab}}(P_{\min})$ and ε , AE_{\star}^2 has regret bounded by*

$$\mathbb{E}^{M^*}[\text{Reg}(T)] \leq (1 + \varepsilon) \mathbf{g}^* \cdot \log(T) + \text{poly}\left(S, A, H, \frac{1}{\varepsilon}, \frac{1}{\Delta_{\min}^*}, \log \frac{1}{P_{\min}}, n_{\varepsilon/36}^*, \log \log(T)\right) \cdot \log^{6/7}(T).$$

7. Here C_{exp}^H denote the uniform exploration coefficient as defined in [Definition A.4](#), but with $D_{\text{KL}}(\cdot \| \cdot)$ replaced with $D_{\text{H}}^2(\cdot, \cdot)$. To prove [Proposition A.8](#), we show that a variant of [Proposition A.1](#) still holds with this alternate definition of $C_{\text{exp}}(\mathcal{M}, \varepsilon)$.

To our knowledge, [Corollary A.6](#) is the first guarantee for tabular RL that achieves the instance-optimal rate while obtaining lower-order terms that scale only polynomially in problem parameters. As noted previously, existing approaches to instance-optimal regret in tabular RL ([Ok et al., 2018](#); [Dong and Ma, 2022](#)) have lower-order terms that scale exponentially in problem parameters, and as a result are truly asymptotic in nature.

While [Corollary A.6](#) is stated in terms of n_ε^* for the sake of generality, as we show in [Appendix G.5](#), if M^* has rewards that are sufficiently small (in particular, if it satisfies $\mathbb{E}_{r_h \sim R_h^{M^*}(s,a)}[r_h] < 1/H^2$ for all (s, a, h)), then we can bound

$$n_\varepsilon^* \leq c \cdot \frac{\mathbf{g}^*}{\Delta_{\min}^*} \cdot \left(1 + \frac{\mathbf{g}^*}{\varepsilon(\Delta_{\min}^*)^2} \right).$$

In this case, [Corollary A.6](#) scales polynomially in all standard problem parameters.

Let us remark that the prior work of [Dong and Ma \(2022\)](#) does not require that $P_h^M(s' | s, a) \geq P_{\min}$ as we do (the work of [Ok et al. \(2018\)](#) only holds for ergodic MDPs, itself a very strong assumption). However, the lower-order term obtained in [Dong and Ma \(2022\)](#) scales polynomially in the inverse probability of observing the trajectory that occurs with minimum (non-zero) probability. In general, this will scale exponentially in H , and inversely with the probability of the transition with minimum (non-zero) probability occurring, that is $\min_{s,a,s',h:P_h^M(s'|s,a)>0} P_h^M(s'|s,a)$. Thus, while we must impose the stronger condition that all transitions occur with some probability P_{\min} , our bounds only scale logarithmically in this quantity, and polynomially in S, A , and H , a significant improvement over [Dong and Ma \(2022\)](#). Understanding whether it is possible to remove the additional restrictions we impose while still obtaining reasonable finite-time performance is an interesting direction for future work.

As far as we are aware, there is no prior work on instance-optimal algorithms for RL settings with general function approximation. While we have only instantiated [Theorem A.2](#) and AE_*^2 in the tabular RL setting, the tools we have developed can also be applied to RL with general model classes. Exploring the application of AE_*^2 to, for example, bilinear classes ([Du et al., 2021](#)) is an exciting avenue for future work.

A.8. Overview of Analysis

To close this section we briefly sketch the proof of the regret bound for AE^2 ([Theorem A.1](#)); the proof of the regret bound for AE_*^2 ([Theorem A.2](#)) builds on these ideas, but is slightly more involved. See [Appendix F](#) for full proofs.

Let us refer to the *exploit phase* as the subset of rounds t in which [Line 6](#) of AE^2 is reached, and refer to the *explore phase* as the subset of rounds in which [Line 8](#) is reached. We focus on bounding the regret in the explore phase—it can be shown ([Lemma F.1](#)) that in the exploit phase, where the if statement on [Line 6](#) is true, $\pi_{\widehat{M}^s} = \pi_*$ for all but $O(\log \log T)$ rounds, so that the regret incurred in this phase is at most $O(\log \log T)$.

Let s_T denote the total number of rounds in the explore phase up to time T . Fix an explore round $s \in [s_T]$. We bound the regret $\Delta^*(p^s)$ by considering three cases.

Case 1: $M^* \in \mathcal{M} \setminus \mathcal{M}_{\varepsilon/6}^{\text{gl}}(\lambda^s; n_{\max})$. In this case, λ^s is not an optimal (normalized) allocation for M^* , but we can use the AEC to argue that the information gained by the algorithm is large. In

particular, since p^s plays ω^s with probability at least $1 - q$, we can bound

$$\begin{aligned}
 & \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M^*(\pi))]]} \\
 & \leq \frac{1}{1 - q} \cdot \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim \omega^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M^*(\pi))]]} \\
 & \stackrel{(a)}{\leq} \frac{1}{1 - q} \cdot \min_{\lambda, \omega \in \Delta_{\Pi}} \sup_{M \in \mathcal{M} \setminus \mathcal{M}_{\varepsilon/6}^{\text{gl}}(\lambda; n_{\max})} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim \omega^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M(\pi))]]} \\
 & \lesssim \frac{1}{1 - q} \cdot \overline{\text{aEC}}_{\varepsilon/12}(\mathcal{M})
 \end{aligned}$$

as long as n_{\max} is chosen appropriately; here, (a) follows because $M^* \in \mathcal{M} \setminus \mathcal{M}_{\varepsilon/6}^{\text{gl}}(\lambda^s; n_{\max})$ by assumption in this case, and by the choice of λ^s and ω^s given in Eq. (18). Rearranging this gives

$$1 \lesssim \frac{1}{1 - q} \cdot \overline{\text{aEC}}_{\varepsilon/12}(\mathcal{M}) \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M^*(\pi))]].$$

This reflects that, when $M^* \in \mathcal{M} \setminus \mathcal{M}_{\varepsilon/6}^{\text{gl}}(\lambda^s; n_{\max})$, our choice of p^s ensures that $\widehat{M} \sim \xi^s$ and M^* can be distinguished, with the amount of information gained lower bounded by $O(\overline{\text{aEC}}_{\varepsilon/12}(\mathcal{M})^{-1})$

Adding and subtracting $\frac{2}{1 - q} \cdot \overline{\text{aEC}}_{\varepsilon/12}(\mathcal{M}) \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M^*(\pi))]]$ to $\Delta^*(p^s)$, and using that $\Delta^*(p^s) \leq 1$ always, we then have that the instantaneous regret in this case is bounded by

$$\begin{aligned}
 \Delta^*(p^s) &= \Delta^*(p^s) - \frac{2}{1 - q} \cdot \overline{\text{aEC}}_{\varepsilon/12}(\mathcal{M}) \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M^*(\pi))]] \\
 &\quad + \frac{2}{1 - q} \cdot \overline{\text{aEC}}_{\varepsilon/12}(\mathcal{M}) \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M^*(\pi))]] \\
 &\leq -1 + \frac{2}{1 - q} \cdot \overline{\text{aEC}}_{\varepsilon/12}(\mathcal{M}) \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M^*(\pi))]].
 \end{aligned}$$

Summing over s , it follows that the total regret in this case is bounded by

$$\sum_{s=1}^{s_T} \Delta^*(p^s) \cdot \mathbb{I}\{s \text{ in Case 1}\} \leq \frac{2}{1 - q} \cdot \overline{\text{aEC}}_{\varepsilon/12}(\mathcal{M}) \cdot \mathbf{Est}_{\text{KL}}(s_T) - \sum_{s=1}^{s_T} \mathbb{I}\{s \text{ in Case 1}\}.$$

Furthermore, since regret is always non-negative—that is, $\Delta^*(p) \geq 0$ for all p —rearranging this inequality leads to a bound on the total number of times this case can occur:

$$\sum_{s=1}^{s_T} \mathbb{I}\{s \text{ in Case 1}\} \leq \frac{2}{1 - q} \cdot \overline{\text{aEC}}_{\varepsilon/12}(\mathcal{M}) \cdot \mathbf{Est}_{\text{KL}}(s_T).$$

Critically, as given in (17), $\mathbf{Est}_{\text{KL}}(s_T)$ scales at most poly-logarithmically in s_T . Thus, as long as s_T is at most $O(\log T)$, the total regret incurred in this case (as well as the total number of times this case can occur), will be at most $O(\overline{\text{aEC}}_{\varepsilon/12}(\mathcal{M}) \cdot \log \log T)$.

Case 2: $M^* \in \mathcal{M}_{\varepsilon/6}^{\text{gl}}(\lambda^s; n_{\max})$ and $\pi_* \in \pi_{\widehat{M}^s}$. In this case, we have that λ^s is a Graves-Lai optimal allocation for M^* . Thus, it follows that

$$\Delta^*(\lambda^s) \leq (1 + \varepsilon/6)\mathbf{g}^*/n^* \quad \text{and} \quad \inf_{M \in \mathcal{M}^{\text{alt}}(M^*)} \mathbb{E}_{\pi \sim \lambda^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \geq (1 - \varepsilon/6)/n^*$$

for some $n^* \leq n_{\max}$. This implies that, for any $M \in \mathcal{M}^{\text{alt}}(M^*)$, we can bound

$$\begin{aligned} \Delta^*(p^s) &= \Delta^*(p^s) - (1 + \varepsilon)\mathbf{g}^* \mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] + (1 + \varepsilon)\mathbf{g}^* \mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \\ &\lesssim (1 + \varepsilon/6)\mathbf{g}^*/n^* - (1 + \varepsilon)(1 - \varepsilon/6)\mathbf{g}^*/n^* + (1 + \varepsilon)\mathbf{g}^* \mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \\ &\lesssim -\varepsilon\mathbf{g}^*/n^* + (1 + \varepsilon)\mathbf{g}^* \mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \\ &\leq -\varepsilon\mathbf{g}^*/n_{\max} + (1 + \varepsilon)\mathbf{g}^* \mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \end{aligned}$$

Since this bound holds uniformly for all $M \in \mathcal{M}^{\text{alt}}(M^*)$, it follows that the total regret in this case can be bounded as

$$\begin{aligned} \sum_{s=1}^{s_T} \Delta^*(p^s) \cdot \mathbb{I}\{s \text{ in Case 2}\} &\lesssim (1 + \varepsilon)\mathbf{g}^* \cdot \sum_{s=1}^{s_T} \inf_{M \in \mathcal{M}^{\text{alt}}(M^*)} \mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \cdot \mathbb{I}\{\pi_* \in \pi_{\widehat{M}^s}\} \\ &\quad - \frac{\varepsilon\mathbf{g}^*}{n_{\max}} \sum_{s=1}^{s_T} \mathbb{I}\{s \text{ in Case 2}\}. \end{aligned}$$

To bound this, the key observation is that, if we explore at round s , then it must be the case that, for all $\pi_{\widehat{M}^s} \in \pi_{\widehat{M}^s}$, there exists some $M \in \mathcal{M}^{\text{alt}}(\pi_{\widehat{M}^s})$ such that $\sum_{i=1}^{s-1} \mathbb{E}_{\widehat{M} \sim \xi^i} [\log \frac{\mathbb{P}^{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}^{M, \pi^i}(r^i, o^i)}] \leq \log(T \log T)$. Using [Assumption A.1](#) and [Assumption A.2](#) to move from $\widehat{M} \sim \xi^i$ to M^* and to relate the observed log-likelihood ratios to the KL divergence, we can furthermore show that

$$\begin{aligned} &\inf_{M \in \mathcal{M}^{\text{alt}}(M^*)} \sum_{s=1}^{s_T} \mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \cdot \mathbb{I}\{\pi_* \in \pi_{\widehat{M}^s}\} \\ &\lesssim \inf_{M \in \mathcal{M}^{\text{alt}}(M^*)} \sum_{s=1}^{s_T} \mathbb{E}_{\widehat{M} \sim \xi^s} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^s}(r^s, o^s)}{\mathbb{P}^{M, \pi^s}(r^s, o^s)} \right] + \sqrt{s_T} \cdot \mathbf{Est}_{\text{KL}}(s_T) \\ &\lesssim \log(T \log T) + \sqrt{s_T} \cdot \mathbf{Est}_{\text{KL}}(s_T). \end{aligned}$$

This allows us to bound

$$\begin{aligned} &(1 + \varepsilon)\mathbf{g}^* \cdot \sum_{s=1}^{s_T} \inf_{M \in \mathcal{M}^{\text{alt}}(M^*)} \mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \cdot \mathbb{I}\{\pi_* \in \pi_{\widehat{M}^s}\} \\ &\lesssim (1 + \varepsilon)\mathbf{g}^* \cdot \log T + \mathbf{g}^* \sqrt{s_T} \cdot \mathbf{Est}_{\text{KL}}(s_T). \end{aligned}$$

Thus, as long as $s_T = O(\log T)$, we can bound the total regret incurred in Case by $(1 + \varepsilon)\mathbf{g}^* \cdot \log T + o(\log T)$. Using that regret is always lower bounded by 0 in the same fashion as Case 1, we can further use this to bound the total number of times that Case 2 occurs by

$$\sum_{s=1}^{s_T} \mathbb{I}\{s \text{ in Case 2}\} \leq \frac{n_{\max}}{\varepsilon\mathbf{g}^*} \left(\mathbf{g}^* \cdot \log T + \mathbf{g}^* \sqrt{s_T} \cdot \mathbf{Est}_{\text{KL}}(s_T) \right).$$

The intuition for this case is that, since we are playing a Graves-Lai allocation for M^* , the regret will scale with g^* , the instance-optimal rate, and, furthermore, the allocation will allow us to distinguish M^* from alternatives $M \in \mathcal{M}^{\text{alt}}(M^*)$. Using that the total estimation error is bounded, and that we only enter the explore phase if there exists some $M \in \mathcal{M}^{\text{alt}}(M^*)$ that we cannot distinguish from M^* , this ultimately implies that the total number of times this phase occurs, and therefore the total regret incurred by this phase, is bounded.

Case 3: $M^* \in \mathcal{M}_{\varepsilon/6}^{\text{gl}}(\lambda^s; n_{\max})$ and $\pi_* \notin \pi_{\widehat{M}^s}$. In this case, we bound $\Delta^*(p^s)$ by adding and subtracting $\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(\widehat{M}(\pi) \parallel M^*(\pi))]]$ in the same fashion as Case 1. Since $\pi_* \notin \pi_{\widehat{M}^s}$, it can be shown that ξ^s must place $\Omega(\Delta_{\min})$ probability mass on $M \in \mathcal{M}^{\text{alt}}(M^*)$, allowing us to lower bound $\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(\widehat{M}(\pi) \parallel M^*(\pi))]]$ using the same reasoning as in Case 2. In total, we can show that the regret in this case is bounded by $O(\frac{g^*}{\Delta_{\min}} \cdot \mathbf{Est}_{\text{KL}}(s_T))$, and the number of times this case can occur is at most $O(\frac{n_{\max}}{\varepsilon \Delta_{\min}} \cdot \mathbf{Est}_{\text{KL}}(s_T))$.

Concluding the Proof. Combining all three cases, we have shown that the regret of the explore phase is bounded by

$$(1 + \varepsilon)g^* \cdot \log T + \widetilde{O}\left(\frac{1}{1 - q} \cdot \overline{\text{aEC}}_{\varepsilon/12}(\mathcal{M}) \cdot \mathbf{Est}_{\text{KL}}(s_T) + g^* \sqrt{s_T} \cdot \mathbf{Est}_{\text{KL}}(s_T) + \frac{g^*}{\Delta_{\min}} \cdot \mathbf{Est}_{\text{KL}}(s_T)\right).$$

Furthermore, using our bounds on the number of times each case can occur, one can show that $s_T = O(\log T)$. Since $\mathbf{Est}_{\text{KL}}(s_T)$ is at most polylogarithmic in s_T , it follows that the regret is bounded as

$$(1 + \varepsilon)g^* \cdot \log T + \widetilde{O}\left(\overline{\text{aEC}}_{\varepsilon/12}(\mathcal{M}) + \overline{\text{aEC}}_{\varepsilon/12}^{1/2}(\mathcal{M}) \cdot \log^{1/2} T\right),$$

as stated in [Theorem A.1](#).

Appendix B. Lower Bounds for Learning the Optimal Allocation

The AE^2 algorithm achieves instance-optimal regret by explicitly learning an ε -optimal Graves-Lai allocation for the underlying model M^* . In this section, we introduce an abstract formulation for the problem of learning an optimal allocation ([Appendix B.1](#)), and provide lower bounds which show that the Allocation-Estimation Coefficient is a fundamental limit for this task ([Appendix B.2](#)). We then present several examples illustrating lower bounds on the Allocation-Estimation Coefficient ([Appendix B.3](#)), and discuss how our lower bounds relate to the problem of minimizing regret ([Appendix B.4](#)).

Additional Notation. Throughout this section we will also make use of the following definition:

$$\Lambda(M; \varepsilon, n_{\max}) := \left\{ \lambda \in \Delta_{\Pi} : \exists n \in (0, n_{\max}] \text{ s.t. } \Delta^M(\lambda) \leq \frac{(1 + \varepsilon)g^M}{n}, \right. \quad (28)$$

$$\left. \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \mathbb{E}_{\pi \sim \lambda} [D_{\text{KL}}(M(\pi) \parallel M'(\pi))] \geq \frac{1 - \varepsilon}{n} \right\}.$$

That is, $\Lambda(M; \varepsilon, n_{\max})$ denotes the set of normalized allocations that are Graves-Lai optimal for M with tolerance ε , and have normalization factor at most n_{\max} . Unless otherwise stated, the results in this section do not make use of [Assumption 1.3](#) or [Assumption A.4](#).

B.1. Learning the Optimal Allocation: Minimax Formulation

We consider the following protocol, which captures the task of learning an optimal Graves-Lai allocation for an unknown model M^* .

- For $t = 1, \dots, T$, sample $\pi^t \sim p^t$ and observe (r^t, o^t) .
- Based on the entire history $\mathcal{H}^T = (\pi^1, r^1, o^1), \dots, (\pi^T, r^T, o^T)$, output a normalized allocation $\hat{\lambda} \in \Delta_{\Pi}$. The allocation may be randomized according to a distribution $q \in \Delta_{\Delta_{\Pi}}$.

We formalize an *algorithm* for this task as a pair $\mathbb{A} = (p, q)$, where $q(\cdot \mid \mathcal{H}^T)$ is the distribution over $\hat{\lambda}$ given the history, and $p = \{p^t\}_{t \in [T]}$ is a sequence of exploration distributions of the form $p^t(\cdot \mid \mathcal{H}^{t-1})$. We let $\mathbb{P}^{M, \mathbb{A}}(\cdot)$ denote the law of \mathcal{H}^T when M is the underlying model and \mathbb{A} is the algorithm, and let $\mathbb{E}^{M, \mathbb{A}}[\cdot]$ denote the corresponding expectation. The goal of the algorithm is to ensure that $\hat{\lambda}$ is an ε -optimal allocation for M^* with high probability, i.e.

$$\mathbb{P}^{M^*, \mathbb{A}}\left(\hat{\lambda} \in \Lambda(M^*; \varepsilon, n_{\max})\right) \geq 1 - \delta$$

for some failure probability $\delta > 0$ and normalization factor $n_{\max} > 0$

Intuitively, learning an optimal Graves-Lai allocation is closely related to achieving instance-optimal regret, but there are some subtle technical differences which we discuss in detail in the sequel. We study the former task because we find it to be more amenable to non-asymptotic lower bounds, and because it captures the behavior of “natural” algorithms such as AE^2 and essentially every existing asymptotically optimal algorithm we are aware of.

Minimax framework. To provide lower bounds on the complexity of learning an optimal Graves-Lai allocation, we consider a minimax framework. Our main quantity of interest will be:

$$T^{\text{gl}}(\mathcal{M}; \varepsilon, n_{\max}, \delta) = \inf_{\mathbb{A}} \inf \left\{ T \in \mathbb{N} \mid \mathbb{P}^{M, \mathbb{A}}\left(\hat{\lambda} \in \Lambda(M; \varepsilon, n_{\max})\right) \geq 1 - \delta, \forall M \in \mathcal{M} \right\}. \quad (29)$$

This represents the earliest time $T \in \mathbb{N}$ for which there exists an algorithm that learns an ε -optimal allocation with probability at least $1 - \delta$, and does so uniformly for all $M \in \mathcal{M}$. Recall that for our upper bounds ([Theorem A.1](#)), the Allocation-Estimation Coefficient gives a bound on the lower-order terms in the regret (reflecting the time required to learn an ε -optimal allocation) that holds *uniformly* for all models in the class \mathcal{M} . To understand the optimality of uniform bounds of this type, a minimax framework is natural. This framework also naturally complements recent non-asymptotic algorithms for linear models such as ([Tirinzi et al., 2020](#); [Kirschner et al., 2021](#)), where the complexity of exploration is captured by problem-dependent quantities such as the feature dimension, which are bounded uniformly for all models in the class. Nonetheless, exploring other notions of optimality (for example, instance-dependent complexity) for learning the Graves-Lai allocation is an interesting direction for future research.

Note that the quantity (29) does not place any constraint on the regret of the algorithm under consideration. It will also be useful to consider the notion

$$\begin{aligned} & T^{\text{gl}}(\mathcal{M}; \varepsilon, n_{\max}, \delta, R) \\ &= \inf_{\mathbb{A}} \inf \left\{ T \in \mathbb{N} \mid \mathbb{P}^{M, \mathbb{A}}\left(\hat{\lambda} \in \Lambda(M; \varepsilon, n_{\max})\right) \geq 1 - \delta, \mathbb{E}^{M, \mathbb{A}}[\mathbf{Reg}(T)] \leq R \cdot \log(T), \forall M \in \mathcal{M} \right\}, \end{aligned} \quad (30)$$

which captures the minimax complexity of learning the Graves-Lai allocation, subject to the constraint that the algorithm achieves logarithmic regret throughout the learning process. This notion is particularly well suited to complement the upper bounds achieved by algorithms such as AE^2 .

We remark that the restriction to allocations with normalization factor no more than n_{\max} in the definitions above is natural for several reasons. First, without such a restriction, the returned allocation can place an arbitrarily small amount of mass on informative actions, and an arbitrarily large amount of mass on the optimal action, since any Graves-Lai optimal allocation is still Graves-Lai optimal if the amount of mass on the optimal action is increased arbitrarily. While technically Graves-Lai optimal, such allocations do not reflect an allocation one could play over a finite time horizon in order to certify the optimal decision π_* —as any algorithm with finite-time guarantees must do—and therefore do not reflect the cost such an algorithm must pay to learn an allocation. Second, given knowledge of the optimal action, any Graves-Lai optimal allocation which has a normalization factor larger than n_{\max} can be transformed into a Graves-Lai optimal allocation with normalization factor n_{\max} by adjusting the mass on the optimal action, assuming n_{\max} is taken to be sufficiently large (in particular, as large as $n_{\max}(\mathcal{M}, \varepsilon)$; cf. Eq. (19)). Therefore, if we restrict our attention to the task to learning the Graves-Lai allocation for a subclass of models which agree on the optimal action (as we do in Theorem B.2), any lower bound for learning an allocation with normalization n_{\max} also applies to learning an unrestricted allocation, for large enough n_{\max} . Finally, AE^2 itself plays allocations with bounded normalization, which as we note in Appendix A.2, does not affect the optimality of its performance—allocations with bounded normalization are always sufficient.

B.2. Main Result

We state two lower bounds. The first scales with the version of the Allocation-Estimation Coefficient appearing in our upper bound (Theorem A.1), but leads to a lower bound on T^{gl} , while the second lower bound scales with the AEC for a restriction \mathcal{M} , but is exponentially stronger in the sense that it provides a similar lower bound on $\log(T^{\text{gl}})$. We remark briefly that the regularity conditions of Appendix A.1 are not required to hold here, unless otherwise stated.

Theorem B.1 (Main lower bound—weak variant). *Let $\varepsilon > 0$, $n_{\max} > 0$, and $\mathcal{M}_0 \subseteq \mathcal{M}$ be given, and set $\delta := \frac{\varepsilon}{2} \cdot \min\{1, \inf_{M \in \mathcal{M}_0} \mathbf{g}^M / n_{\max}\}$. Unless*

$$T > \frac{\delta}{8} \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0, \bar{M}),$$

any algorithm must have, for some $M \in \mathcal{M}_0$:

$$\mathbb{P}^{M, \mathbb{A}} \left[\widehat{\lambda} \notin \Lambda(M; \varepsilon, n_{\max}) \right] \geq \frac{\delta}{6}.$$

Stated equivalently, Theorem B.1 implies that for any $\varepsilon > 0$, if we set $\delta = c \cdot \varepsilon \cdot \inf_{M \in \mathcal{M}_0} \mathbf{g}^M / n_{\max}$ for sufficiently small numerical constant c , then

$$T^{\text{gl}}(\mathcal{M}; \varepsilon, n_{\max}, \delta) \geq \delta \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_{2\varepsilon}(\mathcal{M}_0, \bar{M}).$$

Remark B.1 (Choice of \mathcal{M}_0). *Theorem B.1* is stated with respect to an arbitrary subset, \mathcal{M}_0 , of \mathcal{M} . While one could simply choose $\mathcal{M}_0 \leftarrow \mathcal{M}$, in some cases it is advantageous to choose $\mathcal{M}_0 \subsetneq \mathcal{M}$. In particular, note that our lower bound scale with $\inf_{M \in \mathcal{M}_0} \mathbf{g}^M$. For classes \mathcal{M} where $\inf_{M \in \mathcal{M}} \mathbf{g}^M = 0$ —which will be the case, for example, if \mathcal{M} corresponds to a model class where reward means form a compact set—it is advantageous to restrict \mathcal{M}_0 so that it corresponds only to instances with $\mathbf{g}^M > 0$.

We remark as well that the restriction of the lower bound to \mathcal{M}_0 is somewhat analogous to our upper bound *Theorem A.2*, which provides guarantees in terms of the AEC of a restriction of \mathcal{M} , \mathcal{M}^* . We may therefore choose $\mathcal{M}_0 \leftarrow \mathcal{M}^*$ to obtain lower bounds matching the scaling of *Theorem A.2*. In the following section we provide several examples of how \mathcal{M}_0 can be chosen to yield intuitive lower bounds.

Lastly, let us mention that \mathcal{M}_0 may also be chosen to yield lower bounds that have a more instance-dependent flavor. For example, given some instance M^* , we could choose \mathcal{M}_0 to correspond to all instances identical to M^* up to a permutation of the decisions, in which case *Theorem B.1* will yield lower bounds on the performance of any algorithm on a permutation of M^* , rather than over the entire class.

Remark B.2. *The lower bound in Theorem B.1, for $\mathcal{M}_0 = \mathcal{M}$, scales with the quantity*

$$\sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_\varepsilon(\mathcal{M}, \bar{M}) \geq \sup_{\bar{M} \in \text{co}(\mathcal{M})} \text{aec}_\varepsilon(\mathcal{M}, \bar{M}) \geq \sup_{\xi \in \Delta_{\mathcal{M}}} \overline{\text{aec}}_\varepsilon(\mathcal{M}, \xi) = \overline{\text{aec}}_\varepsilon(\mathcal{M}),$$

which at first glance might appear to be larger than the version of the AEC appearing in our upper bounds. However, there is no contradiction, because these quantities can be shown to be equivalent (up to problem-dependent parameters) under the assumptions with which *Theorem A.1* is proven.

Our second lower bound yields a lower bound on $\log(T^{\text{gl}})$ as opposed to T^{gl} —a significantly stronger result—but scales with the AEC for a restricted class. To state the result, for $\bar{M} \in \mathcal{M}^+$ and $\mathcal{M}_0 \subseteq \mathcal{M}$, define

$$\mathcal{M}_0^{\text{opt}}(\bar{M}) = \{M \in \mathcal{M}_0 \mid \pi_M \subseteq \pi_{\bar{M}}, D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi)) = 0 \ \forall \pi \in \pi_{\bar{M}}\}.$$

This represents the set of models M where 1) the optimal decisions for M are also optimal for \bar{M} and 2) playing the an optimal decision reveals no information that can distinguish M and \bar{M} . Our second lower bound scales with the AEC for $\mathcal{M}_0^{\text{opt}}(\bar{M})$, and is restricted to algorithms with low regret.

Theorem B.2 (Main lower bound—strong variant). *Let $\varepsilon > 0$, $n_{\max} > 0$, and $\mathcal{M}_0 \subseteq \mathcal{M}$ be given, and define $\delta = \frac{\varepsilon}{2} \cdot \min\{1, \inf_{M \in \mathcal{M}_0} \mathbf{g}^M / n_{\max}\}$. Unless*

$$\sup_{M \in \mathcal{M}_0} \frac{\mathbf{g}^M}{\Delta_{\min}^M} \cdot \log(T) \geq \Omega(\delta^2) \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_{2\varepsilon}^M(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}),$$

there is no algorithm that simultaneously ensures that

1. $\mathbb{E}^{M, \mathbb{A}}[\mathbf{Reg}(T)] \leq 2 \cdot \mathbf{g}^M \log(T), \ \forall M \in \mathcal{M}_0.$
2. $\mathbb{P}^{M, \mathbb{A}}[\hat{\lambda} \notin \Lambda(M; \varepsilon, n_{\max})] \leq \frac{\delta}{12}, \ \forall M \in \mathcal{M}_0.$

Equivalently, [Theorem B.2](#) implies that for any $\varepsilon > 0$, if we set $\delta = c \cdot \varepsilon \cdot \inf_{M \in \mathcal{M}_0} \mathbf{g}^M / n_{\max}$ for a sufficiently small numerical constant c , then for $R := 2 \sup_{M \in \mathcal{M}_0} \mathbf{g}^M$,

$$\log(T^{\text{gl}}(\mathcal{M}; \varepsilon, n_{\max}, \delta, R)) \geq \frac{\delta^2 \cdot \inf_{M \in \mathcal{M}_0} \Delta_{\min}^M}{R} \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}).$$

B.3. Examples

As we discuss in the sequel, the dependence on the Allocation-Estimation Coefficient in our lower bounds (particularly [Theorem B.2](#)) qualitatively matches the dependence on the AEC in our main upper bounds, [Theorems A.1](#) and [A.2](#). We defer a detailed discussion comparing our upper and lower bounds for a moment, and make matters concrete by considering three examples: the informative arm example from the introduction ([Example 1.1](#)), multi-armed bandits, and tabular reinforcement learning. We defer proofs for all examples to [Appendix H.4](#).

Example B.1 (Searching for an Informative Arm (revisited)). Let \mathcal{M} be the class constructed in [Example 1.1](#) with parameters $N, A \in \mathbb{N}$ and $\beta \in (0, 1/2)$. Let \mathcal{M}_0 denote the restriction of \mathcal{M} to models with $\Delta_{\min}^M \geq \Delta$ for some $\Delta \in (0, 1/6)$ (so that π_M is unique), and $f^M(\pi_M) < 1$. As long as $N \geq 16$ and $\beta \log(1 + \beta A) \geq 2\Delta/(A - 1)$, we have

$$\sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}) \geq \frac{N}{4\beta}. \quad (31)$$

For this construction, we have $\Delta_{\min}^M \geq \Delta$ and $\Omega(1) \leq \mathbf{g}^M \leq O(\beta^{-1})$ for all $M \in \mathcal{M}_0$. Thus, [Theorem B.2](#) implies that any algorithm with $\mathbb{E}^{M, \mathbb{A}}[\mathbf{Reg}(T)] \leq 2\mathbf{g}^M \log(T)$ for all $M \in \mathcal{M}_0$ must fail to learn an ε -optimal allocation with probability $\delta = \Omega(\varepsilon/n_{\max})$ unless

$$\sup_{M \in \mathcal{M}} \frac{\mathbf{g}^M}{\Delta_{\min}^M} \cdot \log(T) \gtrsim \frac{\varepsilon^2}{n_{\max}^2} \cdot \frac{N}{\beta},$$

which implies that $\log(T^{\text{gl}}(\mathcal{M}_0; \varepsilon, n_{\max}, \delta, R)) \gtrsim \varepsilon^2/n_{\max}^2 \cdot N$ for $R = 2 \sup_{M \in \mathcal{M}} \mathbf{g}^M = O(\beta^{-1})$. Any Graves-Lai allocation need only take at most $O(\beta^{-1})$ pulls to eliminate alternative instances, so an appropriate choice of n_{\max} is $O(\beta^{-1})$, yielding $\log(T^{\text{gl}}(\mathcal{M}_0; \varepsilon, n_{\max}, \delta, R)) \gtrsim \beta^2 \varepsilon^2 \cdot N$. While the dependence on the parameter $n_{\max}, \varepsilon, \Delta > 0$ here is certainly loose, this formalizes the intuition sketched in the introduction: any algorithm that learns an optimal allocation must explore $\Omega(N)$ times, yet any algorithm that achieves near-instance-optimal regret $\mathbb{E}^{M, \mathbb{A}}[\mathbf{Reg}(T)] \leq 2\mathbf{g}^M \log(T) \lesssim \beta^{-1} \log(T)$ can play a sub-optimal decision no more than roughly $\beta^{-1} \log(T)/\Delta$ times, leading to the constraint that

$$\log(T) \gtrsim N.$$

We can also apply [Theorem B.1](#) to show that any algorithm must fail to learn an ε -optimal allocation with probability $\delta = \Omega(\varepsilon/n_{\max})$ unless $T \gtrsim \varepsilon/n_{\max} \cdot \frac{N}{\beta}$. \triangleleft

Example B.2 (Finite-Armed Bandit). Let $A \geq 6$ and $\Delta \in (0, 1/2)$ be given and set $\Pi = [A]$. Let \mathcal{M} be the set of all multi-armed bandit instances with Gaussian noise:

$$\mathcal{M} = \left\{ M(\pi) = \mathcal{N}(f^M(\pi), 1/2) \mid f^M \in [0, 1]^A \right\} \quad (32)$$

and let $\mathcal{M}_0 = \{M \in \mathcal{M} : |\pi_M| = 1, \Delta^M(\pi) \in [\Delta, 2\Delta] \text{ for } \pi \neq \pi_M, f^M(\pi_M) < 1\}$ denote the set of all bandit instances where suboptimal arms have gaps on order Δ . Then for all $\varepsilon \in (0, 1/32)$,

$$\sup_{\bar{M} \in \mathcal{M}} \text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}) \geq c \cdot \varepsilon^{-2} \cdot \frac{A}{\Delta^2}, \quad (33)$$

where $c > 0$ is an absolute constant.

It can be shown that, by the construction of \mathcal{M}_0 , $\mathbf{g}^M = \Theta(A/\Delta)$ for all $M \in \mathcal{M}_0$, so $\delta \propto \varepsilon \cdot \min\{1, \frac{A}{\Delta n_{\max}}\}$. [Theorem B.1](#) then implies that any algorithm must fail to learn an ε -optimal allocation with probability $\delta = \Omega(\varepsilon \cdot \min\{1, \frac{A}{\Delta n_{\max}}\})$ unless

$$T \gtrsim \min\left\{1, \frac{A}{\Delta n_{\max}}\right\} \cdot \frac{A}{\varepsilon \Delta^2}.$$

Any Graves-Lai allocation need only take $O(\frac{A}{\Delta^2})$ pulls to eliminate all alternate instances, so a reasonable choice of n_{\max} is therefore $O(\frac{A}{\Delta^2})$. With this choice of n_{\max} , we have $T^{\text{gl}}(\mathcal{M}; \varepsilon, n_{\max}, \delta) \gtrsim \frac{A}{\varepsilon \Delta}$. We remark that the scaling on all parameters here is natural. Intuitively, we would expect that we need to pull each arm at least once to learn a near-optimal allocation, yielding an $\Omega(A)$ scaling. In addition, note that for multi-armed bandits, the optimal allocation places mass $\propto \frac{1}{\Delta^2}$ on arms with gap Δ . To correctly estimate this proportion requires an accurate estimate of Δ , which becomes increasingly difficult as Δ becomes smaller, yielding an $\Omega(\frac{1}{\Delta})$ scaling. Finally, as we decrease ε , we require that the returned allocation becomes closer to a truly optimal allocation, and we would therefore expect an $\Omega(\frac{1}{\varepsilon})$ scaling.

Note that the result derived by applying [Theorem B.1](#) above only gives a lower bound on T . To obtain a lower bound on $\log(T)$, we combine [Theorem B.2](#) and [Eq. \(33\)](#) with $n_{\max} = O(\frac{A}{\Delta^2})$ as above, which implies that any algorithm with $\mathbb{E}^{M, \mathbb{A}}[\mathbf{Reg}(T)] \leq 2\mathbf{g}^M \log(T)$ for all $M \in \mathcal{M}_0$ must fail to learn an ε -optimal allocation with probability $\delta = \Omega(\varepsilon \cdot \min\{1, \frac{A}{\Delta n_{\max}}\})$ unless

$$\sup_{M \in \mathcal{M}_0} \frac{\mathbf{g}^M}{\Delta_{\min}^M} \cdot \log(T) \gtrsim A,$$

or equivalently $\log(T^{\text{gl}}(\mathcal{M}; \varepsilon, n_{\max}, \delta, R)) \gtrsim \frac{A}{\sup_{M \in \mathcal{M}_0} \mathbf{g}^M / \Delta_{\min}^M}$ for $R = 2 \sup_{M \in \mathcal{M}} \mathbf{g}^M$. To see why such scaling is natural, note that any algorithm which has $\mathbb{E}^{M, \mathbb{A}}[\mathbf{Reg}(T)] \leq 2\mathbf{g}^M \log(T)$ for each instance M can afford to explore (that is, play a suboptimal decision) at most $2 \frac{\mathbf{g}^M}{\Delta_{\min}^M} \log(T)$ times, or their regret could exceed $2\mathbf{g}^M \log(T)$. However, no algorithm has any hope of learning an optimal allocation unless they play every arm at least once, so taking at least A pulls of suboptimal arms seems unavoidable, and we therefore would expect that we must have $2 \frac{\mathbf{g}^M}{\Delta_{\min}^M} \log(T) \gtrsim A$, which is precisely the necessary scaling shown here.

We make two remarks on this $\log(T)$ lower bound. First, note that due to the presence of the δ^2 term in [Eq. \(33\)](#) (which we believe to be loose), the lower bound we derive by applying [Theorem B.2](#) does not scale with the parameter ε^{-1} , as one might hope. Second, as noted, for $M \in \mathcal{M}_0$ for \mathcal{M}_0 chosen as in [Example B.2](#), we have $\mathbf{g}^M = \Omega(A/\Delta)$, in which cases the dependence on A cancels, and the lower bound becomes trivial. Note that this is somewhat to be expected. [Theorem B.2](#) will give a trivial lower bound whenever $\sup_{M \in \mathcal{M}_0} \mathbf{g}^M$ is much larger than the AEC. Recall that in our

upper bound, [Theorem A.2](#), the leading order $g^* \cdot \log T$ term will dominate the lower-order terms once

$$g^* \cdot \log T \gtrsim \Omega(\text{aec}_\varepsilon(\mathcal{M})).$$

Therefore, if g^* is much larger than the AEC, [Theorem A.2](#) simply gives an upper bound scaling as approximately $g^* \cdot \log T$, as long as $\log T = \Omega(1)$. The interpretation in this setting is that the complexity of learning the Graves-Lai allocation is dominated by the regret incurred by *playing* the Graves-Lai allocation, which we know is necessary from [Proposition 1.1](#), and therefore we would not expect lower-order terms to be significant components of the regret, as reflected by [Theorem A.2](#). However, as we have already shown, In settings such as [Example B.1](#) where this is not the case and the AEC is much larger than g^* , [Theorem B.2](#) will give a non-trivial lower bound which reflects the difficulty of learning the optimal allocation. \triangleleft

Example B.3 (Tabular Reinforcement Learning). Let $S, A, H \in \mathbb{N}$ and $\Delta \in (0, 1/2)$ be given, and assume that $SA \geq 24$ and $H \geq \log_2(S/2)$. Let \mathcal{M} be the set of all tabular MDPs with 1) $|\mathcal{S}| = S$, $|\mathcal{A}| = A$ and horizon H , and 2) Gaussian rewards with variance $\sigma^2 = 1/2$ (cf. [Appendix A.7](#)). Let \mathcal{M}_0 be the result of restricting \mathcal{M} in the same fashion as [Example B.2](#). Then for all $\varepsilon \in (0, 1/32)$,

$$\sup_{\bar{M} \in \mathcal{M}} \text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}) \geq c \cdot \varepsilon^{-2} \cdot \frac{SA}{\Delta^2}, \quad (34)$$

where $c > 0$ is an absolute constant.

It can be shown that, by the construction of \mathcal{M}_0 , $g^M \geq \Omega(SA/\Delta)$ for all $M \in \mathcal{M}_0$. Thus, analogous to the multi-armed bandit example, [Theorem B.1](#) and [Eq. \(34\)](#) imply that any algorithm must fail to learn an ε -optimal allocation with probability $\delta = \Omega(\varepsilon \cdot \min\{1, \frac{SA}{\Delta n_{\max}}\})$ unless

$$T \gtrsim \min\left\{1, \frac{SA}{\Delta n_{\max}}\right\} \cdot \frac{SA}{\varepsilon \Delta^2}.$$

Choosing $n_{\max} = O(\frac{SA}{\Delta^2})$ gives $T^{\text{gl}}(\mathcal{M}; \varepsilon, n_{\max}, \delta) \gtrsim \frac{SA}{\varepsilon \Delta}$. With this same choice of n_{\max} , [Theorem B.2](#) implies that any algorithm with $\mathbb{E}^{M, A}[\mathbf{Reg}(T)] \leq 2g^M \log(T)$ for all $M \in \mathcal{M}_0$ must fail to learn an ε -optimal allocation with probability $\delta = \Omega(\varepsilon \Delta)$ unless

$$\sup_{M \in \mathcal{M}_0} \frac{g^M}{\Delta_{\min}^M} \cdot \log(T) \gtrsim SA,$$

or equivalently $\log(T^{\text{gl}}(\mathcal{M}; \varepsilon, n_{\max}, \delta, R)) \gtrsim \frac{SA}{\sup_{M \in \mathcal{M}_0} g^M / \Delta_{\min}^M}$ for $R = 2 \sup_{M \in \mathcal{M}_0} g^M$. \triangleleft

B.4. Discussion and Interpretation

Our lower bounds show that the Allocation-Estimation Coefficient serves as a fundamental limit on the sample complexity required to learn an approximate Graves-Lai allocation. In particular, they capture phenomena such as the necessity of searching for an informative arm in [Example 1.1](#) that are missed by purely asymptotic analyses. To the best of our knowledge, our lower bounds represent the first attempt to systematically understand the sample complexity of learning the Graves-Lai allocation in a general decision making framework. As such, they are somewhat coarse (in particular, the dependence on parameters such as ε , n_{\max} , and Δ_{\min}^M is almost certainly loose), and they are best thought of as a starting point for further research. In what follows, we provide additional interpretation of the results, and highlight some of the most interesting remaining questions.

Regret versus learning the optimal allocation. [Theorems B.1](#) and [B.2](#) lower bound the sample complexity required to learn an ε -optimal Graves-Lai allocation. Intuitively, this task is closely related to achieving instance-optimal regret. Our analysis of AE^2 shows that it is *sufficient*, and many prior works aim to directly estimate the optimal allocation as well. However, it is unclear to what extent learning the optimal allocation is *necessary* to achieve instance-optimal regret.

In more detail, it is quite straightforward to show that if an algorithm achieves instance-optimal regret, its empirical frequencies act as an optimal allocation *in expectation*.

Lemma B.1. *Let $\varepsilon \in (0, 2)$, and suppose that [Assumption A.4](#) holds. Fix $T \in \mathbb{N}$ and consider an algorithm \mathbb{A} such that for all $M \in \mathcal{M}$,*

$$\mathbb{E}^{M, \mathbb{A}}[\mathbf{Reg}(T)] \leq (1 + \varepsilon)\mathbf{g}^M \cdot \log(T).$$

For each $M \in \mathcal{M}$, define $\eta^M \in \mathbb{R}_+^\Pi$ via $\eta^M(\pi) = \mathbb{E}^{M, \mathbb{A}}\left[\frac{T(\pi)}{\log(T)}\right]$, where $T(\pi)$ denotes the number of pulls of decision π , and define $\lambda^M = \eta^M / \|\eta^M\|_1$. Then if

$$\log(T) \geq \frac{6}{\varepsilon} \log\left(\sup_{M \in \mathcal{M}} \frac{2\mathbf{g}^M}{\Delta_{\min}^M} \cdot \log(T)\right),$$

we have that for all $M \in \mathcal{M}$,

$$\lambda^M \in \Lambda(M; \varepsilon). \quad (35)$$

This result gives a guarantee on the expected frequencies of any instance-optimal algorithm, but does not give any guarantee for the realized frequencies. As such, without further assumptions on the algorithm under consideration, it is unclear whether instance-optimal regret implies that it is possible to learn an optimal allocation with high or even *constant* probability. We cannot currently rule out the existence of pathological algorithms for which η^M is optimal in expectation, yet the empirical arm frequencies deviate from the mean with moderate probability. Nonetheless, if one is willing to make stronger assumptions on the algorithm under consideration—in particular, that the *second moment* of regret is controlled—then it is possible to derive lower bounds on regret directly.

Theorem B.3 (Simplified version of [Theorem H.3](#)). *Let the time horizon $T \in \mathbb{N}$ and $\varepsilon \in (0, 1/2)$ be given, and suppose that [Assumptions A.2](#) and [A.4](#) hold. Suppose there exists an algorithm \mathbb{A} with the property that for all $M \in \mathcal{M}$: 1) $\mathbb{E}^{M, \mathbb{A}}[\mathbf{Reg}(T)] \leq (1 + \varepsilon)\mathbf{g}^M \log(T)$, 2) $\sqrt{\mathbb{E}^{M, \mathbb{A}}[(\mathbf{Reg}(T))^2]} \leq 2\mathbf{g}^M \log(T)$, and 3) for all $\pi \in \Pi$, if $\mathbb{E}^{M, \mathbb{A}}[T(\pi)] \neq 0$, then $\mathbb{E}^{M, \mathbb{A}}[T(\pi)] \geq 1$. Then if we define $\delta = \varepsilon \cdot \min\{1, \inf_{M \in \mathcal{M}_0} \frac{\mathbf{g}^M}{3\mathbf{g}^M / \Delta_{\min}^M + n_\varepsilon^M}\}$, it must be the case that*

$$\log^3(T) \geq \frac{\delta^2}{C} \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_{4\varepsilon}^{\mathcal{M}}(\mathcal{M}^{\text{opt}}(\bar{M}), \bar{M}).$$

for $C \leq O\left(\left(\sup_{M \in \mathcal{M}} \frac{\mathbf{g}^M}{\Delta_{\min}^M}\right)^4 \cdot \frac{V_{\mathcal{M}}^2 \log(\delta^{-1})}{\varepsilon^2}\right)$.

See [Appendix H.5](#) for a full statement and details. The idea behind the proof is to 1) show (via robust mean estimation) that any instance-optimal algorithm with well-behaved tails can be used to estimate the optimal allocation with high probability (with a small blowup in time horizon), and then 2) appeal to [Theorem B.2](#). More work is required to understand whether 1) we can prove lower bounds on regret directly, and 2) whether it is possible to show that low regret and learning the optimal allocation are equivalent in a stronger sense.

Comparing upper and lower bounds. Keeping the differences between regret minimization and learning the optimal allocation in mind, let us highlight that the lower bound on T provided by [Theorem B.2](#) seems to qualitatively match the upper bound from [Theorems A.1](#) and [A.2](#). In particular, ignoring problem-dependent parameters and polylogarithmic factors, the upper bound [Theorem A.1](#) scales, for every model $M \in \mathcal{M}$, as

$$\mathbb{E}^M[\mathbf{Reg}(T)] \leq (1 + \varepsilon)g^M \log(T) + \tilde{O}^+(\text{aec}_\varepsilon(\mathcal{M})).$$

In order for this bound to simplify to, say,

$$\mathbb{E}^M[\mathbf{Reg}(T)] \leq (1 + 2\varepsilon)g^M \log(T),$$

we need

$$g^M \cdot \log(T) \geq \tilde{\Omega}^+(1) \cdot \frac{\text{aec}_\varepsilon(\mathcal{M})}{\varepsilon},$$

which has similar scaling to the lower bound

$$\log(T) \gtrsim \tilde{\Omega}^+(1) \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M})$$

from [Theorem B.2](#). As discussed in the prequel, the former result is concerned with regret, while the latter considers the task of learning the optimal allocation, but the scaling $\log(T) \gtrsim \text{aec}_\varepsilon(\mathcal{M})$ seems to be fundamental for both. Of course, beyond the gap between regret and learning the allocation, there is still much room to improve the dependence on problem-dependent parameters in both results.

Comparing [Theorem B.1](#) and [Theorem B.2](#). [Theorem B.1](#) and [Theorem B.2](#) exhibit an interesting dichotomy: [Theorem B.1](#) places no constraints on the regret of the algorithm under consideration, and gives a lower of the form $T \gtrsim \tilde{\Omega}^+(1) \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}_0, \bar{M})$, while [Theorem B.2](#) gives a lower bound of the form $\log(T) \gtrsim \tilde{\Omega}^+(1) \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M})$, or equivalently $T \gtrsim \exp(\tilde{\Omega}^+(1) \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_\varepsilon(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}))$; the latter lower bound is exponentially stronger, with the caveat that 1) the class \mathcal{M}_0 is replaced with the subclass $\mathcal{M}_0^{\text{opt}}(\bar{M})$ and 2) [Theorem B.2](#) assumes that the algorithm achieve nearly-instance optimal regret for every model in \mathcal{M}_0 ($\mathbb{E}^{M, \Delta}[\mathbf{Reg}] \leq 2 \cdot g^M \log(T)$). In what follows, we argue that this tradeoff is fundamental.

- First, let us consider the role of the assumption $\mathbb{E}^{M, \Delta}[\mathbf{Reg}] \leq 2 \cdot g^M \log(T)$. Without this assumption, the lower bound from [Theorem B.1](#) is qualitatively tight: if the algorithm explores optimally for every round $t \in [T]$, it gains roughly $(\sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}_0, \bar{M}))^{-1}$ units of information per round, which is sufficient to identify an optimal allocation as soon as $T \gtrsim \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}_0, \bar{M})$.
- On the other hand, if we require that $\mathbb{E}^{M, \Delta}[\mathbf{Reg}] \leq 2 \cdot g^M \log(T)$, then for each $M \in \mathcal{M}_0$, the algorithm can afford to explore (i.e., play a non-optimal action) at most $2 \cdot \frac{g^M}{\Delta_{\min}^M} \log(T)$ times. This changes the “effective” time horizon for exploration to $T' = 2 \cdot \frac{g^M}{\Delta_{\min}^M} \log(T)$, but only if we restrict to models for which playing an optimal decision gives no information. This is precisely what the subclass $\mathcal{M}_0^{\text{opt}}(\bar{M})$ captures: models $M \in \mathcal{M}_0$ for which decisions that are optimal for \bar{M} lead to no information. Combining these insights leads to the lower bound $T' \approx \frac{g^M}{\Delta_{\min}^M} \log(T) \gtrsim \Omega(1) \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M})$ in [Theorem B.2](#).

We remark in passing that the definition of $\mathcal{M}_0^{\text{opt}}(\bar{M})$, which places the constraint that $D_{\text{KL}}(\bar{M}(\pi) \| M(\pi)) = 0$, $\forall \pi \in \pi_{\bar{M}}$ is somewhat coarse. We expect that both [Theorem B.2](#) and [Theorem A.1/Theorem A.2](#) can be improved to scale with

$$\mathcal{M}_0^{\text{opt}}(\bar{M}; \alpha) = \{M \in \mathcal{M}_0 \mid \pi_M \subseteq \pi_{\bar{M}}, D_{\text{KL}}(\bar{M}(\pi) \| M(\pi)) \leq \alpha^2 \forall \pi \in \pi_{\bar{M}}\}$$

for $\alpha \approx 1/\sqrt{T}$; the intuition is that we get $\Omega(T)$ rounds worth of information on π_M for “free”, which facilitates accurate estimation.

Minimax versus instance-dependent lower bounds. As mentioned in the prequel, our lower bounds have a (constrained) minimax flavor. Specifically, [Theorem B.2](#) shows that if T is not sufficiently large, then for any algorithm, there must exist a “worst-case” model $M \in \mathcal{M}$ for which the algorithm either 1) fails to achieve (approximately) instance-optimal regret or 2) fails to learn an ε -optimal Graves-Lai allocation. While this is quite different from a classical minimax analysis, and certainly is closely connected to instance-optimality, an interesting direction for future work is to develop a fully instance-dependent understanding of the complexity of learning Graves-Lai allocations.

Appendix C. Additional Related Work

In this section, we discuss further related work not already covered in detail.

Asymptotic guarantees for general decision making. For the general decision making framework we consider, which allows for arbitrary model classes and subsumes structured bandits and reinforcement, the only prior works we are aware of that achieve the instance-optimal lower bound from [Graves and Lai \(1997\)](#) are [Komiyama et al. \(2015\)](#), which restricts to finite observation spaces, and [Dong and Ma \(2022\)](#), which restricts to finite decision spaces; these works do not provide non-asymptotic guarantees.

Many works provide purely asymptotic instance-optimality guarantees for more specialized settings, including multi-armed bandits ([Lai and Robbins, 1985](#); [Garivier et al., 2016](#); [Lattimore, 2018](#); [Garivier et al., 2019](#)), linear bandits ([Lattimore and Szepesvari, 2017](#); [Hao et al., 2019, 2020](#)), and general structured bandits ([Burnetas and Katehakis, 1996](#); [Magureanu et al., 2014](#); [Combes et al., 2017](#); [Van Parys and Golrezaei, 2020](#); [Degenne et al., 2020b](#)). Some of these works do provide non-asymptotic bounds on regret, but these results generally have lower-order terms that scale linearly in $|\Pi|$, which renders them vacuous until $\log(T) \gtrsim |\Pi|$; we consider such results to be asymptotic in spirit. Along these lines, it is worth discussing [Jun and Zhang \(2020\)](#), which provides non-asymptotic guarantees for structured bandits in which the lower-order terms scale with a quantity K_ψ that aims to capture the number of “effective arms”. While this quantity can improve over $|\Pi|$ in certain situations, it is not clear whether it is well behaved for standard classes of interest (e.g., linear bandits).

Non-asymptotic guarantees for linear bandits. For linear bandits, a number of recent works provide non-asymptotic instance-optimal regret bounds in which lower order terms scale only with the dimension d rather than the number of decisions $|\Pi|$ ([Tirinzoni et al., 2020](#); [Kirschner et al., 2021](#)). These results take advantage of the specialized geometric structure of the linear bandit setting (e.g., existence of optimal design) for exploration, and cannot be directly adapted to general function approximation, but our results can be viewed as generalizing these guarantees.

Reinforcement learning. For reinforcement learning, a number of works—mostly focusing on tabular settings or linear function approximation—provides non-asymptotic guarantees that are instance-dependent, but not necessarily instance-optimal (Simchowitz and Jamieson, 2019; Al Marjani and Proutiere, 2021; Dann et al., 2021; Al Marjani et al., 2021; Wagenmaker et al., 2022b; Wagenmaker and Jamieson, 2022; Wagenmaker and Pacchiano, 2022). For instance-optimality, the results we are aware of are the classical work of Agrawal et al. (1988) and very recent work Dong and Ma (2022), which provides asymptotic guarantees for finite-horizon tabular RL, Ok et al. (2018), which provides asymptotic guarantees for an infinite-horizon setting under ergodicity assumptions, and Tirinzoni et al. (2022) which provides PAC guarantees for deterministic MDPs (though we note that the guarantee of Tirinzoni et al. (2022) is also achieved, up to H factors, by Wagenmaker and Jamieson (2022)).

Instance-optimal PAC guarantees. Our discussion has largely centered on regret, which is the focus of our work. For the PAC setting, where the goal is to identify the optimal decision (or a near-optimal decision) as quickly as possible, a number of recent works have employed similar techniques to derive instance-optimal algorithms for settings such as multi-armed and structured bandits (Kaufmann et al., 2016; Garivier and Kaufmann, 2016; Russo, 2016; Degenne and Koolen, 2019; Degenne et al., 2019, 2020a). While many of these works are asymptotic in nature, in specialized settings such as multi-armed bandits (Jamieson et al., 2014), linear bandits (Fiez et al., 2019; Katz-Samuels et al., 2020), and linear dynamical systems (Wagenmaker et al., 2021), recent work has shown that the optimal rates are achievable in finite-time.

Complexity of learning the Graves-Lai allocation. The Allocation-Estimation Coefficient aims to capture the sample complexity required to learn an ε -optimal Graves-Lai allocation. To the best of our knowledge, our work is the first to study the complexity of learning the allocation with general function approximation, but a small body of work has studied the complexity in simple settings such as top- k bandits (Simchowitz et al., 2017; Chen et al., 2017), and graph bandits (Marinov et al., 2022a,b).

Minimax regret. While the focus of this work has been on instance-optimality, a large body of work exists on *minimax* optimality, where the goal is to perform optimally on the *hardest* instance within a class. This line of work has established worst-case optimal (or nearly optimal) rates in settings such as multi-armed bandits (Auer et al., 2002; Audibert and Bubeck, 2009), linear bandits (Dani et al., 2008; Abbasi-Yadkori et al., 2011), tabular reinforcement learning (Dann et al., 2019; Zhang et al., 2021), and reinforcement learning with function approximation (Zhou et al., 2021; Du et al., 2021). The recent line of work Foster et al. (2021, 2022b, 2023) shows that, in the interactive decision-making setting considered in this work, the minimax-optimal rates are governed by a quantity known as the *Decision-Estimation Coefficient*. While our work takes inspiration and bears some similarity with this work, we remark that the techniques necessary to establish instance-optimality are significantly more intricate. It is also worth stating that it is always possible to bound the DEC by the AEC; the converse, however, is not true.

Part II

Proofs

Appendix D. Additional Notation

Mathematical Notation	Definition
$D_{\text{KL}}(\cdot \ \cdot)$	KL divergence
$D_{\text{H}}(\cdot, \cdot)$	Hellinger distance
$D_{\text{TV}}(\cdot, \cdot)$	Total variation distance
$D(\cdot \ \cdot)$	General divergence
$\Delta_{\mathcal{X}}$	Set of probability distributions over \mathcal{X}
DMSO Notation	
M	Model
M^*	Ground truth model
\mathcal{M}	Set of models
\mathcal{M}_0	Arbitrary subset of \mathcal{M}
π, Π	Decision π , set of all decisions Π
r, \mathcal{R}	Reward r , set of all rewards \mathcal{R}
o, \mathcal{O}	Observation o , set of all observations \mathcal{O}
$\mathbb{E}^{M, \pi}[\cdot], \mathbb{P}^{M, \pi}[\cdot]$	Expectation and distribution of $(r, o) \sim M(\pi)$
$\mathbb{E}^M[\cdot], \mathbb{P}^M[\cdot]$	Expectation and distribution induced over histories on M
$f^M(\pi)$	Expected reward playing π on M , $f^M(\pi) = \mathbb{E}^{M, \pi}[r]$
π_M	Optimal decision of model M , $\pi_M \in \arg \max_{\pi \in \Pi} f^M(\pi)$
$\boldsymbol{\pi}_M$	Set of optimal decisions of model M
$\Delta^M(\pi)$	Gap of decision π on model M , $\Delta^M(\pi) = f^M(\pi_M) - f^M(\pi)$
Δ_{\min}^M	Minimum gap on model M (see Eq. (2))
\mathcal{M}^+	Set of all possible models, $\mathcal{M}^+ = \{M : \Pi \rightarrow \Delta_{\mathcal{R} \times \mathcal{O}} \mid f^M(\pi) \in [0, 1]\}$
$\text{Reg}(T)$	Regret after T rounds (see Eq. (1))
L_{KL}	Lipschitz constant of KL divergence (see Assumption A.1)
$V_{\mathcal{M}}$	Sub-gaussian parameter of log-likelihood ratio (see Assumption A.2)
$N_{\text{cov}}(\mathcal{M}, \rho, \mu)$	(ρ, μ) covering number of \mathcal{M} (see Definition A.1)
$d_{\text{cov}}, C_{\text{cov}}$	Bounds on covering number (see Assumption A.3)
n_{ε}^M	Information content of optimal decision on M with tolerance ε (see Definition A.2)
$n_{\varepsilon}^{\mathcal{M}}$	Maximum information content of optimal decision on \mathcal{M} , $n_{\varepsilon}^{\mathcal{M}} = \sup_{M \in \mathcal{M}} n_{\varepsilon}^M$
$\mathcal{M}_{x, y}$	$\mathcal{M}_{x, y} = \{M \in \mathcal{M} : \Delta_{\min}^M \geq x, n_{\varepsilon}^M \leq y\}$
Δ_{\star}	$\Delta_{\star} = \min\{\Delta_{\min}^{\star}, 1/n_{\varepsilon/36}^{\star}\}$
\mathcal{M}^*	Restriction of \mathcal{M} induced by M^* (see Eq. (24))

Graves-Lai Notation	
$\text{glc}(\mathcal{M}, M)$	Graves-Lai Coefficient for class \mathcal{M} , model M (see Eq. (3))
$\mathbf{g}^M, \mathbf{g}^*$	Graves-Lai Coefficient for model M, M^* ; $\mathbf{g}^M = \text{glc}(\mathcal{M}, M)$, $\mathbf{g}^* = \text{glc}(\mathcal{M}, M^*)$
$\underline{\mathbf{g}}^{\mathcal{M}}$	Minimum non-zero Graves-Lai Coefficient on \mathcal{M} , $\underline{\mathbf{g}}^{\mathcal{M}} := \min_{M \in \mathcal{M}: \mathbf{g}^M > 0} \mathbf{g}^M$
$\mathcal{M}^{\text{alt}}(M)$	Alternate set for model M , $\mathcal{M}^{\text{alt}}(M) = \{M' \in \mathcal{M} \mid \boldsymbol{\pi}_M \cap \boldsymbol{\pi}_{M'} = \emptyset\}$
$\mathcal{M}_{\text{alt}}^*$	Alternate set for M^* , $\mathcal{M}_{\text{alt}}^* = \mathcal{M}^{\text{alt}}(M^*)$
η	Allocation, $\eta \in \mathbb{R}_+^{\Pi}$
λ	Normalized allocation, $\lambda \in \Delta_{\Pi}$
ω	Exploration distribution, $\omega \in \Delta_{\Pi}$
$\Lambda(M; \varepsilon)$	Set of ε -optimal normalized Graves-Lai allocations for model M (see Eq. (8))
$\Lambda(M; \varepsilon, n_{\max})$	Set of ε -optimal normalized Graves-Lai allocations for model M with normalization factor at most n_{\max} (see Eq. (28))
$\mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda)$	Models for which λ is an ε -optimal Graves-Lai allocation (see Eq. (9))
$\mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda; n_{\max})$	Models for which λ is an ε -optimal Graves-Lai allocation with normalization factor at most n_{\max} (see Eq. (37))
$I^M(\eta; \mathcal{M})$	Information content of η on M with respect to \mathcal{M} (see Eq. (36))
$I^M(\eta)$	$I^M(\eta) = I^M(\eta; \mathcal{M})$
$T^{\text{gl}}(\mathcal{M}_0; \varepsilon, n_{\max}, \delta)$	Minimum time to learn Graves-Lai allocation over \mathcal{M}_0 (see Eq. (29))
AEC Notation	
$\text{aec}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}_0, \bar{M})$	AEC with tolerance ε , model set \mathcal{M}_0 , reference model \bar{M} (see Eq. (10))
$\text{aec}_{\varepsilon}(\mathcal{M}, \bar{M})$	$\text{aec}_{\varepsilon}(\mathcal{M}, \bar{M}) = \text{aec}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}, \bar{M})$
$\text{aec}_{\varepsilon}(\mathcal{M})$	$\text{aec}_{\varepsilon}(\mathcal{M}) = \sup_{\bar{M} \in \text{co}(\mathcal{M})} \text{aec}_{\varepsilon}(\mathcal{M}, \bar{M})$
$\bar{\text{aec}}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}_0, \xi)$	AEC with randomized estimator ξ (see Eq. (20))
$\bar{\text{aec}}_{\varepsilon}(\mathcal{M}, \xi)$	$\bar{\text{aec}}_{\varepsilon}(\mathcal{M}, \xi) = \bar{\text{aec}}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}, \xi)$
$\bar{\text{aec}}_{\varepsilon}(\mathcal{M})$	$\bar{\text{aec}}_{\varepsilon}(\mathcal{M}) = \sup_{\xi \in \Delta_{\mathcal{M}}} \bar{\text{aec}}_{\varepsilon}(\mathcal{M}, \xi)$
$\text{aec}_{\varepsilon}^{\text{D}}(\mathcal{M}, \bar{M})$	AEC defined with respect to general divergence (see Eq. (38))
$\bar{\text{aec}}_{\varepsilon}^{\text{D}}(\mathcal{M}, \xi)$	AEC with randomized estimator, general divergence (see Eq. (39))
Uniform Exploration Notation	
$C_{\text{exp}}^{\xi}(\varepsilon)$	Uniform exploration coefficient with respect to ξ at scale ε (see Definition A.4)
$p_{\text{exp}}^{\xi}(\varepsilon)$	Uniform exploration distribution with respect to ξ at scale ε
$C_{\text{exp}}(\mathcal{M}, \varepsilon)$	Uniform exploration coefficient for class \mathcal{M} at scale ε
$C_{\text{exp}}^{\text{D}, \xi}(\varepsilon)$	Uniform exploration coefficient, general divergence (see Definition E.1)
$p_{\text{exp}}^{\text{D}, \xi}(\varepsilon)$	Uniform exploration distribution, general divergence
$C_{\text{exp}}^{\text{D}}(\mathcal{M}, \varepsilon)$	Uniform exploration coefficient for class \mathcal{M} , general divergence
Estimation Notation	
Alg_{KL}	Estimation oracle
$\text{Est}_{\text{KL}}(s)$	Cumulative KL estimation error (see Definition A.3)
$\text{Est}_{\text{D}}(s)$	Cumulative estimation error with respect to divergence D (see Eq. (40))
$\widehat{\text{Est}}_{\text{D}}(s)$	Cumulative estimation error with arguments flipped (see Eq. (41))

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

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

and (squared) Hellinger distance as

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

Interactive decision making. We formalize probability spaces in the same fashion as Foster et al. (2021, 2022b). Decisions are associated with a measurable space (Π, \mathcal{P}) , rewards are associated with the space $(\mathcal{R}, \mathcal{R})$, and observations are associated with the space $(\mathcal{O}, \mathcal{O})$. The history after round t is denoted by $\mathcal{H}^t = (\pi^1, r^1, o^1), \dots, (\pi^t, r^t, o^t)$. We define

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

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

When the algorithm is clear from context, we let \mathbb{P}^M denote the law it induces on \mathcal{H}^T when $M : \Pi \rightarrow \Delta_{\mathcal{R} \times \mathcal{O}}$ is the underlying model, and let $\mathbb{E}^M[\cdot]$ denote the corresponding expectation. We will also overload notation somewhat and let $\mathbb{P}^{M, \pi}$ the density of $(r, o) \sim M(\pi)$.

Notation for complexity measures and allocations. We let $T(\pi)$ denote the number of times decision π is taken up to time T . For $\eta \in \mathbb{R}_+^{\Pi}$, we define

$$I^M(\eta; \mathcal{M}) = \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi \in \Pi} \eta(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)), \quad (36)$$

so that we can write $\text{glc}(\mathcal{M}, M^*) = \inf_{\eta \in \mathbb{R}_+^{\Pi}} \{ \sum_{\pi \in \Pi} \eta(\pi) \Delta^{M^*}(\pi) \mid I^{M^*}(\eta; \mathcal{M}) \geq 1 \}$. We abbreviate $I^M(\eta) = I^M(\eta; \mathcal{M})$ whenever the class \mathcal{M} is clear from context. We will occasionally overload notation and write $\Delta^M(\eta) = \sum_{\pi \in \Pi} \eta(\pi) \Delta^M(\pi)$ for $\eta \in \mathbb{R}_+^{\Pi}$. We also let $\mathcal{M}_{\text{alt}}^*(M^*) := \mathcal{M}^{\text{alt}}(M^*)$ and

$$\mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda; \mathfrak{n}_{\text{max}}) := \{ M \in \mathcal{M} : \lambda \in \Lambda(M; \varepsilon, \mathfrak{n}_{\text{max}}) \}. \quad (37)$$

Recall the definition

$$\Lambda(M; \varepsilon) = \left\{ \lambda \in \Delta_{\Pi} \mid \exists \mathfrak{n} \in \mathbb{R}_+ \text{ s.t. } \mathbb{E}_{\pi \sim \lambda} [\Delta^M(\pi)] \leq \frac{(1 + \varepsilon) \mathfrak{g}^M}{\mathfrak{n}}, \right. \\ \left. \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \mathbb{E}_{\pi \sim \lambda} [D_{\text{KL}}(M(\pi) \| M'(\pi))] \geq \frac{1 - \varepsilon}{\mathfrak{n}} \right\}.$$

For a given $\lambda \in \Lambda(M; \varepsilon)$, we refer to the $\mathfrak{n} \in \mathbb{R}_+$ which realizes

$$\mathbb{E}_{\pi \sim \lambda} [\Delta^M(\pi)] \leq \frac{(1 + \varepsilon) \mathfrak{g}^M}{\mathfrak{n}} \quad \text{and} \quad \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \mathbb{E}_{\pi \sim \lambda} [D_{\text{KL}}(M(\pi) \| M'(\pi))] \geq \frac{1 - \varepsilon}{\mathfrak{n}}$$

as the *normalization factor* of λ .

General divergences. While in the main text we have focused on results that hold for the KL divergence, throughout the appendix we will consider more general divergences $D(\cdot \parallel \cdot)$. In particular, rather than fixing the divergence in Eq. (18) to the KL divergence, we utilize divergence D . We also perform online estimation with respect to D rather than with respect to the KL divergence. While D could be an arbitrary non-negative function, we make the following assumptions on it.

First, we replace [Assumption A.1](#) with the following more general assumption.

Assumption D.1. *For all $M, M', M'' \in \mathcal{M}$, and $\pi \in \Pi$, there exists some L_{KL} such that*

$$|D_{\text{KL}}(M(\pi) \parallel M''(\pi)) - D_{\text{KL}}(M'(\pi) \parallel M''(\pi))| \leq L_{\text{KL}} \sqrt{D(M(\pi) \parallel M'(\pi))}.$$

Note that, by Jensen's inequality, [Assumption D.1](#) immediately implies that, for $\xi \in \Delta(\mathcal{M})$,

$$\left| D_{\text{KL}}(M(\pi) \parallel M''(\pi)) - \mathbb{E}_{\bar{M} \sim \xi} [D_{\text{KL}}(\bar{M}(\pi) \parallel M''(\pi))] \right| \leq L_{\text{KL}} \sqrt{\mathbb{E}_{\bar{M} \sim \xi} [D(\bar{M}(\pi) \parallel M(\pi))]}.$$

We in addition make the following assumption, which we note is met for the KL divergence.

Assumption D.2. *For all $M, M' \in \text{co}(\mathcal{M})$, and π , we have*

$$D_{\text{H}}^2(M(\pi), M'(\pi)) \leq D(M(\pi) \parallel M'(\pi)).$$

Furthermore, $D(\cdot \parallel \cdot)$ is convex in its second argument.

A direct consequence of this assumption is that, when rewards are observed and bounded in $[0, 1]$, we have

$$|f^M(\pi) - f^{M'}(\pi)| \leq \sqrt{D(M(\pi) \parallel M'(\pi))}.$$

Throughout the appendix, we assume that [Assumption D.1](#) and [Assumption D.2](#) hold for our divergence D .

We also generalize the definition of the Allocation-Estimation Coefficient to account for general divergences as

$$\text{aec}_{\varepsilon}^{\text{D}}(\mathcal{M}, \bar{M}) := \inf_{\lambda, \omega \in \Delta_{\Pi}} \sup_{M \in \mathcal{M} \setminus \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda)} \frac{1}{\mathbb{E}_{\pi \sim \omega} [D(\bar{M}(\pi) \parallel M(\pi))]} \quad (38)$$

and

$$\overline{\text{aec}}_{\varepsilon}^{\text{D}}(\mathcal{M}; \xi) := \inf_{\lambda, \omega \in \Delta_{\Pi}} \sup_{M \in \mathcal{M} \setminus \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda)} \frac{1}{\mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{\pi \sim \omega} [D(\bar{M}(\pi) \parallel M(\pi))]]}. \quad (39)$$

Other variants of the AEC are generalized similarly with a superscript D. Our guarantees will also depend on a notion of estimation error for general divergences, given by

$$\mathbf{Est}_{\text{D}}(s) := \sum_{i=1}^s \mathbb{E}_{\widehat{M} \sim \xi^i} [\mathbb{E}_{\pi \sim p^i} [D(M^*(\pi) \parallel \widehat{M}(\pi))]]. \quad (40)$$

It will also be convenient to work with the following notion of estimation error:

$$\widehat{\mathbf{Est}}_{\text{D}}(s) := \sum_{i=1}^s \mathbb{E}_{\widehat{M} \sim \xi^i} [\mathbb{E}_{\pi \sim p^i} [D(\widehat{M}(\pi) \parallel M^*(\pi))]]. \quad (41)$$

Appendix E. Technical Tools

E.1. Online Learning

In this section, we state online estimation guarantees for variants of the Tempered Aggregation algorithm of [Chen et al. \(2022\)](#). Throughout the section we abbreviate $\mathbb{E}_{t-1}[\cdot] = \mathbb{E}[\cdot \mid \mathcal{H}^{t-1}, p^t, \xi^t]$. We recall that $\mathbb{P}^{M, \pi}(r, o)$ denotes the density over rewards and observations (r, o) under $M(\pi)$

Algorithm 4 Tempered Aggregation

- 1: **input:** Finite class \mathcal{M} .
- 2: Initialize $\xi^1 \leftarrow \text{Unif}(\mathcal{M})$.
- 3: **for** $t = 1, 2, 3, \dots$ **do**
- 4: Receive (π^t, r^t, o^t) .
- 5: Update estimator:

$$\xi^{t+1}(M) \propto \xi^t(M) \cdot \exp\left(\frac{1}{2} \log \mathbb{P}^{M, \pi^t}(r^t, o^t)\right) \quad \forall M \in \mathcal{M}.$$

Proposition E.1 (Tempered Aggregation, Finite Class Setting). *Assume $|\mathcal{M}| \leq \infty$ and $M^* \in \mathcal{M}$. Then [Algorithm 4](#) produces estimates $(\xi^t)_{t=1}^T$ which satisfy, with probability at least $1 - \delta$,*

$$\sum_{t=1}^T \mathbb{E}_{M \sim \xi^t} [\mathbb{E}_{\pi \sim p^t} [D_{\mathbb{H}}^2(M^*(\pi), M(\pi))]] \leq 2 \log \frac{|\mathcal{M}|}{\delta}.$$

Proof of Proposition E.1. We follow closely the proof of Theorem C.1 of [Chen et al. \(2022\)](#). Define the random variable

$$A^t := -\log \mathbb{E}_{M \sim \xi^t} \left[\exp\left(\beta \log \frac{\mathbb{P}^{M, \pi^t}(r^t, o^t)}{\mathbb{P}^{M^*, \pi^t}(r^t, o^t)}\right) \right].$$

We have

$$\begin{aligned} \mathbb{E}_{t-1}[\exp(-A^t)] &= \mathbb{E}_{t-1} \left[\mathbb{E}_{M \sim \xi^t} \left[\exp\left(\frac{1}{2} \log \frac{\mathbb{P}^{M, \pi^t}(r^t, o^t)}{\mathbb{P}^{M^*, \pi^t}(r^t, o^t)}\right) \right] \right] \\ &= \sum_{M \in \mathcal{M}} \xi^t(M) \mathbb{E}_{t-1} \left[\exp\left(\frac{1}{2} \log \frac{\mathbb{P}^{M, \pi^t}(r^t, o^t)}{\mathbb{P}^{M^*, \pi^t}(r^t, o^t)}\right) \right] \\ &= \sum_{M \in \mathcal{M}} \xi^t(M) \mathbb{E}_{t-1} \left[\mathbb{E}_{o \sim M^*(\pi^t)} \left[\sqrt{\frac{\mathbb{P}^{M, \pi^t}(r^t, o^t)}{\mathbb{P}^{M^*, \pi^t}(r^t, o^t)}} \right] \right] \\ &= \sum_{M \in \mathcal{M}} \xi^t(M) \cdot \left(1 - \frac{1}{2} \mathbb{E}_{\pi \sim p^t} [D_{\mathbb{H}}^2(M^*(\pi), M(\pi))] \right) \end{aligned}$$

where the last equality holds by the definition of the Hellinger distance. This implies that

$$1 - \mathbb{E}_{t-1}[\exp(-A^t)] = \frac{1}{2} \mathbb{E}_{M \sim \xi^t} [\mathbb{E}_{\pi \sim p^t} [D_{\mathbb{H}}^2(M^*(\pi), M(\pi))]].$$

By Lemma A.4 of [Foster et al. \(2021\)](#), we have that with probability at least $1 - \delta$,

$$\begin{aligned} \sum_{t=1}^T A^t + \log \frac{1}{\delta} &\geq \sum_{t=1}^T -\log \mathbb{E}_{t-1}[\exp(-A^t)] \\ &\geq \sum_{t=1}^T (1 - \mathbb{E}_{t-1}[\exp(-A^t)]) \\ &= \frac{1}{2} \sum_{t=1}^T \mathbb{E}_{M \sim \xi^t} [\mathbb{E}_{\pi \sim p^t} [D_{\text{H}}^2(M^*(\pi), M(\pi))]]. \end{aligned}$$

We turn to upper bounding $\sum_{t=1}^T A^t$. Following the proof of Theorem C.1 of [Chen et al. \(2022\)](#), we have

$$\sum_{t=1}^T A^t = -\log \left(\sum_{M \in \mathcal{M}} \xi^1(M) \exp \left(\sum_{t=1}^T \frac{1}{2} \log \frac{\mathbb{P}^{M, \pi^t}(r^t, o^t)}{\mathbb{P}^{M^*, \pi^t}(r^t, o^t)} \right) \right).$$

Since $M^* \in \mathcal{M}$, we can then bound

$$\sum_{t=1}^T A^t \leq -\log \left(\xi^1(M^*) \exp \left(\sum_{t=1}^T \frac{1}{2} \log \frac{\mathbb{P}^{M^*, \pi^t}(r^t, o^t)}{\mathbb{P}^{M^*, \pi^t}(r^t, o^t)} \right) \right) = \log |\mathcal{M}|.$$

Combining these expressions gives the result. \square

Proposition E.2 (Tempered Aggregation, Infinite Class Setting). *Let \mathcal{M}_{cov} denote a (ρ, μ) -cover of \mathcal{M} with covering number $N_{\text{cov}}(\mathcal{M}, \rho, \mu)$. If we apply [Algorithm 4](#) to \mathcal{M}_{cov} , we have that whenever $M^* \in \mathcal{M}$, with probability at least $1 - \delta - T\mu$,*

$$\sum_{t=1}^T \mathbb{E}_{M \sim \xi^t} [\mathbb{E}_{\pi \sim p^t} [D_{\text{H}}^2(M^*(\pi), M(\pi))]] \leq 2 \log \frac{N_{\text{cov}}(\mathcal{M}, \rho, \mu)}{\delta} + T\rho.$$

Proof of Proposition E.2. Defining A^t as in [Proposition E.1](#), the first part of the proof is identical to that of [Proposition E.1](#), and we conclude that, with probability at least $1 - \delta$,

$$\sum_{t=1}^T A^t + \log \frac{1}{\delta} \geq \frac{1}{2} \sum_{t=1}^T \mathbb{E}_{M \sim \xi^t} [\mathbb{E}_{\pi \sim p^t} [D_{\text{H}}^2(M^*(\pi), M(\pi))]]$$

and

$$\sum_{t=1}^T A^t = -\log \left(\sum_{M \in \mathcal{M}_{\text{cov}}} \xi^1(M) \exp \left(\sum_{t=1}^T \frac{1}{2} \log \frac{\mathbb{P}^{M, \pi^t}(r^t, o^t)}{\mathbb{P}^{M^*, \pi^t}(r^t, o^t)} \right) \right).$$

Let \mathcal{E} denote the event associated with \mathcal{M}_{cov} , as defined in [Definition A.1](#), which satisfies $\sup_{M' \in \mathcal{M}} \sup_{\pi} \mathbb{P}^{M'}(\mathcal{E} | \pi) \leq \mu$. Let $\tilde{M} \in \mathcal{M}_{\text{cov}}$ denote the element in the cover which satisfies

$$\left| \log \frac{\mathbb{P}^{M^*, \pi}(r, o)}{\mathbb{P}^{\tilde{M}, \pi}(r, o)} \right| = \left| \log \mathbb{P}^{M^*, \pi}(r, o) - \log \mathbb{P}^{\tilde{M}, \pi}(r, o) \right| \leq \rho,$$

for all (r, o, π) with $\sup_{M' \in \mathcal{M}} \mathbb{P}^{M', \pi}(r, o \mid \mathcal{E}) > 0$. We then have

$$\sum_{t=1}^T \log \frac{\mathbb{P}^{M, \pi^t}(r^t, o^t)}{\mathbb{P}^{M^*, \pi^t}(r^t, o^t)} = \sum_{t=1}^T \left(\log \frac{\mathbb{P}^{M, \pi^t}(r^t, o^t)}{\mathbb{P}^{\bar{M}, \pi^t}(r^t, o^t)} + \log \frac{\mathbb{P}^{\bar{M}, \pi^t}(r^t, o^t)}{\mathbb{P}^{M^*, \pi^t}(r^t, o^t)} \right)$$

so we can bound

$$\sum_{t=1}^T A^t \leq \log |\mathcal{M}_{\text{cov}}| + \frac{1}{2} \sum_{t=1}^T \log \frac{\mathbb{P}^{M^*, \pi^t}(r^t, o^t)}{\mathbb{P}^{\bar{M}, \pi^t}(r^t, o^t)}$$

which gives that with probability at least $1 - \delta$,

$$\sum_{t=1}^T \mathbb{E}_{M \sim \xi^t} [\mathbb{E}_{\pi \sim p^t} [D_{\text{H}}^2(M^*(\pi), M(\pi))]] \leq 2 \log |\mathcal{M}_{\text{cov}}| + 2 \log \frac{1}{\delta} + \sum_{t=1}^T \log \frac{\mathbb{P}^{M^*, \pi^t}(r^t, o^t)}{\mathbb{P}^{\bar{M}, \pi^t}(r^t, o^t)}.$$

Let \mathcal{E}_t denote the event that \mathcal{E} occurs at step t . Denote the event $\mathcal{E}_1 := \cap_{t=1}^T \mathcal{E}_t$ and

$$\mathcal{E}_2 := \left\{ \sum_{t=1}^T \mathbb{E}_{M \sim \xi^t} [\mathbb{E}_{\pi \sim p^t} [D_{\text{H}}^2(M^*(\pi), M(\pi))]] \leq 2 \log |\mathcal{M}_{\text{cov}}| + 2 \log \frac{1}{\delta} + T\rho \right\}.$$

Then

$$\mathbb{P}^{M^*}[\mathcal{E}_2] \leq \mathbb{P}^{M^*}[\mathcal{E}_2 \cap \mathcal{E}_1] + \mathbb{P}^{M^*}[\mathcal{E}_1^c].$$

By definition of \mathcal{E}_t and a union bound we have $\mathbb{P}^{M^*}[\mathcal{E}_1^c] \leq T\mu$. Furthermore, on the event \mathcal{E}_1 we can bound, for each $t \leq T$,

$$\log \frac{\mathbb{P}^{M^*, \pi^t}(r^t, o^t)}{\mathbb{P}^{\bar{M}, \pi^t}(r^t, o^t)} \leq \rho.$$

Thus, it follows that on \mathcal{E}_1 , $\sum_{t=1}^T \log \frac{\mathbb{P}^{M^*, \pi^t}(r^t, o^t)}{\mathbb{P}^{\bar{M}, \pi^t}(r^t, o^t)} \leq T\rho$, which implies that $\mathbb{P}^{M^*}[\mathcal{E}_2 \cap \mathcal{E}_1] \leq \delta$. The result follows. \square

E.2. Properties of Graves-Lai Program

In this section, we establish some basic properties of the Graves-Lai coefficient $\mathbf{g}^M \equiv \text{glc}(\mathcal{M}, M)$. Through out the section, we omit dependence on the class \mathcal{M} for various quantities of interest whenever it is clear from context. Throughout, we will use the fact that whenever $\Delta_{\min}^M > 0$, π_M is unique.

E.2.1. BASIC PROPERTIES OF GRAVES-LAI PROGRAM

Lemma E.1. *For any $M \in \mathcal{M}$ and $n > 0$, we have*

$$\frac{\mathbf{g}^M}{n} \leq \inf_{\lambda \in \Delta_{\Pi}} \left\{ \Delta^M(\lambda) : I^M(\lambda) \geq \frac{1}{n} \right\}.$$

Proof of Lemma E.1. Assume the contrary. Then there exists some $\tilde{\lambda}$ such that

$$\Delta^M(\tilde{\lambda}) < \mathbf{g}^M/n \quad \text{and} \quad I^M(\tilde{\lambda}) \geq 1/n.$$

However, since both $\Delta^M(\lambda)$ and $I^M(\lambda)$ are linear in rescaling of λ , this implies that

$$\Delta^M(n\tilde{\lambda}) < \mathbf{g}^M \quad \text{and} \quad I^M(n\tilde{\lambda}) \geq 1.$$

By definition we have

$$\mathbf{g}^M = \inf_{\eta \in \mathbb{R}_+^\Pi} \Delta^M(\eta) \quad \text{s.t.} \quad I^M(\eta) \geq 1.$$

This is a contradiction, so the desired conclusion follows. \square

Lemma E.2. For any $M \in \mathcal{M}^+$ with $\Delta_{\min}^M > 0$, we can bound

$$\mathbf{g}^M \leq C_{\text{exp}}^\xi \left(\frac{1}{4} (\Delta_{\min}^M)^2 \right)$$

for $\xi = \mathbb{I}_M$.

Proof of Lemma E.2. By definition we have

$$\begin{aligned} \mathbf{g}^M &= \inf_{\eta \in \mathbb{R}_+^\Pi} \Delta^M(\eta) \quad \text{s.t.} \quad \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi} \eta(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) \geq 1 \\ &\leq \inf_{\eta \in \mathbb{R}_+^\Pi} \|\eta\|_1 \quad \text{s.t.} \quad \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi} \eta(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) \geq 1. \end{aligned}$$

Let $\xi \in \Delta_{\mathcal{M}}$ denote the distribution with $\xi(M) = 1$. Let and let $p_{\text{exp}} := p_{\text{exp}}^\xi(\varepsilon)$ denote the uniform exploration distribution defined with respect to ξ , and $C_{\text{exp}}^\xi(\varepsilon)$ the corresponding uniform exploration coefficient, for some ε to be chosen.

Consider some $M' \in \mathcal{M}^{\text{alt}}(M)$. Since $\mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{\pi \sim p} [D_{\text{KL}}(\bar{M}(\pi) \| M''(\pi))]] = \mathbb{E}_{\pi \sim p} [D_{\text{KL}}(M(\pi) \| M''(\pi))]$ for all M'' and $D_{\text{KL}}(M(\pi) \| M(\pi)) = 0$ for all π , it follows from the definition of the uniform exploration coefficient that

$$\mathbb{E}_{p_{\text{exp}}} [D_{\text{KL}}(M(\pi) \| M'(\pi))] \leq 1/C_{\text{exp}}^\xi(\varepsilon) \implies \sup_{p \in \Delta_\Pi} \mathbb{E}_p [D_{\text{KL}}(M(\pi) \| M'(\pi))] \leq \varepsilon$$

or, alternatively,

$$\sup_{p \in \Delta_\Pi} \mathbb{E}_p [D_{\text{KL}}(M(\pi) \| M'(\pi))] > \varepsilon \implies \mathbb{E}_{p_{\text{exp}}} [D_{\text{KL}}(M(\pi) \| M'(\pi))] > 1/C_{\text{exp}}^\xi(\varepsilon).$$

If $M' \in \mathcal{M}^{\text{alt}}(M)$, then it follows that $\pi_M \notin \pi_{M'}$. Take some $\pi_{M'} \in \pi_{M'}$. By definition we have $f^M(\pi_M) \geq f^M(\pi_{M'}) + \Delta_{\min}^M$ and $f^{M'}(\pi_{M'}) \geq f^{M'}(\pi_M)$. Thus,

$$\begin{aligned} \Delta_{\min}^M &\leq f^M(\pi_M) - f^M(\pi_{M'}) + f^{M'}(\pi_{M'}) - f^{M'}(\pi_M) \\ &\leq |f^M(\pi_M) - f^{M'}(\pi_M)| + |f^{M'}(\pi_{M'}) - f^M(\pi_{M'})| \\ &\leq \sqrt{D(M(\pi_M) \| M'(\pi_M))} + \sqrt{D(M(\pi_{M'}) \| M'(\pi_{M'}))}. \end{aligned}$$

This implies that there exists some π such that $D(M(\pi) \| M'(\pi)) \geq (\Delta_{\min}^M/2)^2$, so we can lower bound

$$\sup_{p \in \Delta_{\Pi}} \mathbb{E}_p[D_{\text{KL}}(M(\pi) \| M'(\pi))] \geq \frac{1}{4}(\Delta_{\min}^M)^2.$$

Thus, setting $\varepsilon = \frac{1}{4}(\Delta_{\min}^M)^2$, we have that

$$\mathbb{E}_{p_{\text{exp}}}[D_{\text{KL}}(M(\pi) \| M'(\pi))] > 1/C_{\text{exp}}^{\xi}(\frac{1}{4}(\Delta_{\min}^M)^2).$$

It follows that the allocation $\eta = C_{\text{exp}}^{\xi}(\frac{1}{4}(\Delta_{\min}^M)^2) \cdot p_{\text{exp}}$ realizes

$$\inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi} \eta(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) = \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} C_{\text{exp}}^{\xi}(\frac{1}{4}(\Delta_{\min}^M)^2) \cdot \mathbb{E}_{p_{\text{exp}}}^M[D_{\text{KL}}(M(\pi) \| M'(\pi))] \geq 1,$$

which proves the result. \square

Lemma E.3. *Assume $\mathbf{g}^M > 0$, $\Delta_{\min}^M > 0$, and $\mathfrak{n}_{1/4}^M < \infty$. Then it must be the case that*

$$\mathbf{g}^M \geq \Delta_{\min}^M \cdot \frac{1}{\max_{M' \in \mathcal{M}, \pi \in \Pi} D_{\text{KL}}(M(\pi) \| M'(\pi))}.$$

Proof of Lemma E.3. Recall that

$$\mathbf{g}^M = \inf_{\eta \in \mathbb{R}_+^{\Pi}} \Delta^M(\eta) \quad \text{s.t.} \quad \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi} \eta(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) \geq 1.$$

By the definition of $\mathfrak{n}_{1/4}^M$, for any allocation $\eta \in \mathbb{R}_+^{\Pi}$ satisfying $\inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi} \eta(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) \geq 3/4$, the allocation $\tilde{\eta}$ defined as $\tilde{\eta}(\pi) = \eta(\pi)$ for $\pi \neq \pi_M$, and $\tilde{\eta}(\pi_M) = \mathfrak{n}_{1/2}^M$ satisfies

$$\inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi} \tilde{\eta}(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) \geq 1/2,$$

and furthermore $\Delta^M(\eta) = \Delta^M(\tilde{\eta})$. It follows that

$$\begin{aligned} \mathbf{g}^M &\geq \inf_{\eta \in \mathbb{R}_+^{\Pi}} \Delta^M(\eta) \quad \text{s.t.} \quad \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi} \eta(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) \geq 3/4 \\ &\geq \inf_{\eta \in \mathbb{R}_+^{\Pi}} \Delta^M(\eta) \quad \text{s.t.} \quad \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi} \eta(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) \geq 1/2, \quad \eta(\pi_M) \leq \mathfrak{n}_{1/4}^M. \end{aligned}$$

If for all $M' \in \mathcal{M}^{\text{alt}}(M)$ we have $D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) > 0$, this implies that we can distinguish M from M' by playing only π_M . Furthermore, by what we have just shown, this can be achieved by playing π_M at most $2\mathfrak{n}_{1/4}^M$ times. It follows that, if $D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) > 0$ for all $M' \in \mathcal{M}^{\text{alt}}(M)$, then $\mathbf{g}^M = 0$. Thus, if $\mathbf{g}^M > 0$, there must exist some $M' \in \mathcal{M}^{\text{alt}}(M)$ such that $D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) = 0$.

We then have

$$\mathbf{g}^M \geq \inf_{\eta \in \mathbb{R}_+^{\Pi}} \Delta^M(\eta) \quad \text{s.t.} \quad \sum_{\pi \neq \pi_M} \eta(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) \geq 1$$

$$\begin{aligned}
 &= \inf_{\eta \in \mathbb{R}_+^\Pi, \eta(\pi_M)=0} \Delta^M(\eta) \quad \text{s.t.} \quad \sum_{\pi \neq \pi_M} \eta(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) \geq 1 \\
 &\geq \inf_{\eta \in \mathbb{R}_+^\Pi, \eta(\pi_M)=0} \Delta_{\min}^M \cdot \|\eta\|_1 \quad \text{s.t.} \quad \max_{\pi} D_{\text{KL}}(M(\pi) \| M'(\pi)) \cdot \|\eta\|_1 \geq 1 \\
 &= \Delta_{\min}^M \cdot \frac{1}{\max_{\pi} D_{\text{KL}}(M(\pi) \| M'(\pi))}.
 \end{aligned}$$

The result follows. \square

E.2.2. PROPERTIES OF THE INFORMATION CONTENT OF OPTIMAL DECISIONS

Lemma E.4. *Let $\varepsilon \in [0, 1/2)$ and $\bar{n} > 0$ be given. We can bound, for any function $g(\omega, M)$ and $\mathcal{M}_0 \subseteq \mathcal{M}$ with $\inf_{M \in \mathcal{M}_0} \Delta_{\min}^M > 0$,*

$$\inf_{\omega, \lambda \in \Delta_\Pi} \sup_{M \in \mathcal{M}_0 \setminus \mathcal{M}_{2\varepsilon}^{\text{gl}}(\lambda; \bar{n})} g(\omega, M) \leq \inf_{\omega, \lambda \in \Delta_\Pi} \sup_{M \in \mathcal{M}_0 \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)} g(\omega, M)$$

as long as

$$\bar{n} \geq \max_{M \in \mathcal{M}_0} \max \left\{ n_\varepsilon^M, \frac{4\mathbf{g}^M}{\Delta_{\min}^M}, \frac{2\mathbf{g}^M}{\zeta \Delta_{\min}^M} \right\},$$

where

$$\zeta := \min_{M \in \mathcal{M}_0: \mathbf{g}^M > 0} \min \left\{ \frac{\mathbf{g}^M}{\mathbf{g}^M + 2n_\varepsilon^M}, \frac{\Delta_{\min}^M \varepsilon}{4} \right\}.$$

Proof of Lemma E.4. Note that each of these expressions only depend on λ through $\mathcal{M}_{2\varepsilon}^{\text{gl}}(\lambda; \bar{n})$ and $\mathcal{M}_\varepsilon^{\text{gl}}(\lambda)$, respectively. To prove the result, it therefore suffices to show that, for every $\lambda \in \Delta_\Pi$, there exists $\lambda' \in \Delta_\Pi$ such that $\mathcal{M}_0 \cap \mathcal{M}_\varepsilon^{\text{gl}}(\lambda) \subseteq \mathcal{M}_0 \cap \mathcal{M}_{2\varepsilon}^{\text{gl}}(\lambda'; \bar{n})$.

Fix $\lambda \in \Delta_\Pi$. Consider $M \in \mathcal{M}_0 \cap \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)$. By definition we have that there exists some $n > 0$ such that

$$\Delta^M(\lambda) \leq (1 + \varepsilon)\mathbf{g}^M/n \quad \text{and} \quad I^M(\lambda) \geq (1 - \varepsilon)/n,$$

where here $I^M(\lambda) = I^M(\lambda; \mathcal{M})$. We consider two cases.

Case 1: $\lambda = \mathbb{I}_\pi$ for some $\pi \in \Pi$. First, suppose that $\pi \neq \pi_M$. Then we have

$$\Delta_{\min}^M \leq \Delta^M(\lambda) \leq (1 + \varepsilon)\mathbf{g}^M/n \implies n \leq (1 + \varepsilon)\mathbf{g}^M/\Delta_{\min}^M.$$

It follows that as long as $\bar{n} \geq (1 + \varepsilon)\mathbf{g}^M/\Delta_{\min}^M$, then $M \in \mathcal{M}_{2\varepsilon}^{\text{gl}}(\lambda; \bar{n})$.

Now, suppose that $\pi = \pi_M$. In this case, by the definition of n_ε^M , we immediately have that it suffices to take $n = n_\varepsilon^M$. Thus, for $\lambda = \mathbb{I}_\pi$, we have $\mathcal{M}_{2\varepsilon}^{\text{gl}}(\lambda) = \mathcal{M}_\varepsilon^{\text{gl}}(\lambda; \bar{n})$ as long as

$$\bar{n} \geq \max \{ n_\varepsilon^M, (1 + \varepsilon)\mathbf{g}^M/\Delta_{\min}^M \}.$$

Case 2a: $\lambda \neq \mathbb{I}_\pi$ for any $\pi \in \Pi$. Fix some $\zeta \in (0, 1/2)$ to be chosen. Suppose that there exists some π' such that $\lambda(\pi') \geq 1 - \zeta$, and note that there can exist at most one such π' . Define λ' as

$$\lambda'(\pi') = 1 - \zeta, \quad \lambda'(\pi) = \frac{\zeta}{1 - \lambda(\pi')} \lambda(\pi), \quad \forall \pi \neq \pi'$$

and note that $\lambda' \in \Delta_\Pi$. Our goal is to show that $\mathcal{M}_\varepsilon^{\text{gl}}(\lambda) \subseteq \mathcal{M}_{2\varepsilon}^{\text{gl}}(\lambda'; \bar{\mathbf{n}})$.

Suppose that $\pi_M = \pi'$. Then

$$\Delta^M(\lambda') = \frac{\zeta}{1 - \lambda(\pi')} \cdot \Delta^M(\lambda) \leq \frac{\zeta}{1 - \lambda(\pi')} \cdot \frac{(1 + \varepsilon)\mathbf{g}^M}{\mathbf{n}}.$$

Denote $\mathbf{n}' := (\frac{\zeta}{1 - \lambda(\pi')} \cdot \frac{1}{\mathbf{n}})^{-1}$. Since $I^M(\lambda) \geq (1 - \varepsilon)/\mathbf{n}$, we have $I^M(\mathbf{n}\lambda) \geq 1 - \varepsilon$. Then, by the definition of \mathbf{n}_ε^M , we have

$$\inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi \neq \pi_M} \mathbf{n}\lambda(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) + \mathbf{n}_\varepsilon^M D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) \geq 1 - 2\varepsilon.$$

However, note that

$$\begin{aligned} I^M(\mathbf{n}'\lambda') &= \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi \neq \pi_M} \frac{\zeta \mathbf{n}'}{1 - \lambda(\pi_M)} \lambda(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) + \mathbf{n}'(1 - \zeta) D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) \\ &= \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi \neq \pi_M} \mathbf{n}\lambda(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) + \frac{(1 - \zeta)(1 - \lambda(\pi'))\mathbf{n}}{\zeta} \cdot D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) \\ &\stackrel{(a)}{\geq} \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi \neq \pi_M} \mathbf{n}\lambda(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) + \mathbf{n}_\varepsilon^M D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) \\ &\stackrel{(b)}{\geq} 1 - 2\varepsilon \end{aligned}$$

where (a) follows as long as

$$\frac{(1 - \zeta)(1 - \lambda(\pi'))\mathbf{n}}{\zeta} \geq \mathbf{n}_\varepsilon^M, \quad (42)$$

and (b) follows from what we have just shown. Rearranging, we have that [Eq. \(42\)](#) is equivalent to

$$\frac{(1 - \lambda(\pi'))\mathbf{n}}{(1 - \lambda(\pi'))\mathbf{n} + \mathbf{n}_\varepsilon^M} \geq \zeta.$$

Note that by [Lemma E.1](#) and the definition of λ , we have

$$\frac{(1 - \varepsilon)\mathbf{g}^M}{\mathbf{n}} \leq \inf_{\tilde{\lambda} \in \Delta_\Pi} \left\{ \Delta^M(\tilde{\lambda}) : I^M(\tilde{\lambda}) \geq \frac{1 - \varepsilon}{\mathbf{n}} \right\} \leq \Delta^M(\lambda),$$

which implies that

$$(1 - \varepsilon)\mathbf{g}^M \leq \Delta^M(\lambda) \cdot \mathbf{n} \leq (1 - \lambda(\pi_M)) \cdot \mathbf{n}.$$

As the function $\frac{x}{x+n_\varepsilon^M}$ is increasing in x , a sufficient choice of ζ for this M is then

$$\min\left\{\frac{\mathbf{g}^M}{\mathbf{g}^M + 2n_\varepsilon^M}, 3/8\right\} \geq \zeta$$

Thus, for such a ζ , we have that $M \in \mathcal{M}_\varepsilon^{\text{gl}}(\lambda'; n')$, which implies that $M \in \mathcal{M}_{2\varepsilon}^{\text{gl}}(\lambda'; n')$. Note that $(1 - \lambda(\pi'))n \leq (1 + \varepsilon)\mathbf{g}^M/\Delta_{\min}^M$ in this case, so we have that $M \in \mathcal{M}_{2\varepsilon}^{\text{gl}}(\lambda'; \bar{n})$ as long as

$$\bar{n} \geq \frac{2\mathbf{g}^M}{\zeta\Delta_{\min}^M}.$$

Consider now the case where $\pi_M \neq \pi'$. In this case, defining λ' as before, we can bound

$$\Delta^M(\lambda') \leq \zeta + (1 - \zeta)\Delta^M(\pi') \leq \zeta + \lambda(\pi')\Delta^M(\pi') \leq \zeta + \Delta^M(\lambda) \leq \zeta + (1 + \varepsilon)\mathbf{g}^M/n.$$

Furthermore,

$$\begin{aligned} I^M(\lambda') &= \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi \neq \pi'} \frac{\zeta}{1 - \lambda(\pi')} \lambda(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) + (1 - \zeta) D_{\text{KL}}(M(\pi') \| M'(\pi')) \\ &\geq \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi \neq \pi'} (1 - \zeta) \lambda(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) + (1 - \zeta) \lambda(\pi') D_{\text{KL}}(M(\pi') \| M'(\pi')) \\ &= (1 - \zeta) I^M(\lambda) \\ &\geq (1 - \zeta)(1 - \varepsilon)/n. \end{aligned}$$

Since $\pi' \neq \pi_M$, we can lower bound $\Delta^M(\lambda) \geq (1 - \zeta)\Delta_{\min}^M$, so

$$(1 - \zeta)\Delta_{\min}^M \leq \Delta^M(\lambda) \leq (1 + \varepsilon)\mathbf{g}^M/n \implies n \leq \frac{(1 + \varepsilon)\mathbf{g}^M}{(1 - \zeta)\Delta_{\min}^M} \leq \frac{4\mathbf{g}^M}{\Delta_{\min}^M}$$

We can therefore bound $\zeta \leq \varepsilon\mathbf{g}^M/n$ as long as

$$\zeta \leq \frac{\Delta_{\min}^M \cdot \varepsilon}{4}.$$

Consider ζ that satisfies this inequality. Then $\Delta^M(\lambda') \leq (1 + 2\varepsilon)\mathbf{g}^M/n$. We can also lower bound $(1 - \zeta)(1 - \varepsilon) \geq (1 - 2\varepsilon)$ as long as $\zeta \leq \frac{\varepsilon}{1 - \varepsilon} \leq \varepsilon$. Thus, as long as

$$\zeta \leq \min\left\{1, \frac{\Delta_{\min}^M}{4}\right\} \cdot \varepsilon \quad \text{and} \quad \bar{n} \geq \frac{4\mathbf{g}^M}{\Delta_{\min}^M},$$

we have that $M \in \mathcal{M}_{2\varepsilon}^{\text{gl}}(\lambda'; \bar{n})$.

Case 2b: $\lambda \neq \mathbb{I}_\pi$ for any $\pi \in \Pi$. Finally, it remains to handle the case then there does not exist π' such that $\lambda(\pi') \geq 1 - \zeta$. In this case, we can always lower bound

$$\zeta\Delta_{\min}^M \leq \Delta^M(\lambda) \leq (1 + \varepsilon)\mathbf{g}^M/n \implies n \leq \frac{(1 + \varepsilon)\mathbf{g}^M}{\zeta\Delta_{\min}^M},$$

so as long as $\bar{n} \geq \frac{(1 + \varepsilon)\mathbf{g}^M}{\zeta\Delta_{\min}^M}$, we have $M \in \mathcal{M}_{2\varepsilon}^{\text{gl}}(\lambda; \bar{n})$.

Since we are in the regime where $\lambda \neq \mathbb{I}_\pi$, it must be the case that if $M \in \mathcal{M}_0 \cap \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)$, then $\mathbf{g}^M > 0$. Thus, a sufficient choice of ζ is

$$\zeta = \min_{M \in \mathcal{M}_0: \mathbf{g}^M > 0} \min\left\{\frac{\mathbf{g}^M}{\mathbf{g}^M + 2n_\varepsilon^M}, \frac{\Delta_{\min}^M \varepsilon}{4}\right\}.$$

Concluding the Proof. To show the result, we need that \bar{n} is large enough for each $M \in \mathcal{M}_0$. The argument above shows that it suffices to take

$$\bar{n} \geq \max_{M \in \mathcal{M}_0} \max \left\{ n_\varepsilon^M, \frac{4g^M}{\Delta_{\min}^M}, \frac{2g^M}{\zeta \Delta_{\min}^M} \right\}$$

for

$$\zeta := \min_{M \in \mathcal{M}_0: g^M > 0} \min \left\{ \frac{g^M}{g^M + 2n_\varepsilon^M}, \frac{\Delta_{\min}^M \varepsilon}{4} \right\}.$$

This proves the result. □

Lemma E.5. *For every $M \in \mathcal{M}$ with $\Delta_{\min}^M > 0$, there exists some $\lambda \in \Lambda(M; \varepsilon)$ with normalization factor n satisfying*

$$n \leq g^M / \Delta_{\min}^M + n_\varepsilon^M,$$

i.e., we have $M \in \mathcal{M}_\varepsilon^{\text{gl}}(\lambda; n)$.

Proof of Lemma E.5. Consider some allocation $\eta \in \mathbb{R}_+^\Pi$ such that

$$\Delta^M(\eta) \leq g^M \quad \text{and} \quad I^M(\eta) \geq 1. \quad (43)$$

Let η' denote the allocation satisfying $\eta'(\pi) = \eta(\pi)$ for $\pi \neq \pi_M$, and $\eta'(\pi_M) = 0$. Note that $\Delta^M(\eta) \geq \Delta_{\min}^M \|\eta'\|_1$, which implies that

$$\|\eta'\|_1 \leq g^M / \Delta_{\min}^M.$$

Let η'' denote the allocation satisfying $\eta''(\pi) = \eta(\pi)$ for $\pi \neq \pi_M$ and $\eta''(\pi_M) = n_\varepsilon^M$. Then by definition of n_ε^M , and since η satisfies Eq. (43), we have $I^M(\eta'') \geq 1 - \varepsilon$. Furthermore, it is straightforward to see that $\Delta^M(\eta'') \leq g^M$. This implies that $\eta'' / \|\eta''\|_1 \in \Lambda(M; \varepsilon)$ with normalization factor $\|\eta''\|_1$. However, we can bound

$$\|\eta''\|_1 = \|\eta'\|_1 + n_\varepsilon^M \leq g^M / \Delta_{\min}^M + n_\varepsilon^M.$$

This proves the result. □

E.2.3. BOUNDING ALLOCATION-ESTIMATION COEFFICIENT VIA UNIFORM EXPLORATION COEFFICIENT

In this section we prove a generalized version of [Proposition A.1](#). In particular, rather than specializing to the KL divergence, we consider a general divergence D . We define the uniform exploration coefficient with respect to D as follows.

Definition E.1 (Uniform Exploration Coefficient, General Divergences). *For a randomized estimator $\xi \in \Delta_{\mathcal{M}}$ and divergence $D(\cdot \| \cdot)$, we define the uniform exploration coefficient with respect to ξ at scale $\varepsilon > 0$ as the value of the following program:*

$$C_{\text{exp}}^{\text{D},\xi}(\varepsilon) := \min_{C \in \mathbb{R}_+, p \in \Delta_{\Pi}} \left\{ C \mid \forall M, M' \in \mathcal{M} : \begin{array}{l} \max_{M'' \in \{M, M'\}} \mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{\pi \sim p} [D(\bar{M}(\pi) \| M''(\pi))]] \leq 1/C \\ \implies \max_{p' \in \Delta_{\Pi}} \mathbb{E}_{\pi \sim p'} [D(M(\pi) \| M'(\pi))] \leq \varepsilon \end{array} \right\}.$$

We define $p_{\text{exp}}^{\text{D},\xi}(\varepsilon)$ as the minimizing distribution for this program, and let

$$C_{\text{exp}}^{\text{D}}(\mathcal{M}, \varepsilon) := \sup_{\xi \in \Delta_{\mathcal{M}}} C_{\text{exp}}^{\text{D},\xi}(\varepsilon)$$

denote the uniform exploration constant for class \mathcal{M} .

Lemma E.6 (Formal version of [Proposition A.1](#)). *Let $\varepsilon \in [0, 1/2)$ and $\mathcal{M}_0 \subseteq \mathcal{M}$ be given, and assume that $\inf_{M \in \mathcal{M}_0} \mathbf{g}^M > 0$, $\inf_{M \in \mathcal{M}_0} \Delta_{\min}^M > 0$, $\sup_{M \in \mathcal{M}_0} \mathbf{n}_{1/4}^M < \infty$, and [Assumptions A.2, D.1](#) and [D.2](#) hold. Then for any $\xi \in \Delta_{\mathcal{M}_0}$, we can bound*

$$\min_{\lambda, \omega \in \Delta_{\Pi}} \sup_{M \in \mathcal{M}_0 \setminus \mathcal{M}_{\xi}^{\text{gl}}(\lambda)} \frac{1}{\mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{\omega} [D(\bar{M}(\pi) \| M(\pi))]]} \leq C_{\text{exp}}^{\text{D}}(\mathcal{M}_0, \delta)$$

for any $\delta > 0$ satisfying

$$\sqrt{\delta} \leq \min_{M \in \mathcal{M}_0} \min \left\{ \min \left\{ \frac{1}{81L_{\text{KL}}}, \frac{\Delta_{\min}^M}{34V_{\mathcal{M}}} \right\} \cdot \frac{\varepsilon}{2\mathbf{g}^M / \Delta_{\min}^M + \mathbf{n}_{\varepsilon/36}^M}, \frac{\Delta_{\min}^M}{3} \right\}.$$

Proof of Lemma E.6. Let $\delta > 0$ be some tolerance to be chosen. Let \widetilde{M} denote some $M \in \mathcal{M}_0$ such that $\mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{p_{\text{exp}}} [D(\bar{M}(\pi) \| M(\pi))]] \leq 1/C_{\text{exp}}$, where we abbreviate $C_{\text{exp}} := C_{\text{exp}}^{\text{D}}(\mathcal{M}_0, \delta)$ and $p_{\text{exp}} := p_{\text{exp}}^{\text{D},\xi}(\delta)$ is a distribution that achieves the value of $C_{\text{exp}}^{\text{D}}(\mathcal{M}_0, \delta)$ for ξ ; if such an \widetilde{M} does not exist, we let $\widetilde{M} = \arg \min_{M \in \mathcal{M}_0} \mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{p_{\text{exp}}} [D(\bar{M}(\pi) \| M(\pi))]]$. Let $\varepsilon' > 0$ be some value to be chosen, and let $\tilde{\lambda} \in \Lambda(\widetilde{M}; \varepsilon')$ denote the allocation in $\Lambda(\widetilde{M}; \varepsilon')$ with smallest normalizing factor \mathbf{n} . Let $\tilde{\mathbf{n}}$ denote the value of this normalizing factor, then:

$$\Delta^{\widetilde{M}}(\tilde{\lambda}) \leq (1 + \varepsilon') \mathbf{g}^{\widetilde{M}} / \tilde{\mathbf{n}} \quad \text{and} \quad I^{\widetilde{M}}(\tilde{\lambda}) \geq (1 - \varepsilon') / \tilde{\mathbf{n}}.$$

We can bound:

$$\begin{aligned} & \min_{\lambda \in \Delta_{\Pi}} \min_{\omega \in \Delta_{\Pi}} \sup_{M \in \mathcal{M}_{\xi}^{\text{gl}}(\lambda)^c \cap \mathcal{M}_0} \frac{1}{\mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{\omega} [D(\bar{M}(\pi) \| M(\pi))]]} \\ & \leq \sup_{M \in \mathcal{M}_{\xi}^{\text{gl}}(\tilde{\lambda})^c \cap \mathcal{M}_0} \frac{1}{\mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{p_{\text{exp}}} [D(\bar{M}(\pi) \| M(\pi))]]} \\ & \leq \sup_{M \in \mathcal{M}_{\xi}^{\text{gl}}(\tilde{\lambda})^c \cap \mathcal{M}_0} \frac{\mathbb{I}\{\mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{p_{\text{exp}}} [D(\bar{M}(\pi) \| M(\pi))]] \leq 1/C_{\text{exp}}\}}{\mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{p_{\text{exp}}} [D(\bar{M}(\pi) \| M(\pi))]]} + C_{\text{exp}}, \end{aligned} \quad (44)$$

where here we take $\frac{0}{0} = 0$. If $\mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{p_{\text{exp}}} [D(\bar{M}(\pi) \| \widetilde{M}(\pi))]] > 1/C_{\text{exp}}$, then by definition of \widetilde{M} , for all $M \in \mathcal{M}_0$ we have

$$\mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{p_{\text{exp}}} [D(\bar{M}(\pi) \| M(\pi))]] > 1/C_{\text{exp}},$$

so we can simply bound Eq. (44) $\leq C_{\text{exp}}$. Going forward, we assume this is not the case, so that $\mathbb{E}_{\bar{M} \sim \xi}[\mathbb{E}_{p_{\text{exp}}}[D(\bar{M}(\pi) \parallel \widetilde{M}(\pi))]] \leq 1/C_{\text{exp}}$. Our goal is to show that, for small enough δ , $\widetilde{\lambda} \in \Lambda(\bar{M}; \varepsilon)$ for every other $M \in \mathcal{M}_0$ with $\mathbb{E}_{\bar{M} \sim \xi}[\mathbb{E}_{p_{\text{exp}}}[D(\bar{M}(\pi) \parallel \widetilde{M}(\pi))]] \leq 1/C_{\text{exp}}$, so that $M \in \mathcal{M}_{\varepsilon}^{\text{gl}}(\widetilde{\lambda})$. This will imply that there does not exist $M \in \mathcal{M}_{\varepsilon}^{\text{gl}}(\widetilde{\lambda})^c \cap \mathcal{M}_0$ with $\mathbb{E}_{\bar{M} \sim \xi}[\mathbb{E}_{p_{\text{exp}}}[D(\bar{M}(\pi) \parallel \widetilde{M}(\pi))]] \leq 1/C_{\text{exp}}$, which further implies that Eq. (44) $\leq C_{\text{exp}}$.

Fix any $M \in \mathcal{M}_0$. We note that by the definition of p_{exp} , if $\mathbb{E}_{\bar{M} \sim \xi}[\mathbb{E}_{p_{\text{exp}}}[D(\bar{M}(\pi) \parallel M(\pi))]] \leq 1/C_{\text{exp}}$, then

$$\sup_{p \in \Delta_{\Pi}} \mathbb{E}_p[D(\widetilde{M}(\pi) \parallel M(\pi))] \leq \delta,$$

This implies in particular that, for each π , $D(\widetilde{M}(\pi) \parallel M(\pi)) \leq \delta$.

Step 1: $\mathbb{E}_{\bar{M} \sim \xi}[\mathbb{E}_{p_{\text{exp}}}[D(\bar{M}(\pi) \parallel M(\pi))]] \leq 1/C_{\text{exp}}$ **implies** $\pi_M = \pi_{\bar{M}}$. As noted, if

$$\mathbb{E}_{\bar{M} \sim \xi}[\mathbb{E}_{p_{\text{exp}}}[D(\bar{M}(\pi) \parallel M(\pi))]] \leq 1/C_{\text{exp}},$$

we have $D(\widetilde{M}(\pi) \parallel M(\pi)) \leq \delta$ for all π . Assume that $\pi_M \neq \pi_{\bar{M}}$ (note that since $\inf_{M \in \mathcal{M}_0} \Delta_{\min}^M > 0$ by assumption, all $M \in \mathcal{M}_0$ have unique optimal). By definition we have $f^{\bar{M}}(\pi_{\bar{M}}) \geq f^{\bar{M}}(\pi_M) + \Delta_{\min}^{\bar{M}}$ and $f^M(\pi_M) \geq f^M(\pi_{\bar{M}})$. Thus,

$$\begin{aligned} \Delta_{\min}^{\bar{M}} &\leq f^{\bar{M}}(\pi_{\bar{M}}) - f^{\bar{M}}(\pi_M) + f^M(\pi_M) - f^M(\pi_{\bar{M}}) \\ &\leq |f^{\bar{M}}(\pi_{\bar{M}}) - f^{\bar{M}}(\pi_M)| + |f^M(\pi_M) - f^M(\pi_{\bar{M}})| \\ &\leq \sqrt{D(\widetilde{M}(\pi_{\bar{M}}) \parallel M(\pi_{\bar{M}}))} + \sqrt{D(\widetilde{M}(\pi_M) \parallel M(\pi_M))}. \end{aligned}$$

This implies that there exists some π such that $D(\widetilde{M}(\pi) \parallel M(\pi)) \geq (\Delta_{\min}^{\bar{M}}/2)^2$. Assuming

$$\delta \leq \min_{M \in \mathcal{M}_0} (\Delta_{\min}^{\bar{M}}/3)^2, \quad (45)$$

this is a contradiction. Thus, it must be the case that $\pi_M = \pi_{\bar{M}}$, as long as Eq. (45) is satisfied.

Step 2: $\mathbb{E}_{\bar{M} \sim \xi}[\mathbb{E}_{p_{\text{exp}}}[D(\bar{M}(\pi) \parallel M(\pi))]] \leq 1/C_{\text{exp}}$ **implies** $\widetilde{\lambda} \in \Lambda(\bar{M}; \varepsilon)$. Under Assumption D.2, we can bound, $\forall \pi \in \Pi$,

$$|f^{\bar{M}}(\pi) - f^M(\pi)| \leq \sqrt{D(\widetilde{M}(\pi) \parallel M(\pi))} \leq \sqrt{\delta}.$$

This implies that, for any $\lambda \in \Delta_{\Pi}$,

$$|\Delta^{\bar{M}}(\lambda) - \Delta^M(\lambda)| \leq |f^{\bar{M}}(\pi_M) - f^M(\pi_M)| + \sum_{\pi} \lambda_{\pi} |f^{\bar{M}}(\pi) - f^M(\pi)| \leq 4\sqrt{\delta}.$$

In addition, under Assumption D.1, we have

$$\begin{aligned} D_{\text{KL}}(M(\pi) \parallel M'(\pi)) &\geq D_{\text{KL}}(\widetilde{M}(\pi) \parallel M'(\pi)) - L_{\text{KL}} \sqrt{D(\widetilde{M}(\pi) \parallel M(\pi))} \\ &\geq D_{\text{KL}}(\widetilde{M}(\pi) \parallel M'(\pi)) - 2L_{\text{KL}} \sqrt{\delta}. \end{aligned}$$

This implies, for any $\lambda \in \Delta_\Pi$,

$$\begin{aligned} I^M(\lambda) &= \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi} \lambda(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) \\ &\geq \inf_{M' \in \mathcal{M}^{\text{alt}}(M)} \sum_{\pi} \lambda(\pi) D_{\text{KL}}(\widetilde{M}(\pi) \| M'(\pi)) - 2L_{\text{KL}}\sqrt{\delta} \\ &= I^{\widetilde{M}}(\lambda) - 2L_{\text{KL}}\sqrt{\delta} \end{aligned}$$

where the final equality uses that, given what we have already shown, $\pi_M = \pi_{\widetilde{M}}$, so that $\mathcal{M}^{\text{alt}}(M) = \mathcal{M}^{\text{alt}}(\widetilde{M})$. Repeating the calculation in the other direction, we get that $|I^M(\lambda) - I^{\widetilde{M}}(\lambda)| \leq 2L_{\text{KL}}\sqrt{\delta}$.

We next relate \mathbf{g}^M to $\mathbf{g}^{\widetilde{M}}$. By definition we have

$$(1 + \varepsilon')\mathbf{g}^M/\tilde{n} \geq \inf_{\lambda \in \Delta_\Pi} \Delta^M(\lambda) \quad \text{s.t.} \quad I^M(\lambda) \geq (1 - \varepsilon')/\tilde{n}.$$

Applying our perturbation bounds we can lower bound this as

$$\begin{aligned} &\geq \inf_{\lambda \in \Delta_\Pi} \Delta^{\widetilde{M}}(\lambda) - 4\sqrt{\delta} \quad \text{s.t.} \quad I^{\widetilde{M}}(\lambda) \geq (1 - \varepsilon')/\tilde{n} - 2L_{\text{KL}}\sqrt{\delta} \\ &\geq \frac{\mathbf{g}^{\widetilde{M}}}{((1 - \varepsilon')/\tilde{n} - 2L_{\text{KL}}\sqrt{\delta})^{-1}} - 4\sqrt{\delta} \end{aligned}$$

where the last inequality follows from [Lemma E.1](#). This implies that

$$\mathbf{g}^{\widetilde{M}} \leq ((1 - \varepsilon')/\tilde{n} - 2L_{\text{KL}}\sqrt{\delta})^{-1}(1 + \varepsilon') \cdot \frac{\mathbf{g}^M}{\tilde{n}} + 4((1 - \varepsilon')/\tilde{n} - 2L_{\text{KL}}\sqrt{\delta})^{-1}\sqrt{\delta}. \quad (46)$$

Assuming that

$$\sqrt{\delta} \leq \frac{\varepsilon' - 2(\varepsilon')^2}{2(1 + 2\varepsilon')L_{\text{KL}}\tilde{n}},$$

some algebra shows that

$$\text{Eq. (46)} \leq (1 + 2\varepsilon')^2\mathbf{g}^M + 4(1 + 2\varepsilon')\tilde{n}\sqrt{\delta}.$$

Now we can bound

$$\begin{aligned} \Delta^M(\tilde{\lambda}) &\leq \Delta^{\widetilde{M}}(\tilde{\lambda}) + 4\sqrt{\delta} \\ &\leq (1 + \varepsilon')\mathbf{g}^{\widetilde{M}}/\tilde{n} + 4\sqrt{\delta} \\ &\leq (1 + 2\varepsilon')^3\mathbf{g}^M/\tilde{n} + 4(1 + 2\varepsilon')^2\sqrt{\delta} + 4\sqrt{\delta} \end{aligned}$$

and

$$I^M(\tilde{\lambda}) \geq I^{\widetilde{M}}(\tilde{\lambda}) - 2L_{\text{KL}}\sqrt{\delta} \geq (1 - \varepsilon')/\tilde{n} - 2L_{\text{KL}}\sqrt{\delta}.$$

If ε' and δ are small enough so that

$$(1 + 2\varepsilon')^3 \leq 1 + \varepsilon/2 \quad \text{and} \quad 4(1 + 2\varepsilon')^2\sqrt{\delta} + 4\sqrt{\delta} \leq \varepsilon\mathbf{g}^M/2\tilde{n}$$

and

$$1 - \varepsilon' \geq 1 - \varepsilon/2 \quad \text{and} \quad 2L_{\text{KL}}\sqrt{\delta} \leq \varepsilon/2\tilde{n},$$

then $\Delta^M(\tilde{\lambda}) \leq (1 + \varepsilon)\mathbf{g}^M/\tilde{n}$ and $I^M(\tilde{\lambda}) \geq (1 - \varepsilon)/\tilde{n}$, which implies that $\tilde{\lambda} \in \Lambda(M; \varepsilon)$ with scaling factor \tilde{n} .

Step 3: Condition on δ . Altogether, we have assumed that δ satisfies [Eq. \(45\)](#) and, that for some $M \in \mathcal{M}_0$ with $\mathbb{E}_{\bar{M} \sim \xi}[\mathbb{E}_{p_{\text{exp}}}[D(\bar{M}(\pi) \| M(\pi))]] \leq 1/C_{\text{exp}}$, we have

$$\sqrt{\delta} \leq \frac{\varepsilon' - 2(\varepsilon')^2}{2(1 + 2\varepsilon')L_{\text{KL}}\tilde{n}}, \quad 4(1 + 2\varepsilon')^2\sqrt{\delta} + 4\sqrt{\delta} \leq \varepsilon\mathbf{g}^M/2\tilde{n}, \quad 2L_{\text{KL}}\sqrt{\delta} \leq \varepsilon/2\tilde{n} \quad (47)$$

and

$$(1 + 2\varepsilon')^3 \leq 1 + \varepsilon/2, \quad 1 - \varepsilon' \geq 1 - \varepsilon/2.$$

Some algebra shows that, as long as $\varepsilon \leq 1$, it suffices to take $\varepsilon' = \varepsilon/36$ to satisfy the latter two conditions. Furthermore, some calculation shows that a sufficient condition for [Eq. \(47\)](#) to be met is that

$$\sqrt{\delta} \leq \min \left\{ \frac{1}{81L_{\text{KL}}}, \frac{\mathbf{g}^M}{17} \right\} \cdot \frac{\varepsilon}{\tilde{n}}.$$

By [Lemma E.5](#) and our choice of \tilde{n} , we can bound

$$\tilde{n} \leq 2\mathbf{g}^{\bar{M}}/\Delta_{\min}^{\bar{M}} + n_{\varepsilon/36}^{\bar{M}}$$

so it suffices that we take

$$\sqrt{\delta} \leq \min \left\{ \frac{1}{81L_{\text{KL}}}, \frac{\mathbf{g}^M}{17} \right\} \cdot \varepsilon \left(\frac{2\mathbf{g}^{\bar{M}}}{\Delta_{\min}^{\bar{M}}} + n_{\varepsilon}^{\bar{M}} \right)^{-1}.$$

As \bar{M} was chosen to an arbitrary model in \mathcal{M}_0 with $\mathbb{E}_{\bar{M} \sim \xi}[\mathbb{E}_{p_{\text{exp}}}[D(\bar{M}(\pi) \| \bar{M}(\pi))]] \leq 1/C_{\text{exp}}$, we take it to minimize $\mathbf{g}^{\bar{M}}$ over this constraint. It suffices then that

$$\sqrt{\delta} \leq \min \left\{ \frac{1}{81L_{\text{KL}}}, \frac{\mathbf{g}^{\bar{M}}}{17} \right\} \cdot \varepsilon \left(\frac{2\mathbf{g}^{\bar{M}}}{\Delta_{\min}^{\bar{M}}} + n_{\varepsilon/36}^{\bar{M}} \right)^{-1}.$$

Finally, by [Lemma F.13](#) (under [Assumption A.2](#)) and [Lemma E.3](#), we can lower bound $\mathbf{g}^{\bar{M}} \geq \Delta_{\min}^{\bar{M}}/2V_{\mathcal{M}}$. Combining this condition with [Eq. \(45\)](#) gives the result. \square

Appendix F. Proofs from Appendix A

Organization of [Appendix F](#). In this section we prove the main results from [Appendix A](#). We consider a slightly generalized version of the setting in [Appendix A](#), where we allow for divergences other than just the KL divergence, as described below. This section is organized as follows.

- First, in [Appendix F.1](#), we give the proof of our main result, [Theorem A.1](#). We break this proof into two principle components: bounding the regret of AE^2 in the exploit phase ([Section F.1.1](#)), and explore phase ([Section F.1.2](#)). The key results in this section are [Lemma F.3](#), which formalizes the key algorithm intuition given in [Appendix A](#), showing that exploring via the AEC yields low regret, and [Lemma F.4](#), which shows that, to enter the explore phase, the total ‘‘information gain’’ must be bounded as $O(\log T)$, which ultimately yields the optimal leading-order scaling. We combine these results with our estimation guarantees in [Appendix F.1.3](#), where we give the proof of [Theorem A.1](#).

- In [Appendix F.2](#), we extend AE^2 and [Theorem A.1](#) to the case where we assume no lower bound on the minimum gap over the model class, first presenting our main algorithm in this setting, [Algorithm 6](#), and then giving a proof of [Theorem A.2](#). The structure of this section is similar to [Appendix F.1](#)—the primary difference being a slightly different argument to handle the need to adapt to the minimum gap of the ground truth instance.
- Finally, in [Appendix F.3](#) we present our estimation routine with covering, and prove that it achieves low estimation error, and in [Appendix F.4](#) we provide the proofs of miscellaneous results used throughout [Appendix F](#).

F.1. Regret Bound for Uniformly Regular Classes (Theorem A.1)

In this section we prove [Theorem F.1](#), which generalizes [Theorem A.1](#). To do so, we analyze [Algorithm 5](#), which generalizes [Algorithm 2](#) to allow for general divergences.

Throughout, we define

$$\underline{g}^{\mathcal{M}} := \min_{M \in \mathcal{M}: g^M > 0} g^M.$$

Algorithm 5 Allocation Estimation via Adaptive Exploration (AE^2 , general divergences)

- 1: **input:** Optimality tolerance δ , model class \mathcal{M} , estimation oracle Alg_D .
 - 2: Initialize $s \leftarrow 1$, $\varepsilon \leftarrow \frac{\delta}{4+2\delta}$, $n_{\max} \leftarrow n_{\max}(\mathcal{M}, \varepsilon)$, $q \leftarrow \frac{4n_{\max} + \delta \underline{g}^{\mathcal{M}}}{4n_{\max} + 2\delta \underline{g}^{\mathcal{M}}}$.
 - 3: Compute $\xi^1 \leftarrow \text{Alg}_D(\{\emptyset\})$ and $\widehat{M}^1 \leftarrow \mathbb{E}_{M \sim \xi^1}[M]$.
 - 4: **for** $t = 1, 2, 3, \dots$ **do**
 - 5: **if** $\exists \pi_{\widehat{M}^s} \in \pi_{\widehat{M}^s}$ s.t. $\forall M \in \mathcal{M}^{\text{alt}}(\pi_{\widehat{M}^s})$, $\sum_{i=1}^{s-1} \mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}_{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}_{M, \pi^i}(r^i, o^i)} \right] \geq \log(t \log t)$ **then**
 - 6: Play $\pi_{\widehat{M}^s}$. // Exploit
 - 7: **else** // Explore
 - 8: Set $p^s \leftarrow q\lambda^s + (1-q)\omega^s$ for
 - $$\lambda^s, \omega^s \leftarrow \arg \min_{\lambda, \omega \in \Delta_{\Pi}} \sup_{M \in \mathcal{M} \setminus \mathcal{M}_{\varepsilon}^{\text{sl}}(\lambda; n_{\max})} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim \omega} [D(\widehat{M}(\pi) \| M(\pi))]]}. \quad (48)$$
 - 9: Draw $\pi^s \sim p^s$, observe r^s, o^s .
 - 10: Compute estimate $\xi^{s+1} \leftarrow \text{Alg}_D(\{(\pi^i, r^i, o^i)\}_{i=1}^s)$ and let $\widehat{M}^s = \mathbb{E}_{M \sim \xi^s}[M]$.
 - 11: $s \leftarrow s + 1$.
-

F.1.1. BOUNDING REGRET OF EXPLOIT PHASE

We refer to the *exploit phase* as the subset of rounds t in which [Line 6](#) is reached, and refer to the *explore phase* as the subset of rounds in which [Line 8](#) is reached.

Lemma F.1. *The total expected regret incurred by the exploit phase of [Algorithm 5](#) is bounded by $2 \log \log T + 3$.*

Proof of Lemma F.1. Let \mathcal{E}_t denote the event that we exploit at round t , and that $\pi_{\widehat{M}^s} \neq \pi_\star$. Then, since we can incur suboptimality of at most 1 at each round, the total expected regret incurred by the exploit phase is bounded by

$$\sum_{t=1}^T \mathbb{E}^{M^\star} [\mathbb{I}\{\mathcal{E}_t\}].$$

Let $\widetilde{\mathcal{E}}_t$ denote the event

$$\widetilde{\mathcal{E}}_t := \left\{ \forall s \geq 1 : \sum_{i=1}^s \mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}^{M^\star, \pi^i}(r^i, o^i)} \right] < \log(t \log t) \right\}.$$

By Lemma F.2, we have $\mathbb{P}^{M^\star} [\mathbb{I}\{\widetilde{\mathcal{E}}_t^c\}] \leq \frac{1}{t \log t}$, and we can bound

$$\mathbb{E}^{M^\star} [\mathbb{I}\{\mathcal{E}_t\}] \leq \mathbb{E}^{M^\star} [\mathbb{I}\{\mathcal{E}_t \cap \widetilde{\mathcal{E}}_t\}] + \mathbb{E}^{M^\star} [\mathbb{I}\{\widetilde{\mathcal{E}}_t^c\}] \leq \mathbb{E}^{M^\star} [\mathbb{I}\{\mathcal{E}_t \cap \widetilde{\mathcal{E}}_t\}] + \frac{1}{t \log t}.$$

Let s_t denote the exploration round at round t . If we exploit at round t , this implies that for all $M \in \mathcal{M}^{\text{alt}}(\pi_{\widehat{M}^s})$, we have

$$\sum_{i=1}^{s_t-1} \mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}^{M, \pi^i}(r^i, o^i)} \right] \geq \log(t \log t). \quad (49)$$

If $\pi_{\widehat{M}^s} \neq \pi_\star$, then $M^\star \in \mathcal{M}^{\text{alt}}(\pi_{\widehat{M}^s})$, so Eq. (49) must hold for $M \leftarrow M^\star$. This contradicts $\widetilde{\mathcal{E}}_t$, however, so $\mathbb{E}^{M^\star} [\mathbb{I}\{\mathcal{E}_t \cap \widetilde{\mathcal{E}}_t\}] = 0$. Thus, $\mathbb{E}^{M^\star} [\mathbb{I}\{\mathcal{E}_t\}] \leq \frac{1}{t \log t}$, so

$$\sum_{t=1}^T \mathbb{E}^{M^\star} [\mathbb{I}\{\mathcal{E}_t\}] \leq 3 + \sum_{t=3}^T \frac{1}{t \log t} \leq 3 + 2 \int_e^T \frac{1}{t \log t} dt = 3 + 2 \log \log T.$$

□

Lemma F.2. For $\{(r^i, o^i, \xi^i)\}_{i=1}^s$ generated as in Algorithm 5, we have that

$$\mathbb{P}^{M^\star} \left[\exists s \geq 1 : \sum_{i=1}^s \mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}^{M^\star, \pi^i}(r^i, o^i)} \right] \geq \varepsilon \right] \leq e^{-\varepsilon}.$$

Proof of Lemma F.2. Denote

$$X_s := \exp \left(\sum_{i=1}^s \mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}^{M^\star, \pi^i}(r^i, o^i)} \right] \right).$$

We first show that X_s is a supermartingale. Letting \mathcal{F}_{s-1} denote the filtration up to $s-1$, we have

$$\mathbb{E}^{M^\star} [X_s \mid \mathcal{F}_{s-1}] = \exp \left(\sum_{i=1}^{s-1} \mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}^{M^\star, \pi^i}(r^i, o^i)} \right] \right) \cdot \mathbb{E}^{M^\star} \left[\exp \left(\mathbb{E}_{\widehat{M} \sim \xi^s} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^s}(r^s, o^s)}{\mathbb{P}^{M^\star, \pi^s}(r^s, o^s)} \right] \right) \mid \mathcal{F}_{s-1} \right]$$

$$\begin{aligned}
 &= X_{s-1} \cdot \mathbb{E}^{M^*} \left[\exp \left(\mathbb{E}_{\widehat{M} \sim \xi^s} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^s}(r^s, o^s)}{\mathbb{P}^{M^*, \pi^s}(r^s, o^s)} \right] \right) \mid \mathcal{F}_{s-1} \right] \\
 &\stackrel{(a)}{\leq} X_{s-1} \cdot \mathbb{E}^{M^*} \left[\mathbb{E}_{\widehat{M} \sim \xi^s} \left[\exp \left(\log \frac{\mathbb{P}^{\widehat{M}, \pi^s}(r^s, o^s)}{\mathbb{P}^{M^*, \pi^s}(r^s, o^s)} \right) \right] \mid \mathcal{F}_{s-1} \right] \\
 &= X_{s-1} \cdot \mathbb{E}^{M^*} \left[\frac{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{P}^{\widehat{M}, \pi^s}(r^s, o^s)]}{\mathbb{P}^{M^*, \pi^s}(r^s, o^s)} \mid \mathcal{F}_{s-1} \right] \\
 &= X_{s-1}
 \end{aligned}$$

where (a) holds by Jensen's inequality, and the final equality holds since $\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{P}^{\widehat{M}, \pi^s}(\cdot, \cdot)]$ is a valid distribution over $\mathcal{R} \times \mathcal{O}$. Thus, X_s is a supermartingale. Ville's Maximal Inequality then immediately gives that

$$\mathbb{P}^{M^*} [\exists s \geq 1 : X_s \geq e^\varepsilon] \leq \frac{\mathbb{E}^{M^*}[X_1]}{e^\varepsilon}.$$

To complete the proof, using the same calculation as above, we bound

$$\mathbb{E}^{M^*}[X_1] = \mathbb{E}^{M^*} \left[\exp \left(\mathbb{E}_{\widehat{M} \sim \xi^1} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^1}(r^1, o^1)}{\mathbb{P}^{M^*, \pi^1}(r^1, o^1)} \right] \right) \right] \leq \mathbb{E}^{M^*} \left[\mathbb{E}_{\widehat{M} \sim \xi^1} \left[\exp \left(\log \frac{\mathbb{P}^{\widehat{M}, \pi^1}(r^1, o^1)}{\mathbb{P}^{M^*, \pi^1}(r^1, o^1)} \right) \right] \right] \leq 1.$$

□

F.1.2. BOUNDING REGRET OF EXPLORE PHASE

Lemma F.3 (Main Explore Phase Regret Bound). *Let s_T denote the total number of exploration rounds (which is a random variable), and assume that $\delta \in [0, 1/2)$. Then running [Algorithm 5](#), if $g^* > 0$, we can bound*

$$\begin{aligned}
 \mathbb{E}[s_T] &\leq \frac{24n_{\max}^2 + 8n_{\max}g^{\mathcal{M}}}{(\delta g^{\mathcal{M}})^2} \cdot \overline{\text{aec}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}) \cdot \mathbb{E}[\widehat{\text{Est}}_{\text{D}}(s_T)] + \frac{12n_{\max}}{\delta \Delta_{\min}} \cdot \mathbb{E}[\text{Est}_{\text{KL}}(s_T)] \\
 &\quad + \frac{6n_{\max}}{\delta} \cdot \mathbb{E} \left[\sum_{s=1}^{s_T} \inf_{M \in \mathcal{M}^{\text{alt}}(M^*)} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \cdot \mathbb{I}\{\pi_* \in \boldsymbol{\pi}_{\widehat{M}^s}\} \right]
 \end{aligned}$$

and the regret during exploration rounds is bounded as

$$\begin{aligned}
 \mathbb{E} \left[\sum_{s=1}^{s_T} \Delta^*(\pi^s) \right] &\leq \frac{8n_{\max} + 2g^{\mathcal{M}}}{\delta g^{\mathcal{M}}} \cdot \overline{\text{aec}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}) \cdot \mathbb{E}[\widehat{\text{Est}}_{\text{D}}(s_T)] + \frac{2(1+\delta)g^*}{\Delta_{\min}} \cdot \mathbb{E}[\text{Est}_{\text{KL}}(s_T)] \\
 &\quad + (1+\delta)g^* \cdot \mathbb{E} \left[\sum_{s=1}^{s_T} \inf_{M \in \mathcal{M}^{\text{alt}}(M^*)} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \cdot \mathbb{I}\{\pi_* \in \boldsymbol{\pi}_{\widehat{M}^s}\} \right].
 \end{aligned}$$

Proof of Lemma F.3. The expected regret during exploration can be written as $\mathbb{E}[\sum_{s=1}^{s_T} \Delta^*(\pi^s)] = \mathbb{E}[\sum_{s=1}^{s_T} \mathbb{E}_{p^s}[\Delta^*(\pi)]]$. By definition, for exploration rounds, we have $p^s \leftarrow q\lambda^s + (1-q)\omega^s$. For each $s \leq s_T$, we consider three cases to bound the instantaneous expected regret, $\mathbb{E}_{p^s}[\Delta^*(\pi)] = \Delta^*(p^s)$.

Case 1: $M^* \in \mathcal{M} \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda^s; n_{\max})$. Denote such rounds as $\mathcal{S}_{\text{exp}}^1$. Write

$$\Delta^*(p^s) = \left[\Delta^*(p^s) - \gamma^s \mathbb{E}_{\widehat{M} \sim \xi^s} \left[\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))] \right] \right] + \gamma^s \mathbb{E}_{\widehat{M} \sim \xi^s} \left[\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))] \right]$$

for

$$\gamma^s := \frac{1 + \delta}{1 - q} \cdot \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\omega^s} [D(\widehat{M}(\pi) \| M^*(\pi))]]}. \quad (50)$$

In this case we have that

$$\begin{aligned} & \gamma^s \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]] \\ &= \frac{1 + \delta}{1 - q} \cdot \frac{1}{\mathbb{E}_{\omega^s} [\mathbb{E}_{\widehat{M} \sim \xi^s} [D(\widehat{M}(\pi) \| M^*(\pi))]]} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]] \\ &\geq \frac{1 + \delta}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]]} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]] \\ &= 1 + \delta. \end{aligned}$$

Thus, since the suboptimality gap is always bounded by 1, we can bound

$$\Delta^*(p^s) - \gamma^s \mathbb{E}_{\widehat{M} \sim \xi^s} \left[\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))] \right] \leq 1 - \gamma^s \mathbb{E}_{\widehat{M} \sim \xi^s} \left[\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))] \right] \leq -\delta.$$

So,

$$\Delta^*(p^s) \leq -\delta + \gamma^s \mathbb{E}_{\widehat{M} \sim \xi^s} \left[\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))] \right].$$

Case 2: $M^* \in \mathcal{M}_\varepsilon^{\text{gl}}(\lambda^s; n_{\max})$, $\pi_* \in \pi_{M^*}$. Denote such rounds as $\mathcal{S}_{\text{exp}}^2$, and write

$$\Delta^*(p^s) = [\Delta^*(p^s) - (1 + \delta) \mathbf{g}^* \cdot \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))]] + (1 + \delta) \mathbf{g}^* \cdot \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))]$$

for any $M \in \mathcal{M}_{\text{alt}}^*$. In this case, since $M^* \in \mathcal{M}_\varepsilon^{\text{gl}}(\lambda^s; n_{\max})$, we have that $\lambda^s \in \Lambda(M^*; \varepsilon)$. This then implies that

$$\Delta^*(\lambda^s) \leq (1 + \varepsilon) \mathbf{g}^* / n^* \quad \text{and} \quad \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{\lambda^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \geq (1 - \varepsilon) / n^*$$

for some $n^* \leq n_{\max}$. Since $M \in \mathcal{M}_{\text{alt}}^*$, it follows that $\mathbb{E}_{\lambda^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \geq (1 - \varepsilon) / n^*$.

Thus,

$$\begin{aligned} \Delta^*(p^s) - (1 + \delta) \mathbf{g}^* \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] &\leq q [\Delta^*(\lambda^s) - (1 + \delta) \mathbf{g}^* \mathbb{E}_{\lambda^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))]] + 1 - q \\ &\leq q [(1 + \varepsilon) \mathbf{g}^* / n^* - (1 + \delta)(1 - \varepsilon) \mathbf{g}^* / n^*] + 1 - q \\ &= q [2\varepsilon - \delta(1 - \varepsilon)] \cdot \frac{\mathbf{g}^*}{n^*} + 1 - q \\ &\stackrel{(a)}{=} -\frac{q\delta}{2} \frac{\mathbf{g}^*}{n^*} + 1 - q \\ &\leq -\frac{q\delta}{2} \frac{\mathbf{g}^*}{n_{\max}} + 1 - q \end{aligned}$$

$$\stackrel{(b)}{\leq} -\frac{\delta}{4} \frac{\mathbf{g}^*}{n_{\max}},$$

where (a) follows from our choice of $\varepsilon = \frac{\delta/2}{2+\delta}$ and setting of n^* , and (b) follows from our setting of $q = \frac{4n_{\max} + \delta \mathbf{g}^*}{4n_{\max} + 2\delta \mathbf{g}^*}$, since $\mathbf{g}^* > 0$, and some algebra. Thus,

$$\Delta^*(p^s) \leq (1 + \delta) \mathbf{g}^* \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] - \frac{\delta}{4} \frac{\mathbf{g}^*}{n_{\max}}.$$

As this holds for every $M \in \mathcal{M}_{\text{alt}}^*$, we therefore have

$$\Delta^*(p^s) \leq \inf_{M \in \mathcal{M}_{\text{alt}}^*} (1 + \delta) \mathbf{g}^* \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] - \frac{\delta}{4} \frac{\mathbf{g}^*}{n_{\max}}.$$

Case 3: $M^* \in \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda^s; n_{\max})$, $\pi_{\star} \notin \pi_{\widehat{M}^s}$. Denote such rounds as $\mathcal{S}_{\text{exp}}^3$, and write

$$\begin{aligned} \Delta^*(p^s) &= \left[\Delta^*(p^s) - \frac{2(1 + \delta) \mathbf{g}^*}{\Delta_{\min}} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi))] \right] \\ &\quad + \frac{2(1 + \delta) \mathbf{g}^*}{\Delta_{\min}} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi))]]. \end{aligned}$$

Since $M^* \in \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda^s; n_{\max})$, we have that $\lambda^s \in \Lambda(M^*; \varepsilon)$. This then implies that for any $M \in \mathcal{M}_{\text{alt}}^*$:

$$\Delta^*(\lambda^s) \leq (1 + \varepsilon) \mathbf{g}^* / n^* \quad \text{and} \quad \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{\lambda^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \geq (1 - \varepsilon) / n^*$$

for some $n^* \leq n_{\max}$. By [Lemma F.9](#), since $\pi_{\star} \notin \pi_{\widehat{M}^s}$, we can lower bound $\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{I}\{\widehat{M} \in \mathcal{M}_{\text{alt}}^*\}] \geq \frac{1}{2} \Delta_{\min}$. Thus, we have

$$\begin{aligned} \frac{2(1 + \delta) \mathbf{g}^*}{\Delta_{\min}} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\lambda^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi))] &\geq \frac{2(1 + \delta) \mathbf{g}^*}{\Delta_{\min}} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\lambda^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi)) \cdot \mathbb{I}\{\widehat{M} \in \mathcal{M}_{\text{alt}}^*\}]] \\ &\geq \frac{(1 + \delta)(1 - \varepsilon) \mathbf{g}^*}{n^*}. \end{aligned}$$

This implies that

$$\begin{aligned} \Delta^*(p^s) - \frac{2(1 + \delta) \mathbf{g}^*}{\Delta_{\min}} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi))] &\leq q[(1 + \varepsilon) \mathbf{g}^* / n^* - (1 + \delta)(1 - \varepsilon) \mathbf{g}^* / n^*] + 1 - q \\ &\leq -\frac{\delta}{4} \frac{\mathbf{g}^*}{n_{\max}}, \end{aligned}$$

where the final inequality follows by the same argument as in Case 2. Thus,

$$\Delta^*(p^s) \leq \frac{2(1 + \delta) \mathbf{g}^*}{\Delta_{\min}} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi))] - \frac{\delta}{4} \frac{\mathbf{g}^*}{n_{\max}}.$$

Completing the Proof. In total we have

$$\begin{aligned}
 \mathbb{E} \left[\sum_{s=1}^{s_T} \Delta^*(p^s) \right] &\leq \mathbb{E} \left[-\delta |\mathcal{S}_{\text{exp}}^1| - \frac{\delta \mathbf{g}^*}{4n_{\text{max}}} |\mathcal{S}_{\text{exp}}^2 \cup \mathcal{S}_{\text{exp}}^3| + \sum_{s \in \mathcal{S}_{\text{exp}}^1} \gamma^s \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]] \right. \\
 &\quad + (1 + \delta) \mathbf{g}^* \sum_{s \in \mathcal{S}_{\text{exp}}^2} \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \\
 &\quad \left. + \frac{2(1 + \delta) \mathbf{g}^*}{\Delta_{\text{min}}} \sum_{s \in \mathcal{S}_{\text{exp}}^3} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi))]] \right] \\
 &\leq \mathbb{E} \left[-\frac{\delta \mathbf{g}^*}{4n_{\text{max}}} \cdot s_T + \frac{8n_{\text{max}} + 2\mathbf{g}^{\mathcal{M}}}{\delta \mathbf{g}^{\mathcal{M}}} \cdot \overline{\text{aec}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}) \cdot \widehat{\mathbf{Est}}_{\text{D}}(s_T) + \frac{2(1 + \delta) \mathbf{g}^*}{\Delta_{\text{min}}} \cdot \mathbf{Est}_{\text{KL}}(s_T) \right. \\
 &\quad \left. + (1 + \delta) \mathbf{g}^* \sum_{s=1}^{s_T} \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \cdot \mathbb{I}\{\pi_{\star} \in \pi_{\widehat{M}^s}\} \right]
 \end{aligned}$$

where the last inequality follows from [Lemma F.10](#), which bounds

$$\gamma^s \leq \frac{8n_{\text{max}} + 2\mathbf{g}^{\mathcal{M}}}{\delta \mathbf{g}^{\mathcal{M}}} \cdot \overline{\text{aec}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}),$$

and also using that for any $s \in \mathcal{S}_{\text{exp}}^2$ we have $\pi_{\star} \in \pi_{\widehat{M}^s}$. Upper bounding $-\frac{\delta \mathbf{g}^*}{4n_{\text{max}}} \cdot s_T \leq 0$ proves the second claim in the lemma statement.

For the first claim, as regret is always nonnegative, it follows that

$$\begin{aligned}
 0 &\leq \mathbb{E} \left[-\frac{\delta \mathbf{g}^*}{4n_{\text{max}}} \cdot s_T + \frac{8n_{\text{max}} + 2\mathbf{g}^{\mathcal{M}}}{\delta \mathbf{g}^{\mathcal{M}}} \cdot \overline{\text{aec}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}) \cdot \widehat{\mathbf{Est}}_{\text{D}}(s_T) + \frac{2(1 + \delta) \mathbf{g}^*}{\Delta_{\text{min}}} \cdot \mathbf{Est}_{\text{KL}}(s_T) \right. \\
 &\quad \left. + (1 + \delta) \mathbf{g}^* \cdot \sum_{s=1}^{s_T} \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \cdot \mathbb{I}\{\pi_{\star} \in \pi_{\widehat{M}^s}\} \right]
 \end{aligned}$$

which implies

$$\begin{aligned}
 \mathbb{E}[s_T] &\leq \frac{4n_{\text{max}}}{\delta \mathbf{g}^*} \cdot \mathbb{E} \left[\frac{8n_{\text{max}} + 2\mathbf{g}^{\mathcal{M}}}{\delta \mathbf{g}^{\mathcal{M}}} \cdot \overline{\text{aec}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}) \cdot \widehat{\mathbf{Est}}_{\text{D}}(s_T) + \frac{2(1 + \delta) \mathbf{g}^*}{\Delta_{\text{min}}} \cdot \mathbf{Est}_{\text{KL}}(s_T) \right. \\
 &\quad \left. + (1 + \delta) \mathbf{g}^* \cdot \sum_{s=1}^{s_T} \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \cdot \mathbb{I}\{\pi_{\star} \in \pi_{\widehat{M}^s}\} \right].
 \end{aligned}$$

□

Lemma F.4. When running both [Algorithm 5](#) and [Algorithm 6](#), for all $\alpha \in [0, 1)$, we have

$$\begin{aligned}
 &\mathbb{E} \left[\sum_{s=1}^{s_T} \inf_{M \in \mathcal{M}_{\text{alt}}^*(M^*)} s^\alpha \cdot \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \cdot \mathbb{I}\{\pi_{\star} \in \pi_{\widehat{M}^s}\} \right] \\
 &\leq \mathbb{E}[s_T^\alpha] \log T + \mathbb{E} \left[V_{\mathcal{M}} s_T^{1/2+\alpha} \sqrt{1344 d_{\text{cov}} \cdot \log(128 C_{\text{cov}} s_T)} \right]
 \end{aligned}$$

$$\begin{aligned}
 & + (V_{\mathcal{M}} + L_{\text{KL}}) \left(4 \frac{s_T^{1/2+\alpha/2}}{1-\alpha} + \sum_{s=1}^{s_T} s^{1/2+\alpha/2} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]] \right) \\
 & + \mathbb{E}[s_T^\alpha] \log \log T + 7V_{\mathcal{M}}.
 \end{aligned}$$

Proof of Lemma F.4. Let \tilde{s}_T denote the final exploration round for which $\pi_\star \in \pi_{\widehat{M}^s}$. Then, upper bounding the KL divergence by $2V_{\mathcal{M}}$ via Lemma F.13, we have

$$\begin{aligned}
 & \sum_{s=1}^{s_T} \inf_{M \in \mathcal{M}_{\text{alt}}^*} s^\alpha \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \cdot \mathbb{I}\{\pi_\star \in \pi_{\widehat{M}^s}\} \\
 & \leq \tilde{s}_T^\alpha \sum_{s=1}^{\tilde{s}_T-1} \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] + 2V_{\mathcal{M}} s_T^\alpha.
 \end{aligned}$$

Under Assumption D.1 and via Jensen's inequality and AM-GM, we have, for any $\alpha_s > 0$ and $M \in \mathcal{M}$,

$$\begin{aligned}
 \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] & \leq \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M(\pi))]] + L_{\text{KL}} \sqrt{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]]} \\
 & \leq \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M(\pi))]] + L_{\text{KL}} \left(\frac{1}{\alpha_s} + \alpha_s \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]] \right),
 \end{aligned}$$

We now wish to bound

$$\mathbb{E} \left[\tilde{s}_T^\alpha \sum_{s=1}^{\tilde{s}_T-1} \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M(\pi))]] \right].$$

Let $\mathcal{M}_{\text{cov}}^j$ denote an (ρ_j, μ_j) -cover of $\mathcal{M}_{\text{alt}}^*$ and \mathcal{E}_j^s the corresponding event at step $s \leq 2^j$, for some ρ_j and μ_j to be chosen. Let $\mathcal{E}_j := \cap_{s \leq 2^j} \mathcal{E}_j^s$. By definition $\mathbb{P}_{M^*} [\mathcal{E}_j^s] \geq 1 - \mu_j$, so $\mathbb{P}_{M^*} [\mathcal{E}_j] \geq 1 - 2^j \mu_j$. Define an event

$$\begin{aligned}
 A_j & := \left\{ \forall s \leq 2^j, M \in \mathcal{M}_{\text{cov}}^j : (s+1)^\alpha \sum_{i=1}^s \mathbb{E}_{\widehat{M} \sim \xi^i} [\mathbb{E}_{p^i} [D_{\text{KL}}(\widehat{M}(\pi) \| M(\pi))]] \right. \\
 & \leq (s+1)^\alpha \sum_{i=1}^s \mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}^{M, \pi^i}(r^i, o^i)} \right] + V_{\mathcal{M}} (s+1)^{\frac{1}{2}+\alpha} \sqrt{56 \log \frac{2^j \text{N}_{\text{cov}}(\mathcal{M}_{\text{alt}}^*, \rho_j, \mu_j)}{\delta_j}} \\
 & \left. + V_{\mathcal{M}} \cdot \left(4 \frac{(s+1)^{\frac{1+\alpha}{2}}}{1-\alpha} + \sum_{i=1}^s i^{\frac{1+\alpha}{2}} \cdot \mathbb{E}_{\widehat{M} \sim \xi^i} [\mathbb{E}_{\pi \sim p^i} [D(\widehat{M}(\pi) \| M^*(\pi))]] \right) \right\}
 \end{aligned}$$

for some δ_j to be chosen. By invoking Lemma F.12 with $\beta_i = i^{1/2+\alpha/2}/(s+1)^\alpha$ and a union bound, we have $\mathbb{P}[A_j] \geq 1 - \delta_j$ (while Lemma F.12 does not contain the $(s+1)^\alpha$ term, the bound in the expression for A_j simply gives the bound from Lemma F.12 multiplied through with $(s+1)^\alpha$, which is non-random, so this is admissible). We can decompose

$$\mathbb{E} \left[\tilde{s}_T^\alpha \sum_{s=1}^{\tilde{s}_T-1} \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M(\pi))]] \right]$$

$$\begin{aligned}
 &\leq \sum_{j=1}^{\lceil \log T \rceil} \mathbb{E} \left[\tilde{s}_T^\alpha \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \sum_{s=1}^{\tilde{s}_T-1} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M(\pi))]] \cdot \mathbb{I}\{\tilde{s}_T \in [2^{j-1}, 2^j], A_j \cap \mathcal{E}_j\} \right] \\
 &\quad + \sum_{j=1}^{\lceil \log T \rceil} \mathbb{E} \left[\tilde{s}_T^\alpha \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \sum_{s=1}^{\tilde{s}_T-1} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M(\pi))]] \cdot \mathbb{I}\{\tilde{s}_T \in [2^{j-1}, 2^j], A_j^c \cup \mathcal{E}_j^c\} \right].
 \end{aligned}$$

We bound these terms separately. We can bound the second term as

$$\begin{aligned}
 &\sum_{j=1}^{\lceil \log T \rceil} \mathbb{E} \left[\tilde{s}_T^\alpha \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \sum_{s=1}^{\tilde{s}_T-1} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M(\pi))]] \cdot \mathbb{I}\{\tilde{s}_T \in [2^{j-1}, 2^j], A_j^c \cup \mathcal{E}_j^c\} \right] \\
 &\leq \sum_{j=1}^{\lceil \log T \rceil} 2V_{\mathcal{M}} 2^{2j} (\mathbb{P}[A_j^c] + \mathbb{P}[\mathcal{E}_j^c]) \\
 &\leq \sum_{j=1}^{\lceil \log T \rceil} 2V_{\mathcal{M}} 2^{2j} (\delta_j + 2^j \mu_j)
 \end{aligned}$$

where the first inequality follows by bounding [Lemma F.13](#), and $\tilde{s}_T \leq 2^j$, and the second follows by our bound on the probability of A_j . Letting $\delta_j = \frac{1}{2^{3j}}$, $\mu_j = \frac{1}{2^{4j}}$, this term can be bounded by $4V_{\mathcal{M}}$. We turn now to the first term. Fix $j \in [\lceil \log T \rceil]$. By definition of the event A_j , and since $\tilde{s}_T \leq 2^j$, plugging in our choice of δ_j we can bound, for any $M \in \mathcal{M}_{\text{cov}}^j$,

$$\begin{aligned}
 &\mathbb{E} \left[\tilde{s}_T^\alpha \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \sum_{s=1}^{\tilde{s}_T-1} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(\widehat{M}(\pi) \| M(\pi))]] \cdot \mathbb{I}\{\tilde{s}_T \in [2^{j-1}, 2^j], A_j \cap \mathcal{E}_j\} \right] \\
 &\leq \mathbb{E} \left[\left(\tilde{s}_T^\alpha \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \sum_{s=1}^{\tilde{s}_T-1} \mathbb{E}_{\widehat{M} \sim \xi^s} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^s}(r^s, o^s)}{\mathbb{P}^{M, \pi^s}(r^s, o^s)} \right] + V_{\mathcal{M}} \tilde{s}_T^{1/2+\alpha} \sqrt{56 \log(2^{3j} \mathbf{N}_{\text{cov}}(\mathcal{M}_{\text{alt}}^*, \rho_j, \mu_j))} \right. \right. \\
 &\quad \left. \left. + V_{\mathcal{M}} \cdot \left(4 \frac{\tilde{s}_T^{1/2+\alpha/2}}{1-\alpha} + \sum_{s=1}^{\tilde{s}_T} s^{1/2+\alpha/2} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]] \right) \right) \cdot \mathbb{I}\{\tilde{s}_T \in [2^{j-1}, 2^j], \mathcal{E}_j\} \right] \\
 &\leq \mathbb{E} \left[\left(\tilde{s}_T^\alpha \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \sum_{s=1}^{\tilde{s}_T-1} \mathbb{E}_{\widehat{M} \sim \xi^s} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^s}(r^s, o^s)}{\mathbb{P}^{M, \pi^s}(r^s, o^s)} \right] + V_{\mathcal{M}} \tilde{s}_T^{1/2+\alpha} \sqrt{168 \log(8\tilde{s}_T \mathbf{N}_{\text{cov}}(\mathcal{M}_{\text{alt}}^*, \rho_j, \frac{\tilde{s}_T^4}{16}))} \right. \right. \\
 &\quad \left. \left. + V_{\mathcal{M}} \cdot \left(4 \frac{\tilde{s}_T^{1/2+\alpha/2}}{1-\alpha} + \sum_{s=1}^{\tilde{s}_T} s^{1/2+\alpha/2} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]] \right) \right) \cdot \mathbb{I}\{\tilde{s}_T \in [2^{j-1}, 2^j], \mathcal{E}_j\} \right] \\
 &\tag{51}
 \end{aligned}$$

where the second inequality follows since, if $\tilde{s}_T \in [2^{j-1}, 2^j]$, then we can bound $2^j \leq 2\tilde{s}_T$.

Since \tilde{s}_T is an exploration round by definition, we know that for all $\pi_{\widehat{M}^{\tilde{s}_T}} \in \pi_{\widehat{M}^{\tilde{s}_T}}$, there exists some $M' \in \mathcal{M}_{\text{alt}}^{\text{alt}}(\pi_{\widehat{M}^{\tilde{s}_T}})$ such that

$$\sum_{s=1}^{\tilde{s}_T-1} \mathbb{E}_{\widehat{M} \sim \xi^s} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^s}(r^s, o^s)}{\mathbb{P}^{M', \pi^s}(r^s, o^s)} \right] \leq \log(T \log T).$$

By assumption we have $\pi_\star \in \pi_{\widehat{M}^{\tilde{s}_T}}$, so it follows that there exists some (random) $M' \in \mathcal{M}_{\text{alt}}^\star$ such that the above inequality holds. Let $M'' \in \mathcal{M}_{\text{cov}}^j$ denote the model in $\mathcal{M}_{\text{cov}}^j$ such that

$$\left| \log \frac{\mathbb{P}^{M', \pi}(r, o)}{\mathbb{P}^{M'', \pi}(r, o)} \right| = \left| \log \mathbb{P}^{M', \pi}(r, o) - \log \mathbb{P}^{M'', \pi}(r, o) \right| \leq \rho_j,$$

for all (r, o, π) for which $\sup_{\widehat{M} \in \mathcal{M}} \mathbb{P}^{\widehat{M}, \pi}(r, o \mid \mathcal{E}_j) > 0$, which is guaranteed to exist by [Definition A.1](#). Note that by definition of $\mathcal{M}_{\text{cov}}^j$, we have $M'' \in \mathcal{M}_{\text{alt}}^\star$. We then have

$$\begin{aligned} & \tilde{s}_T^\alpha \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^\star} \sum_{s=1}^{\tilde{s}_T-1} \mathbb{E}_{\widehat{M} \sim \xi^s} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^s}(r^s, o^s)}{\mathbb{P}^{M, \pi^s}(r^s, o^s)} \right] \cdot \mathbb{I}\{\tilde{s}_T \in [2^{j-1}, 2^j], \mathcal{E}_j\} \\ & \leq \tilde{s}_T^\alpha \sum_{s=1}^{\tilde{s}_T-1} \mathbb{E}_{\widehat{M} \sim \xi^s} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^s}(r^s, o^s)}{\mathbb{P}^{M'', \pi^s}(r^s, o^s)} \right] \cdot \mathbb{I}\{\tilde{s}_T \in [2^{j-1}, 2^j], \mathcal{E}_j\} \\ & = \tilde{s}_T^\alpha \sum_{s=1}^{\tilde{s}_T-1} \left[\mathbb{E}_{\widehat{M} \sim \xi^s} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^s}(r^s, o^s)}{\mathbb{P}^{M', \pi^s}(r^s, o^s)} \right] + \log \frac{\mathbb{P}^{M', \pi^s}(r^s, o^s)}{\mathbb{P}^{M'', \pi^s}(r^s, o^s)} \right] \cdot \mathbb{I}\{\tilde{s}_T \in [2^{j-1}, 2^j], \mathcal{E}_j\} \\ & \leq \tilde{s}_T^\alpha \log(T \log T) + \rho_j \cdot \tilde{s}_T^{1+\alpha}, \end{aligned}$$

where the inequality holds since on \mathcal{E}_j , we can ensure that $\log \frac{\mathbb{P}^{M', \pi^s}(r^s, o^s)}{\mathbb{P}^{M'', \pi^s}(r^s, o^s)} \leq \rho_j$. Therefore, choosing $\rho_j = 2^{-3j}$, we have

$$\begin{aligned} \text{Eq. (51)} & \leq \mathbb{E} \left[\left(\tilde{s}_T^\alpha \log(T \log T) + \rho_j \cdot \tilde{s}_T^{1+\alpha} + V_{\mathcal{M}} \tilde{s}_T^{1/2+\alpha} \sqrt{168 \tilde{s}_T \log(8 \mathbf{N}_{\text{cov}}(\mathcal{M}_{\text{alt}}^\star, \rho_j, \frac{\tilde{s}_T^4}{16})} \right) \right. \\ & \quad \left. + V_{\mathcal{M}} \cdot \left(4 \frac{\tilde{s}_T^{1/2+\alpha/2}}{1-\alpha} + \sum_{s=1}^{\tilde{s}_T} s^{1/2+\alpha/2} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D(\widehat{M}(\pi) \parallel M^\star(\pi))]] \right) \right) \cdot \mathbb{I}\{\tilde{s}_T \in [2^{j-1}, 2^j]\} \Big] \\ & \leq \mathbb{E} \left[\left(\tilde{s}_T^\alpha \log(T \log T) + 2^{-j} + V_{\mathcal{M}} \tilde{s}_T^{\alpha+1/2} \sqrt{168 \log(8 \tilde{s}_T \mathbf{N}_{\text{cov}}(\mathcal{M}_{\text{alt}}^\star, \frac{\tilde{s}_T^3}{8}, \frac{\tilde{s}_T^4}{16}))} \right) \right. \\ & \quad \left. + V_{\mathcal{M}} \cdot \left(4 \frac{\tilde{s}_T^{1/2+\alpha/2}}{1-\alpha} + \sum_{s=1}^{\tilde{s}_T} s^{1/2+\alpha/2} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D(\widehat{M}(\pi) \parallel M^\star(\pi))]] \right) \right) \cdot \mathbb{I}\{\tilde{s}_T \in [2^{j-1}, 2^j]\} \Big]. \end{aligned}$$

As this holds for each j , and the events $\{\tilde{s}_T \in [2^{j-1}, 2^j]\}$ are disjoint, we can sum over j to bound

$$\begin{aligned} & \sum_{j=1}^{\lceil \log T \rceil} \mathbb{E} \left[\tilde{s}_T^\alpha \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^\star} \sum_{s=1}^{\tilde{s}_T-1} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(\widehat{M}(\pi) \parallel M(\pi))]] \cdot \mathbb{I}\{\tilde{s}_T \in [2^{j-1}, 2^j], A_j \cap \mathcal{E}_j\} \right] \\ & \leq \mathbb{E} \left[\tilde{s}_T^\alpha \log(T \log T) + 1 + V_{\mathcal{M}} \tilde{s}_T^{\alpha+1/2} \sqrt{168 \log(8 \tilde{s}_T \mathbf{N}_{\text{cov}}(\mathcal{M}_{\text{alt}}^\star, \frac{\tilde{s}_T^3}{8}, \frac{\tilde{s}_T^4}{16}))} \right. \\ & \quad \left. + V_{\mathcal{M}} \cdot \left(4 \frac{\tilde{s}_T^{1/2+\alpha/2}}{1-\alpha} + \sum_{s=1}^{\tilde{s}_T} s^{1/2+\alpha/2} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D(\widehat{M}(\pi) \parallel M^\star(\pi))]] \right) \right]. \end{aligned}$$

Finally, under [Assumption A.3](#) and [Lemma F.15](#) we can bound

$$\begin{aligned} \log(8\tilde{s}_T \mathbf{N}_{\text{cov}}(\mathcal{M}_{\text{alt}}^*, \frac{\tilde{s}_T^3}{8}, \frac{\tilde{s}_T^4}{16})) &\leq \log(8\tilde{s}_T) + \log \mathbf{N}_{\text{cov}}(\mathcal{M}, \frac{\tilde{s}_T^3}{8}, \frac{\tilde{s}_T^4}{16}) \\ &\leq \log(8\tilde{s}_T) + d_{\text{cov}} \cdot \log(128C_{\text{cov}}\tilde{s}_T^7) \\ &\leq 8d_{\text{cov}} \cdot \log(128C_{\text{cov}}\tilde{s}_T). \end{aligned}$$

The result follows. \square

F.1.3. COMPLETING THE PROOF

Theorem F.1 (Full version of [Theorem A.1](#)). *For any $\delta \leq 1/2$, if we set $D(\cdot \| \cdot) = D_{\text{KL}}(\cdot \| \cdot)$ and instantiate \mathbf{Alg}_D with [Algorithm 7](#), under [Assumptions 1.1 to 1.3](#) and [A.1 to A.4](#) and if $\mathbf{g}^* > 0$, the expected regret of [Algorithm 5](#) is bounded by*

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq (1 + \delta)\mathbf{g}^* \cdot \log T + C_{\text{aec}} \cdot \overline{\text{aec}}_{\varepsilon/2}(\mathcal{M}) \cdot \log^{3/2}(\log T) + \text{lin}\left(C_{\text{low}}, \sqrt{\log T}, \log^{3/2}(\log T)\right),$$

for

$$\begin{aligned} C_{\text{low}} &:= \text{lin}\left(\mathbf{g}^*, \max_{M \in \mathcal{M}} \mathbf{g}^M, V_{\mathcal{M}}^{13/2}, L_{\text{KL}}^2, d_{\text{cov}}, \log C_{\text{cov}}, \frac{1}{\delta^2}, \frac{1}{\Delta_{\text{min}}^3}, n_{\delta/6}^{\mathcal{M}}, \sqrt{\text{aec}_{\varepsilon/2}(\mathcal{M})}\right) \\ C_{\text{aec}} &:= c \cdot \frac{V_{\mathcal{M}}^2 d_{\text{cov}} \log(C_{\text{cov}}) \cdot \max_{M \in \mathcal{M}} \mathbf{g}^M \cdot (\delta^{-1} + V_{\mathcal{M}} n_{\delta/6}^{\mathcal{M}})}{\delta \Delta_{\text{min}}^3} \cdot \log(C_{\text{low}}) \end{aligned}$$

and where $\text{lin}(\cdot)$ denotes a function linear and poly-logarithmic in its arguments and $c > 0$ is a universal constant.

Proof of Theorem F.1. Letting $\mathcal{T}_{\text{exploit}}$ denote the exploitation rounds, we can bound the total expected regret as

$$\sum_{t=1}^T \mathbb{E}[\Delta^*(\pi^t)] = \mathbb{E}\left[\sum_{t \in \mathcal{T}_{\text{exploit}}} \Delta^*(\pi^t)\right] + \mathbb{E}\left[\sum_{s=1}^{s_T} \Delta^*(\pi^s)\right] \leq 2 \log \log T + 3 + \mathbb{E}\left[\sum_{s=1}^{s_T} \Delta^*(\pi^s)\right],$$

where the inequality follows from [Lemma F.1](#). It remains to bound the regret in the exploration rounds. By [Lemma F.3](#), we can bound this as

$$\begin{aligned} \mathbb{E}\left[\sum_{s=1}^{s_T} \Delta^*(\pi^s)\right] &\leq \frac{8n_{\max} + 2\mathbf{g}^{\mathcal{M}}}{\delta \mathbf{g}^{\mathcal{M}}} \cdot \overline{\text{aec}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}) \cdot \mathbb{E}[\widehat{\mathbf{Est}}_{\text{D}}(s_T)] + \frac{2(1 + \delta)\mathbf{g}^*}{\Delta_{\text{min}}} \cdot \mathbb{E}[\mathbf{Est}_{\text{KL}}(s_T)] \\ &\quad + (1 + \delta)\mathbf{g}^* \cdot \mathbb{E}\left[\sum_{s=1}^{s_T} \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s}[D_{\text{KL}}(M^*(\pi) \| M(\pi))] \cdot \mathbb{I}\{\pi_* \in \boldsymbol{\pi}_{\widehat{M}^s}\}\right]. \end{aligned}$$

Applying [Lemma F.4](#) with $\alpha = 0$, we can bound

$$\mathbb{E}\left[\sum_{s=1}^{s_T} \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s}[D_{\text{KL}}(M^*(\pi) \| M(\pi))] \cdot \mathbb{I}\{\pi_* \in \boldsymbol{\pi}_{\widehat{M}^s}\}\right] \leq \log T + \mathbb{E}\left[V_{\mathcal{M}} \sqrt{1344d_{\text{cov}}s_T \cdot \log(128C_{\text{cov}}s_T)}\right]$$

$$\begin{aligned}
 & + (V_{\mathcal{M}} + L_{\text{KL}}) \left(4\sqrt{s_T} + \sum_{s=1}^{s_T} \sqrt{s} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D(\widehat{M}(\pi) \parallel M^*(\pi))]] \right) + \log \log T + 7V_{\mathcal{M}} \\
 & \leq \log T + V_{\mathcal{M}} \sqrt{1344d_{\text{cov}} \mathbb{E}[s_T] \cdot \log(128C_{\text{cov}} \mathbb{E}[s_T])} \\
 & + (V_{\mathcal{M}} + L_{\text{KL}}) \left(4\sqrt{\mathbb{E}[s_T]} + \mathbb{E} \left[\sum_{s=1}^{s_T} \sqrt{s} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D(\widehat{M}(\pi) \parallel M^*(\pi))]] \right] \right) + \log \log T + 7V_{\mathcal{M}}
 \end{aligned} \tag{52}$$

where the last inequality follows from [Lemma F.16](#)—applied with $\alpha = \beta = 1/2$, $a = 128C_{\text{cov}}$ —and the concavity of \sqrt{x} . Note that the condition of [Lemma F.16](#) is met here since we assume $C_{\text{cov}} \geq 1$. It follows from [Lemma F.7](#) that

$$\begin{aligned}
 & \mathbb{E} \left[\sum_{s=1}^{s_T} \sqrt{s} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D(\widehat{M}(\pi) \parallel M^*(\pi))]] \right] \\
 & \leq \mathbb{E} \left[(2 + 6V_{\mathcal{M}})(20d_{\text{cov}} \cdot \log(64C_{\text{cov}}s_T) + 1) \cdot 5\sqrt{2s_T \log(2s_T)} \right] \\
 & \quad + \mathbb{E}[32(1 + V_{\mathcal{M}}) \log(s_T)] + 8 \\
 & \leq (2 + 6V_{\mathcal{M}})(20d_{\text{cov}} \cdot \log(64C_{\text{cov}} \mathbb{E}[s_T]) + 1) \cdot 5\sqrt{2\mathbb{E}[s_T] \log(2\mathbb{E}[s_T])} \\
 & \quad + 32(1 + V_{\mathcal{M}}) \log \mathbb{E}[s_T] + 8 \\
 & = O\left(V_{\mathcal{M}}d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T]) \sqrt{\mathbb{E}[s_T] \log(\mathbb{E}[s_T])} + V_{\mathcal{M}} \sqrt{\mathbb{E}[s_T]}\right),
 \end{aligned} \tag{53}$$

where the last inequality follows from [Lemma F.16](#). Similarly, by [Lemma F.8](#), we can bound both $\mathbb{E}[\widehat{\mathbf{Est}}_{\text{D}}(s_T)]$ and $\mathbb{E}[\mathbf{Est}_{\text{KL}}(s_T)]$ as

$$\begin{aligned}
 \mathbb{E}[\widehat{\mathbf{Est}}_{\text{D}}(s_T)], \mathbb{E}[\mathbf{Est}_{\text{KL}}(s_T)] & \leq \mathbb{E} \left[(2 + 5V_{\mathcal{M}})(20d_{\text{cov}} \cdot \log(64C_{\text{cov}}s_T) + 1) \cdot \sqrt{\log(2s_T)} \right] \\
 & \quad + \mathbb{E}[32(1 + V_{\mathcal{M}}) \log(s_T)] + 8 \\
 & \leq (2 + 5V_{\mathcal{M}})(20d_{\text{cov}} \cdot \log(64C_{\text{cov}} \mathbb{E}[s_T]) + 1) \cdot \sqrt{\log(2\mathbb{E}[s_T])} \\
 & \quad + 32(1 + V_{\mathcal{M}}) \log(\mathbb{E}[s_T]) + 8 \\
 & = O\left(V_{\mathcal{M}}d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T]) \sqrt{\log(\mathbb{E}[s_T])} + V_{\mathcal{M}} \log(\mathbb{E}[s_T])\right),
 \end{aligned} \tag{54}$$

where the second inequality follows from [Lemma F.16](#) and Jensen's inequality, which lets us pass the expectation through. Together this gives a regret bound of

$$\begin{aligned}
 & \mathbb{E} \left[\sum_{s=1}^{s_T} \Delta^*(\pi^s) \right] \\
 & \leq (1 + \delta) \mathbf{g}^* \cdot \log T + (1 + \delta) \mathbf{g}^* V_{\mathcal{M}} \sqrt{1344d_{\text{cov}} \mathbb{E}[s_T] \cdot \log(128C_{\text{cov}} \mathbb{E}[s_T])} + (1 + \delta) \mathbf{g}^* (\log \log T + 7V_{\mathcal{M}}) \\
 & \quad + (1 + \delta) \mathbf{g}^* (V_{\mathcal{M}} + L_{\text{KL}}) \cdot O\left(\sqrt{\mathbb{E}[s_T]} + V_{\mathcal{M}}d_{\text{cov}} \cdot \log(C_{\text{cov}} \mathbb{E}[s_T]) \sqrt{\mathbb{E}[s_T] \log(\mathbb{E}[s_T])} + V_{\mathcal{M}} \log \mathbb{E}[s_T]\right) \\
 & \quad + \frac{8n_{\max} + 2\mathbf{g}^{\mathcal{M}}}{\delta \underline{\mathbf{g}}^{\mathcal{M}}} \cdot \overline{\text{aEC}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}) \cdot O\left(V_{\mathcal{M}}d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T]) \sqrt{\log(\mathbb{E}[s_T])} + V_{\mathcal{M}} \log(\mathbb{E}[s_T])\right)
 \end{aligned}$$

$$+ \frac{2(1+\delta)\mathbf{g}^{\star}}{\Delta_{\min}} \cdot O\left(V_{\mathcal{M}}d_{\text{cov}} \log(C_{\text{cov}}\mathbb{E}[s_T])\sqrt{\log(\mathbb{E}[s_T])} + V_{\mathcal{M}} \log(\mathbb{E}[s_T])\right)$$

By [Lemma F.3](#) we can bound the total number of exploration rounds by

$$\begin{aligned} \mathbb{E}[s_T] &\leq \frac{24n_{\max}^2 + 8n_{\max}\mathbf{g}^{\mathcal{M}}}{(\delta\mathbf{g}^{\mathcal{M}})^2} \cdot \overline{\text{aEC}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}) \cdot \mathbb{E}[\widehat{\text{Est}}_{\text{D}}(s_T)] + \frac{12n_{\max}}{\delta\Delta_{\min}} \cdot \mathbb{E}[\text{Est}_{\text{KL}}(s_T)] \\ &\quad + \frac{6n_{\max}}{\delta} \cdot \mathbb{E}\left[\sum_{s=1}^{s_T} \inf_{M \in \mathcal{M}_{\text{alt}}^{\star}} \mathbb{E}_{p^s}[D_{\text{KL}}(M^{\star}(\pi) \parallel M(\pi))] \cdot \mathbb{I}\{\pi_{\star} \in \pi_{\widehat{M}^s}\}\right] \\ &\leq \frac{24n_{\max}^2 + 8n_{\max}\mathbf{g}^{\mathcal{M}}}{(\delta\mathbf{g}^{\mathcal{M}})^2} \cdot \overline{\text{aEC}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}) \cdot O\left(V_{\mathcal{M}}d_{\text{cov}} \log(C_{\text{cov}}\mathbb{E}[s_T])\sqrt{\log(\mathbb{E}[s_T])} + V_{\mathcal{M}} \log(\mathbb{E}[s_T])\right) \\ &\quad + \frac{12n_{\max}}{\delta\Delta_{\min}} \cdot O\left(V_{\mathcal{M}}d_{\text{cov}} \log(C_{\text{cov}}\mathbb{E}[s_T])\sqrt{\log(\mathbb{E}[s_T])} + V_{\mathcal{M}} \log(\mathbb{E}[s_T])\right) \\ &\quad + \frac{6n_{\max}}{\delta} \cdot (\log T + V_{\mathcal{M}}\sqrt{1344d_{\text{cov}}\mathbb{E}[s_T]} \cdot \log(128C_{\text{cov}}\mathbb{E}[s_T]) + \log \log T + 7V_{\mathcal{M}}) \\ &\quad + \frac{6n_{\max}}{\delta}(V_{\mathcal{M}} + L_{\text{KL}})\left(4\sqrt{\mathbb{E}[s_T]} + O\left(V_{\mathcal{M}}d_{\text{cov}} \log(C_{\text{cov}}\mathbb{E}[s_T])\sqrt{\mathbb{E}[s_T] \log(\mathbb{E}[s_T])} + V_{\mathcal{M}}\sqrt{\mathbb{E}[s_T]}\right)\right) \end{aligned}$$

where the second inequality follows from [Equations \(52\) to \(54\)](#). Using [Lemma F.17](#) and bounding $\mathbf{g}^{\mathcal{M}} \leq n_{\max}$, we can solve this for $\mathbb{E}[s_T]$ to get

$$\begin{aligned} \mathbb{E}[s_T] &\leq \tilde{O}\left(\frac{n_{\max}}{\delta} \cdot \log T + \frac{V_{\mathcal{M}}d_{\text{cov}} \log C_{\text{cov}} \cdot n_{\max}^2}{(\delta\mathbf{g}^{\mathcal{M}})^2} \cdot \overline{\text{aEC}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}) + \frac{n_{\max}V_{\mathcal{M}}d_{\text{cov}} \log C_{\text{cov}}}{\delta\Delta_{\min}}\right. \\ &\quad \left.+ \frac{n_{\max}^2(V_{\mathcal{M}} + L_{\text{KL}})^2(V_{\mathcal{M}}^2d_{\text{cov}}^2 \log^2 C_{\text{cov}} + V_{\mathcal{M}}^2)}{\delta^2}\right). \end{aligned}$$

Finally, by [Lemma F.13](#) and [Lemma E.3](#) we can lower bound $\mathbf{g}^{\mathcal{M}} \geq \Delta_{\min}/2V_{\mathcal{M}}$. Plugging this into the above expression and using that

$$n_{\max} = n_{\max}(\mathcal{M}, \delta/6) := \frac{32}{\Delta_{\min}^2} \cdot \left(\frac{6}{\delta} + 2V_{\mathcal{M}}n_{\delta/6}^{\mathcal{M}}\right) \cdot \max_{M \in \mathcal{M}} \mathbf{g}^M,$$

gives the result, after simplifying. □

F.2. Regret Bound without Uniform Regularity (Theorem A.2)

In this section, we prove [Theorem A.2](#) (as well as a more general result, [Theorem F.2](#)), which gives regret bounds for a variant of AE^2 , [Algorithm 6](#), AE_+^2 , which does not require uniform regularity, and adapts to the gap $\Delta_{\min}^{\star} := \Delta_{\min}^{M^{\star}}$ for the underlying model M^{\star} . [Algorithm 6](#) is a slightly more general version of [Algorithm 3](#), incorporating general divergences D .

Throughout this section, we let s^{\star} denote the first exploration round s such that $M^{\star} \in \mathcal{M}^{\ell}$ in [Algorithm 3](#). Note that s^{\star} is a deterministic quantity (though, the first round t in which $s = s^{\star}$ is not deterministic).

We remark briefly that, to bound Eq. (55) by a restriction of the AEC, it is critical that our estimator produced at each exploration round s , ξ^s , is only supported on \mathcal{M}^ℓ . To accomplish this, we explicitly generate a cover of \mathcal{M}^ℓ , $\mathcal{M}_{\text{cov}}^\ell$, and run an estimation procedure on this cover. As Algorithm 7, the estimation procedure used to prove Theorem A.1, directly covers all of \mathcal{M} , to prove Theorem A.2 we do not appeal directly to this algorithm, yet we note that the covering and estimation procedure employed by AE_*^2 are essentially identical to that in Algorithm 7, modulo the choice of which set is being covered.

Algorithm 6 Adaptive Exploration for Allocation Estimation for classes without uniform regularity (AE_*^2)

1: **input:** Optimality tolerance δ , divergence D , estimation oracle Alg_D , growth parameters α_q , $\alpha_n, \alpha_{\mathcal{M}}$.

2: $s \leftarrow 1, \ell \leftarrow 1, \varepsilon \leftarrow \frac{\delta}{4+2\delta}$.

3: $q^s \leftarrow 1 - s^{-\alpha_q}, n^s \leftarrow s^{\alpha_n}$.

4: Compute $\xi^1 \leftarrow \text{Alg}_D(\{\emptyset\})$ and $\widehat{M}^1 \leftarrow \mathbb{E}_{M \sim \xi^1}[M]$.

5: **for** $t = 1, 2, 3, \dots$ **do**

6: **if** $s \geq 2^\ell$ **then** // Form active set and cover

7: $\ell \leftarrow \ell + 1$.

8: $\Delta^\ell \leftarrow \arg \min_{\Delta \geq 0} \Delta \quad \text{s.t.} \quad \overline{\text{aec}}_{\varepsilon/2}^{\text{D}, \mathcal{M}}(\mathcal{M}_{\Delta, \frac{1}{\Delta}}) \leq s^{\alpha_{\mathcal{M}}}$.

9: $\mathcal{M}^\ell \leftarrow \mathcal{M}_{\Delta^\ell, \frac{1}{\Delta^\ell}} \cap \{M \in \mathcal{M} : n_\varepsilon^M + \frac{1}{\Delta_{\min}^M} + \frac{4g^M}{\Delta_{\min}^M} + \frac{2n_\varepsilon^M}{g^M} + \frac{4}{\Delta_{\min}^M \varepsilon} \leq \sqrt{n^s}\}$.

10: $\mathcal{M}_{\text{cov}}^\ell \leftarrow (\rho_\ell, \mu_\ell)$ -cover of \mathcal{M}^ℓ for $\rho_\ell \leftarrow 2^{-\ell}, \mu_\ell \leftarrow 2^{-5\ell}, \mathfrak{D}^\ell \leftarrow \emptyset$.

11: **if** $\exists \pi_{\widehat{M}^s} \in \pi_{\widehat{M}^s}$ s.t. $\forall M \in \mathcal{M}^{\text{alt}}(\pi_{\widehat{M}^s}), \sum_{i=1}^{s-1} \mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}_{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}_{M, \pi^i}(r^i, o^i)} \right] \geq \log(t \log t)$ **then**

12: Play $\pi_{\widehat{M}^s}$. // Exploit

13: **else** // Explore

14: Set $p^s \leftarrow q^s \lambda^s + (1 - q^s) \omega^s$ for

$$\lambda^s, \omega^s \leftarrow \arg \min_{\lambda, \omega \in \Delta_\Pi} \sup_{M \in \mathcal{M}^\ell \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda; n^s)} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim \omega} [D(\widehat{M}(\pi) \| M(\pi))]]}. \quad (55)$$

15: Draw $\pi^s \sim p^s$, observe r^s, o^s , set $\mathfrak{D}^\ell \leftarrow \mathfrak{D}^\ell \cup \{(\pi^s, r^s, o^s)\}$.

16: Compute estimate $\xi^{s+1} \leftarrow \text{Alg}_D(\mathfrak{D}^\ell, \mathcal{M}_{\text{cov}}^\ell)$ and $\widehat{M}^{s+1} = \mathbb{E}_{M \sim \xi^{s+1}}[M]$.

17: $s \leftarrow s + 1$.

F.2.1. BOUNDING REGRET OF EXPLORE PHASE

Lemma F.5. *We have*

$$s^* \leq \left(\overline{\text{aec}}_{\varepsilon/2}^{\text{D}, \mathcal{M}}(\mathcal{M}^*) \right)^{\frac{1}{\alpha_{\mathcal{M}}}} + \left(n_\varepsilon^* + \frac{1}{\Delta_{\min}^*} + \frac{4g^*}{\Delta_{\min}^*} + \frac{2n_\varepsilon^*}{g^*} + \frac{4}{\Delta_{\min}^* \varepsilon} \right)^{\frac{2}{\alpha_n}}.$$

Proof of Lemma F.5. We will have $M^* \in \mathcal{M}^\ell$ as soon as $M^* \in \mathcal{M}'$ and $M^* \in \mathcal{M}''$ for

$$\mathcal{M}' \leftarrow \left\{ M \in \mathcal{M} : n_\varepsilon^M + \frac{1}{\Delta_{\min}^M} + \frac{4g^M}{\Delta_{\min}^M} + \frac{2n_\varepsilon^M}{g^M} + \frac{4}{\Delta_{\min}^M \varepsilon} \leq \sqrt{n^{2^\ell-1}} \right\}, \quad \mathcal{M}'' \leftarrow \mathcal{M}_{\Delta^\ell, \frac{1}{\Delta^\ell}}$$

where

$$\Delta^\ell = \arg \min_{\Delta \geq 0} \Delta \quad \text{s.t.} \quad \overline{\text{aec}}_{\varepsilon/2}^{\text{D}, \mathcal{M}}(\mathcal{M}_{\Delta, \frac{1}{\Delta}}) \leq s^{\alpha_{\mathcal{M}}}.$$

Note that if this occurs for some ℓ , then the first exploration round s such that $M^* \in \mathcal{M}^\ell$ is $s = 2^{\ell-1}$. From the definition $n^{2^{\ell-1}} = 2^{\alpha_n(\ell-1)}$, we have $M^* \in \mathcal{M}'$ as soon as

$$\begin{aligned} \sqrt{2^{\alpha_n(\ell-1)}} &\geq n_\varepsilon^* + \frac{1}{\Delta_{\min}^*} + \frac{4\mathbf{g}^*}{\Delta_{\min}^*} + \frac{2n_\varepsilon^*}{\mathbf{g}^*} + \frac{4}{\Delta_{\min}^* \varepsilon} \\ \iff 2^{\ell-1} &\geq \left(n_\varepsilon^* + \frac{1}{\Delta_{\min}^*} + \frac{4\mathbf{g}^*}{\Delta_{\min}^*} + \frac{2n_\varepsilon^*}{\mathbf{g}^*} + \frac{4}{\Delta_{\min}^* \varepsilon} \right)^{2/\alpha_n}. \end{aligned}$$

To have $M^* \in \mathcal{M}''$, we need $\Delta^\ell \leq \Delta_{\min}^*$ and $n_\varepsilon^* \leq \frac{1}{\Delta^\ell}$, which will be the case once

$$\overline{\text{aec}}_{\varepsilon/2}^{\text{D}, \mathcal{M}}(\mathcal{M}^*) \leq 2^{(\ell-1) \cdot \alpha_{\mathcal{M}}} \iff \left(\overline{\text{aec}}_{\varepsilon/2}^{\text{D}, \mathcal{M}}(\mathcal{M}^*) \right)^{\frac{1}{\alpha_{\mathcal{M}}}} \leq 2^{\ell-1}.$$

The result then follows by combining these bounds. \square

Lemma F.6. *Let s_T denote the total number of exploration rounds. For $\delta \leq 1/2$, running [Algorithm 6](#), we can almost surely bound bound*

$$\begin{aligned} s_T &\leq \frac{4}{\delta \mathbf{g}^*} \left(\sum_{s=s^*}^{s_T} 2s^{\alpha_q + \alpha_{\mathcal{M}}} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]] \right. \\ &\quad + \sum_{s=s^*}^{s_T} s^{\alpha_n} \cdot 2\mathbf{g}^* \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \cdot \mathbb{I}\{\pi_* \in \pi_{\widehat{M}^s}\} \\ &\quad \left. + \sum_{s=s^*}^{s_T} s^{3\alpha_n/2} \cdot 4\mathbf{g}^* \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi))]] + 4\mathbf{g}^* \left(\frac{8}{\delta \mathbf{g}^*} \right)^{\frac{1+\alpha_n}{\alpha_q - \alpha_n}} \right) + s^*. \end{aligned}$$

In addition, the regret during exploration rounds is bounded as

$$\begin{aligned} \mathbb{E} \left[\sum_{s=s^*}^{s_T} \Delta^*(\pi^s) \right] &\leq \mathbb{E} \left[\sum_{s=s^*}^{s_T} 2s^{\alpha_q + \alpha_{\mathcal{M}}} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]] \right. \\ &\quad + \sum_{s=s^*}^{s_T} (1 + \delta) \mathbf{g}^* \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \cdot \mathbb{I}\{\pi_* \in \pi_{\widehat{M}^s}\} \\ &\quad \left. + \sum_{s=s^*}^{s_T} s^{\alpha_n/2} \cdot 4\mathbf{g}^* \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi))]] \right] + 4\mathbf{g}^* \left(\frac{8}{\delta \mathbf{g}^*} \right)^{\frac{1}{\alpha_q - \alpha_n}}. \end{aligned}$$

Proof of Lemma F.6. We first prove the bound on the regret during the exploration rounds, then use this result to prove the bound on s_T .

Recall that by definition, for exploration rounds, we have $p^s \leftarrow q^s \lambda^s + (1 - q^s) \omega^s$. We consider three cases to bound the instantaneous expected regret, $\Delta^*(p^s)$, for each $s \geq s^*$.

Case 1: $M^* \in \mathcal{M}^\ell \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda^s; n^s)$. Denote such rounds as $\mathcal{S}_{\text{exp}}^1$. This case follows identically to Case 1 in [Lemma F.3](#) with

$$\gamma^s = \frac{1 + \mathbf{g}^* \delta}{1 + q^s} \cdot \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\omega^s} [D(\widehat{M}(\pi) \| M^*(\pi))]]}. \quad (56)$$

We therefore omit the proof and conclude that

$$\begin{aligned} \Delta^*(p^s) &\leq -\mathbf{g}^* \delta + \gamma^s \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]] \\ &\leq \gamma^s \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]]. \end{aligned}$$

As regret is always lower-bounded by 0, we have $\Delta^*(p^s) \geq 0$, so for rounds $s \in \mathcal{S}_{\text{exp}}^1$, we can also write

$$\mathbf{g}^* \delta \leq \gamma^s \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]]. \quad (57)$$

Case 2: $M^* \in \mathcal{M}_\varepsilon^{\text{gl}}(\lambda^s; n^s)$, $\pi_* \in \pi_{\widehat{M}^s}$. Denote such rounds as $\mathcal{S}_{\text{exp}}^2$, and write

$$\Delta^*(p^s) = [\Delta^*(p^s) - (1 + \delta) \mathbf{g}^* \cdot \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] + (1 + \delta) \mathbf{g}^* \cdot \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))]]$$

for an arbitrary model $M \in \mathcal{M}_{\text{alt}}^*$. In this case, since $M^* \in \mathcal{M}_\varepsilon^{\text{gl}}(\lambda^s; n^s)$, we have that $\lambda^s \in \Lambda(M^*; \varepsilon)$. This implies that

$$\Delta^*(\lambda^s) \leq (1 + \varepsilon) \mathbf{g}^* / n_s^* \quad \text{and} \quad \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{\lambda^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \geq (1 - \varepsilon) / n_s^*$$

for some $n_s^* \leq n^s$.

Since $M \in \mathcal{M}_{\text{alt}}^*$, it follows that $\mathbb{E}_{\lambda^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \geq (1 - \varepsilon) / n_s^*$. Thus,

$$\begin{aligned} \Delta^*(p^s) - (1 + \delta) \mathbf{g}^* \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] &\leq q^s [\Delta^*(\lambda^s) - (1 + \delta) \mathbf{g}^* \mathbb{E}_{\lambda^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))]] + 1 - q^s \\ &\leq q^s [(1 + \varepsilon) \mathbf{g}^* / n_s^* - (1 + \delta)(1 - \varepsilon) \mathbf{g}^* / n_s^*] + 1 - q^s \\ &= q^s [2\varepsilon - \delta(1 - \varepsilon)] \cdot \frac{\mathbf{g}^*}{n_s^*} + 1 - q^s \\ &\stackrel{(a)}{=} -\frac{(1 - s^{-\alpha_q}) \delta}{2} \cdot \frac{\mathbf{g}^*}{n_s^*} + s^{-\alpha_q} \\ &\stackrel{(b)}{\leq} \begin{cases} s^{-\alpha_q} & s < (\frac{8n_s^*}{\delta \mathbf{g}^*})^{1/\alpha_q} \\ -\frac{\delta \mathbf{g}^*}{4n_s^*} & s \geq (\frac{8n_s^*}{\delta \mathbf{g}^*})^{1/\alpha_q} \end{cases} \\ &\stackrel{(c)}{\leq} \begin{cases} s^{-\alpha_q} & s < (\frac{8}{\delta \mathbf{g}^*})^{\frac{1}{\alpha_q - \alpha_n}} \\ -\frac{1}{4} \delta \mathbf{g}^* \cdot s^{-\alpha_n} & s \geq (\frac{8}{\delta \mathbf{g}^*})^{\frac{1}{\alpha_q - \alpha_n}} \end{cases} \\ &\stackrel{(d)}{\leq} 2\mathbf{g}^* \mathbb{I}\{s < (\frac{8}{\delta \mathbf{g}^*})^{\frac{1}{\alpha_q - \alpha_n}}\} - \frac{1}{4} \delta \mathbf{g}^* \cdot s^{-\alpha_n} \end{aligned}$$

where (a) follows from our choice of $\varepsilon = \frac{\delta/2}{2+\delta}$ and $q^s = 1 - s^{-\alpha_q}$, (b) follows from some algebra, (c) uses that $n_s^* \leq n^s$ and $n^s = s^{\alpha_n}$, and (d) follows since we will always have $s^{-\alpha_q} \leq 2\mathbf{g}^* - \frac{1}{4} \delta \mathbf{g}^* \cdot s^{-\alpha_n}$. Thus:

$$\Delta^*(p^s) \leq (1 + \delta) \mathbf{g}^* \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] + 2\mathbf{g}^* \mathbb{I}\{s < (\frac{8}{\delta \mathbf{g}^*})^{\frac{1}{\alpha_q - \alpha_n}}\} - \frac{1}{4} \delta \mathbf{g}^* \cdot s^{-\alpha_n}$$

$$\leq (1 + \delta)g^* \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] + 2g^* \mathbb{I}\{s < (\frac{8}{\delta g^*})^{\frac{1}{\alpha q - \alpha n}}\}.$$

As this holds for all $M \in \mathcal{M}_{\text{alt}}^*$, we have

$$\Delta^*(p^s) \leq (1 + \delta)g^* \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] + 2g^* \mathbb{I}\{s < (\frac{8}{\delta g^*})^{\frac{1}{\alpha q - \alpha n}}\}.$$

Since $\Delta^*(p^s) \geq 0$, this also implies that for $s \in \mathcal{S}_{\text{exp}}^2$:

$$\begin{aligned} \frac{1}{4}\delta g^* &\leq s^{\alpha n} \cdot 2g^* \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] + 2g^* s^{\alpha n} \mathbb{I}\{s < (\frac{8}{\delta g^*})^{\frac{1}{\alpha q - \alpha n}}\} \\ &\leq s^{\alpha n} \cdot 2g^* \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] + 2g^* (\frac{8}{\delta g^*})^{\frac{\alpha n}{\alpha q - \alpha n}} \cdot \mathbb{I}\{s < (\frac{8}{\delta g^*})^{\frac{1}{\alpha q - \alpha n}}\}. \end{aligned} \tag{58}$$

Case 3: $M^* \in \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda^s; n^s)$, $\pi_* \notin \pi_{\widehat{M}^s}$. Denote such rounds as $\mathcal{S}_{\text{exp}}^3$, and write

$$\begin{aligned} \Delta^*(p^s) &= \left[\Delta^*(p^s) - 2(1 + \delta)g^* \sqrt{n^s} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi))] \right] \\ &\quad + 2(1 + \delta)g^* \sqrt{n^s} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi))]. \end{aligned}$$

Since $M^* \in \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda^s)$, we have that $\lambda^s \in \Lambda(M^*; \varepsilon)$. This implies that for any $M \in \mathcal{M}_{\text{alt}}^*$:

$$\Delta^*(\lambda^s) \leq (1 + \varepsilon)g^*/n_s^* \quad \text{and} \quad \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{\lambda^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \geq (1 - \varepsilon)/n_s^*$$

for some $n_s^* \leq n^s$. Assume that we are at epoch ℓ . By construction we have that, for $M \in \mathcal{M}^\ell$, $\frac{1}{\Delta_{\min}^M} \leq \sqrt{n^{2\ell}} \iff \Delta_{\min}^M \geq \frac{1}{\sqrt{n^{2\ell}}}$. Since n^s is increasing in s , this implies that $\Delta_{\min}^M \geq \frac{1}{\sqrt{n^s}}$. As ξ^s is only supported on \mathcal{M}^ℓ , since $\pi_* \notin \pi_{\widehat{M}^s}$, [Lemma F.9](#) implies that $\mathbb{E}_{M \sim \xi^s} [\mathbb{I}\{M \in \mathcal{M}_{\text{alt}}^*\}] \geq \frac{1}{2\sqrt{n^s}}$. Thus, we have

$$\begin{aligned} &2(1 + \delta)g^* \sqrt{n^s} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\lambda^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi))] \\ &\geq 2(1 + \delta)g^* \sqrt{n^s} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\lambda^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi)) \cdot \mathbb{I}\{M \in \mathcal{M}_{\text{alt}}^*\}]] \\ &\geq 2(1 + \delta)g^* \sqrt{n^s} \cdot \frac{1 - \varepsilon}{2\sqrt{n^s} n_s^*} \\ &\geq \frac{(1 + \delta)(1 - \varepsilon)g^*}{n_s^*}. \end{aligned}$$

This implies that

$$\begin{aligned} &\Delta^*(p^s) - 2(1 + \delta)g^* \sqrt{n^s} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi))] \\ &\leq q^s [(1 + \varepsilon)g^*/n_s^* - (1 + \delta)(1 - \varepsilon)g^*/n_s^*] + 1 - q^s \\ &\leq 2g^* \mathbb{I}\{s < (\frac{8}{\delta g^*})^{\frac{1}{\alpha q - \alpha n}}\} - \frac{1}{4}\delta g^* \cdot s^{-\alpha n}, \end{aligned}$$

where the final inequality follows by the same argument as in Case 2. Thus,

$$\Delta^*(p^s) \leq 2(1 + \delta)g^* \sqrt{n^s} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi))] + 2g^* \mathbb{I}\{s < (\frac{8}{\delta g^*})^{\frac{1}{\alpha q - \alpha n}}\} - \frac{1}{4}\delta g^* \cdot s^{-\alpha n}$$

$$\begin{aligned}
 &= s^{\alpha_n/2} \cdot 2(1 + \delta) \mathbf{g}^* \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel \widehat{M}(\pi))]] + 2\mathbf{g}^* \mathbb{I}\{s < (\frac{8}{\delta \mathbf{g}^*})^{\frac{1}{\alpha_q - \alpha_n}}\} - \frac{1}{4} \delta \mathbf{g}^* \cdot s^{-\alpha_n} \\
 &\leq s^{\alpha_n/2} \cdot 2(1 + \delta) \mathbf{g}^* \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel \widehat{M}(\pi))]] + 2\mathbf{g}^* \mathbb{I}\{s < (\frac{8}{\delta \mathbf{g}^*})^{\frac{1}{\alpha_q - \alpha_n}}\}.
 \end{aligned}$$

Just as in Case 2, using that $\Delta^*(p^s) \geq 0$, this also implies that for $s \in \mathcal{S}_{\text{exp}}^3$:

$$\frac{1}{4} \delta \mathbf{g}^* \leq s^{3\alpha_n/2} \cdot 2(1 + \delta) \mathbf{g}^* \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel \widehat{M}(\pi))]] + 2\mathbf{g}^* (\frac{8}{\delta \mathbf{g}^*})^{\frac{\alpha_n}{\alpha_q - \alpha_n}} \cdot \mathbb{I}\{s < (\frac{8}{\delta \mathbf{g}^*})^{\frac{1}{\alpha_q - \alpha_n}}\}. \quad (59)$$

Completing the Proof. In total we have

$$\begin{aligned}
 \sum_{s=s^*}^{s_T} \Delta^*(p^s) &\leq \sum_{s \in \mathcal{S}_{\text{exp}}^1} \gamma^s \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \parallel M^*(\pi))]] + \sum_{s \in \mathcal{S}_{\text{exp}}^2} (1 + \delta) \mathbf{g}^* \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \\
 &\quad + \sum_{s \in \mathcal{S}_{\text{exp}}^3} s^{\alpha_n/2} \cdot 2(1 + \delta) \mathbf{g}^* \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel \widehat{M}(\pi))]] + 4\mathbf{g}^* (\frac{8}{\delta \mathbf{g}^*})^{\frac{1}{\alpha_q - \alpha_n}}.
 \end{aligned}$$

By [Lemma F.11](#), for $s \in \mathcal{S}_{\text{exp}}^1$, we can bound $\gamma^s \leq (1 + \mathbf{g}^* \delta) \cdot s^{\alpha_q + \alpha_{\mathcal{M}}}$. This gives an upper bound on the above of

$$\begin{aligned}
 &\leq \sum_{s=s^*}^{s_T} (1 + \mathbf{g}^*) s^{\alpha_q + \alpha_{\mathcal{M}}} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \parallel M^*(\pi))]] \\
 &\quad + \sum_{s=s^*}^{s_T} (1 + \delta) \mathbf{g}^* \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \cdot \mathbb{I}\{\pi_\star \in \boldsymbol{\pi}_{\widehat{M}^s}\} \\
 &\quad + \sum_{s=s^*}^{s_T} s^{\alpha_n/2} \cdot 4\mathbf{g}^* \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel \widehat{M}(\pi))]] + 4\mathbf{g}^* (\frac{8}{\delta \mathbf{g}^*})^{\frac{1}{\alpha_q - \alpha_n}},
 \end{aligned}$$

which proves the regret bound.

We now bound the number of exploration rounds. Since for every $s \in \mathcal{S}_{\text{exp}}^1$ [Eq. \(57\)](#) holds, for every $s \in \mathcal{S}_{\text{exp}}^2$ [Eq. \(58\)](#) holds, and for every $s \in \mathcal{S}_{\text{exp}}^3$ [Eq. \(59\)](#) holds, combining these inequalities gives

$$\begin{aligned}
 &\frac{1}{4} \delta \mathbf{g}^* |\mathcal{S}_{\text{exp}}^1 \cup \mathcal{S}_{\text{exp}}^2 \cup \mathcal{S}_{\text{exp}}^3| \\
 &\leq \sum_{s \in \mathcal{S}_{\text{exp}}^1} \gamma^s \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \parallel M^*(\pi))]] + \sum_{s \in \mathcal{S}_{\text{exp}}^2} s^{\alpha_n} \cdot 2\mathbf{g}^* \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \\
 &\quad + \sum_{s \in \mathcal{S}_{\text{exp}}^3} s^{3\alpha_n/2} \cdot 4\mathbf{g}^* \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel \widehat{M}(\pi))]] + 4\mathbf{g}^* (\frac{8}{\delta \mathbf{g}^*})^{\frac{1 + \alpha_n}{\alpha_q - \alpha_n}} \\
 &\leq \sum_{s=s^*}^{s_T} (1 + \mathbf{g}^*) s^{\alpha_q + \alpha_{\mathcal{M}}} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \parallel M^*(\pi))]] \\
 &\quad + \sum_{s=s^*}^{s_T} s^{\alpha_n} \cdot 2\mathbf{g}^* \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \cdot \mathbb{I}\{\pi_\star \in \boldsymbol{\pi}_{\widehat{M}^s}\}
 \end{aligned}$$

$$+ \sum_{s=s^*}^{s_T} s^{3\alpha_n/2} \cdot 4\mathbf{g}^* \mathbb{E}_{\widehat{M} \sim \xi_s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel \widehat{M}(\pi))]] + 4\mathbf{g}^* \left(\frac{8}{\delta \mathbf{g}^*}\right)^{\frac{1+\alpha_n}{\alpha_q - \alpha_n}}.$$

Using that $|\mathcal{S}_{\text{exp}}^1 \cup \mathcal{S}_{\text{exp}}^2 \cup \mathcal{S}_{\text{exp}}^3| = s_T - s^*$ and rearranging gives the bound on s_T . \square

F.2.2. COMPLETING THE PROOF

Theorem F.2 (Full version of [Theorem A.2](#)). *Consider [Algorithm 6](#), and suppose we set $\delta \leq 1/2$, $D(\cdot \parallel \cdot) = D_{\text{KL}}(\cdot \parallel \cdot)$, $\alpha_{\mathcal{M}} = 1/2$, $\alpha_{\zeta} = 1/8$, and $\alpha_q = 1/4$, and instantiate \mathbf{Alg}_D with [Algorithm 4](#). Then if [Assumptions 1.1 to 1.3](#) and [A.1 to A.3](#) hold and $\mathbf{g}^* > 0$, the expected regret of is bounded by*

$$\mathbb{E}^{M^*} [\mathbf{Reg}(T)] \leq (1 + \delta) \mathbf{g}^* \log T + C_{\text{aec}} \cdot \left(\overline{\text{aec}}_{\varepsilon/2}^{\mathcal{D}, \mathcal{M}}(\mathcal{M}^*) \right)^3 \cdot \log^{3/2} \log T + C_{\text{low}} \cdot \log^{6/7} T$$

for $\varepsilon \leftarrow \frac{\delta}{4+2\delta}$,

$$C_{\text{aec}} := \tilde{O} \left(\frac{(V_{\mathcal{M}} + L_{\text{KL}}) V_{\mathcal{M}}^3 d_{\text{cov}} \log(C_{\text{cov}})}{\delta \Delta_{\min}^*} \right),$$

and

$$C_{\text{low}} := \tilde{O} \left(\text{poly} \left(V_{\mathcal{M}}, L_{\text{KL}}, n_{\varepsilon}^*, d_{\text{cov}}, \log C_{\text{cov}}, \mathbf{g}^*, \frac{1}{\Delta_{\min}^*}, \frac{1}{\delta}, \log \log T \right) \right).$$

Proof of Theorem F.2. The bound on the regret incurred in the exploit phase follows identically to [Theorem F.1](#), since the exploit test is the same. We turn to bounding the regret in the explore phase. First, since we can incur regret of at most 1 at every round, we bound

$$\mathbb{E} \left[\sum_{s=1}^{s_T} \Delta^{M^*}(p^s) \right] \leq \mathbb{E} \left[\sum_{s=s^*}^{s_T} \Delta^{M^*}(p^s) \right] + s^*$$

By [Lemma F.6](#), we can bound

$$\begin{aligned} \mathbb{E} \left[\sum_{s=s^*}^{s_T} \Delta^*(\pi^s) \right] &\leq \mathbb{E} \left[\sum_{s=s^*}^{s_T} (1 + \mathbf{g}^*) s^{\alpha_q + \alpha_{\mathcal{M}}} \mathbb{E}_{\widehat{M} \sim \xi_s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \parallel M^*(\pi))]] \right. \\ &\quad + \sum_{s=s^*}^{s_T} (1 + \delta) \mathbf{g}^* \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \cdot \mathbb{I}\{\pi_{\star} \in \pi_{\widehat{M}^s}\} \\ &\quad \left. + \sum_{s=s^*}^{s_T} s^{\alpha_n/2} \cdot 4\mathbf{g}^* \mathbb{E}_{\widehat{M} \sim \xi_s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel \widehat{M}(\pi))]] \right] + 4\mathbf{g}^* \left(\frac{8}{\delta \mathbf{g}^*}\right)^{\frac{1}{\alpha_q - \alpha_n}} \\ &\leq \mathbb{E} \left[\sum_{s=s^*}^{s_T} (1 + \mathbf{g}^*) s^{\alpha_q + \alpha_{\mathcal{M}}} \mathbb{E}_{\widehat{M} \sim \xi_s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \parallel M^*(\pi))]] \right. \\ &\quad \left. + \sum_{s=1}^{s_T} (1 + \delta) \mathbf{g}^* \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \cdot \mathbb{I}\{\pi_{\star} \in \pi_{\widehat{M}^s}\} \right] \end{aligned}$$

$$+ \sum_{s=s^*}^{s_T} s^{\alpha_n/2} \cdot 4\mathbf{g}^* \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel \widehat{M}(\pi))]] \Big] + 4\mathbf{g}^* \left(\frac{8}{\delta \mathbf{g}^*}\right)^{\frac{1}{\alpha_q - \alpha_n}}.$$

Applying [Lemma F.4](#) with $\alpha = 0$, we have

$$\begin{aligned} \mathbb{E} \left[\sum_{s=1}^{s_T} \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \cdot \mathbb{I}\{\pi_* \in \pi_{\widehat{M}^s}\} \right] &\leq \log T + \mathbb{E} \left[V_{\mathcal{M}} s_T^{1/2} \cdot \sqrt{1344 d_{\text{cov}} \cdot \log(128 C_{\text{cov}} s_T)} \right] \\ &+ (V_{\mathcal{M}} + L_{\text{KL}}) \left(4s_T^{1/2} + \sum_{s=1}^{s_T} s^{1/2} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D(\widehat{M}(\pi) \parallel M^*(\pi))] \right] \Big] + \log \log T + 7V_{\mathcal{M}}. \end{aligned}$$

Note that for $s \geq s^*$, our estimator is applied to a cover of a set containing M^* . Furthermore, note that the estimation procedure used in [Algorithm 6](#) is, other than the different choice of set to cover, identical to that used in [Algorithm 7](#). It follows that the analysis of [Algorithm 7](#) can be applied to the estimation procedure of [Algorithm 6](#), with only the mild modification accounting for the difference in the size of the cover (as we are covering \mathcal{M}^ℓ instead of \mathcal{M}). However, as we can upper bound the size of the cover of \mathcal{M}^ℓ by the size of the cover of \mathcal{M} via [Lemma F.15](#), this change is inconsequential. Thus, by [Lemma F.7](#), we can bound

$$\begin{aligned} \mathbb{E} \left[\sum_{s=s^*}^{s_T} 2s^{\alpha_q + \alpha_{\mathcal{M}}} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \parallel M^*(\pi))] \right] &\leq O \left(\mathbb{E} \left[V_{\mathcal{M}} d_{\text{cov}} \log(C_{\text{cov}} s_T) s_T^{\alpha_q + \alpha_{\mathcal{M}}} \sqrt{\log(s_T)} \right] \right) \\ &\leq O \left(V_{\mathcal{M}} d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T]) \mathbb{E}[s_T]^{\alpha_q + \alpha_{\mathcal{M}}} \sqrt{\log(\mathbb{E}[s_T])} \right), \end{aligned} \tag{60}$$

where the second inequality uses Jensen's inequality and [Lemma F.16](#) to pass the expectation through, which holds as long as $1/100 \leq \alpha_q + \alpha_{\mathcal{M}} \leq 3/4$. Similarly,

$$\mathbb{E} \left[\sum_{s=s^*}^{s_T} s^{\alpha_n/2} \cdot 4\mathbf{g}^* \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel \widehat{M}(\pi))] \right] \leq O \left(\mathbf{g}^* V_{\mathcal{M}} d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T]) \mathbb{E}[s_T]^{\alpha_n/2} \sqrt{\log(\mathbb{E}[s_T])} \right), \tag{61}$$

where we have again used Jensen's inequality and [Lemma F.16](#) to pass the expectation through, which holds as long as $1/100 \leq \alpha_n/2 \leq 3/4$. For $s \leq s^*$, we do not have $M^* \in \mathcal{M}^\ell$, and therefore the estimation guarantees are vacuous. In this regime, using that the KL divergence is always bounded by $2V_{\mathcal{M}}$ ([Lemma F.13](#)), we can simply upper bound the estimation error by $2V_{\mathcal{M}}$. Thus,

$$\begin{aligned} &\mathbb{E} \left[\sum_{s=1}^{s_T} s^{1/2} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D(\widehat{M}(\pi) \parallel M^*(\pi))] \right] \\ &\leq \mathbb{E} \left[\sum_{s=s^*}^{s_T} s^{1/2} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D(\widehat{M}(\pi) \parallel M^*(\pi))] \right] + 2V_{\mathcal{M}} (s^*)^{3/2} \\ &\leq O \left(V_{\mathcal{M}} d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T]) \mathbb{E}[s_T]^{1/2} \sqrt{\log(\mathbb{E}[s_T])} \right) + 2V_{\mathcal{M}} (s^*)^{3/2}, \end{aligned}$$

which gives, applying [Lemma F.16](#) again:

$$\mathbb{E} \left[\sum_{s=1}^{s_T} \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))] \cdot \mathbb{I}\{\pi_* \in \pi_{\widehat{M}^s}\} \right] \leq \log T + O \left(V_{\mathcal{M}} \mathbb{E}[s_T]^{1/2} \sqrt{d_{\text{cov}} \cdot \log(C_{\text{cov}} \mathbb{E}[s_T])} \right)$$

$$\begin{aligned}
 & + (V_{\mathcal{M}} + L_{\text{KL}}) \left(\mathbb{E}[s_T]^{1/2} + V_{\mathcal{M}} d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T]) \mathbb{E}[s_T]^{1/2} \sqrt{\log(\mathbb{E}[s_T])} \right) \\
 & + 2V_{\mathcal{M}}(s^*)^{3/2} + \log \log T + 7V_{\mathcal{M}}. \tag{62}
 \end{aligned}$$

Combining Eqs. (60) to (62), we have

$$\begin{aligned}
 \mathbb{E} \left[\sum_{s=s^*}^{s_T} \Delta^*(\pi^s) \right] & \leq (1 + \delta) \mathbf{g}^* \log T + O \left(V_{\mathcal{M}} \sqrt{d_{\text{cov}} \mathbb{E}[s_T]} \log(C_{\text{cov}} \mathbb{E}[s_T]) \right. \\
 & \quad + (V_{\mathcal{M}} + L_{\text{KL}}) V_{\mathcal{M}} d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T]) \mathbb{E}[s_T]^{1/2} \sqrt{\log(\mathbb{E}[s_T])} \\
 & \quad + V_{\mathcal{M}} d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T]) \sqrt{\log \mathbb{E}[s_T]} ((1 + \mathbf{g}^*) \mathbb{E}[s_T]^{\alpha_q + \alpha_{\mathcal{M}}} + \varepsilon \mathbf{g}^* \mathbb{E}[s_T]^{\alpha_n/2}) \\
 & \quad \left. + \mathbf{g}^* \left(\frac{1}{\delta \mathbf{g}^*} \right)^{\frac{1}{\alpha_q - \alpha_n}} + \log \log T + V_{\mathcal{M}}(s^*)^{3/2} \right).
 \end{aligned}$$

To control this, it remains to bound s_T . By Lemma F.6 we have the following almost sure bound:

$$\begin{aligned}
 s_T & \leq \underbrace{\frac{4}{\delta \mathbf{g}^*} \left(\sum_{s=s^*}^{s_T} (1 + \mathbf{g}^*) s^{\alpha_q + \alpha_{\mathcal{M}}} \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]] \right)}_{(a)} \\
 & \quad + \underbrace{\sum_{s=s^*}^{s_T} s^{\alpha_n} \cdot 2\mathbf{g}^* \cdot \inf_{M \in \mathcal{M}_{\text{alt}}^*} \mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| M(\pi))] \cdot \mathbb{I}\{\pi_* \in \pi_{\widehat{M}^s}\}}_{(b)} \\
 & \quad + \underbrace{\sum_{s=1}^{s_T} s^{3\alpha_n/2} \cdot 4\mathbf{g}^* \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{p^s} [D_{\text{KL}}(M^*(\pi) \| \widehat{M}(\pi))]] + 4\mathbf{g}^* \left(\frac{4}{\delta \mathbf{g}^*} \right)^{\frac{1+\alpha_n}{\alpha_q - \alpha_n}}}_{(c)} + s^*.
 \end{aligned}$$

We bound the expectation of term (a) as in Eq. (60). To bound term (b), we apply Lemma F.4 with $\alpha = \alpha_n$ to get

$$\begin{aligned}
 \mathbb{E}[(b)] & \leq \mathbb{E}[s_T^{\alpha_n}] \log T + \mathbb{E} \left[V_{\mathcal{M}} s_T^{1/2 + \alpha_n} \cdot \sqrt{1344 d_{\text{cov}} \cdot \log(128 C_{\text{cov}} s_T)} \right. \\
 & \quad \left. + (V_{\mathcal{M}} + L_{\text{KL}}) \left(4 \frac{s_T^{1/2 + \alpha_n/2}}{1 - \alpha_n} + \sum_{s=1}^{s_T} s^{1/2 + \alpha_n/2} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]] \right) \right] \\
 & \quad + \mathbb{E}[s_T^{\alpha_n}] \log \log T + 7V_{\mathcal{M}}.
 \end{aligned}$$

Again applying Lemma F.7 and Lemma F.16, we have

$$\mathbb{E}[(c)] \leq O \left(\mathbf{g}^* V_{\mathcal{M}} d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T]) \mathbb{E}[s_T]^{3\alpha_n/2} \sqrt{\log(\mathbb{E}[s_T])} \right)$$

and

$$\mathbb{E} \left[\sum_{s=1}^{s_T} s^{1/2 + \alpha_n/2} \cdot \mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D(\widehat{M}(\pi) \| M^*(\pi))]] \right] \leq O \left(V_{\mathcal{M}} d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T]) \mathbb{E}[s_T]^{1/2 + \alpha_n/2} \sqrt{\log(\mathbb{E}[s_T])} \right)$$

$$+ 2V_{\mathcal{M}}(s^*)^{3/2+\alpha_n/2},$$

as long as $1/100 \leq \alpha_n \leq 1/4$. We therefore have (using [Lemma F.16](#) to pass the expectation through):

$$\begin{aligned} \mathbb{E}[s_T] &\leq \frac{1}{\delta \mathbf{g}^*} \cdot O\left(\mathbb{E}[s_T]^{\alpha_n} \log T + V_{\mathcal{M}} \mathbb{E}[s_T]^{1/2+\alpha_n} \sqrt{d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T])} + (V_{\mathcal{M}} + L_{\text{KL}}) \mathbb{E}[s_T]^{1/2+\alpha_n/2} \right. \\ &\quad + (V_{\mathcal{M}} + L_{\text{KL}}) V_{\mathcal{M}} d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T]) \mathbb{E}[s_T]^{1/2+\alpha_n/2} \sqrt{\log(\mathbb{E}[s_T])} \\ &\quad + V_{\mathcal{M}} d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T]) (1 + \mathbf{g}^*) \mathbb{E}[s_T]^{\alpha_q + \alpha_{\mathcal{M}}} \sqrt{\log(\mathbb{E}[s_T])} \\ &\quad \left. + \mathbf{g}^* V_{\mathcal{M}} d_{\text{cov}} \log(C_{\text{cov}} \mathbb{E}[s_T]) \mathbb{E}[s_T]^{3\alpha_n/2} \sqrt{\log(\mathbb{E}[s_T])} + \mathbf{g}^* \left(\frac{1}{\delta \mathbf{g}^*}\right)^{\frac{1+\alpha_n}{\alpha_q - \alpha_n}} + V_{\mathcal{M}}(s^*)^{3/2+\alpha_n/2} \right). \end{aligned}$$

We now set $\alpha_{\mathcal{M}} = 1/2$, $\alpha_n = 1/8$, and $\alpha_q = 1/4$, and note that all of the preceding parameter restrictions are satisfied for these choices. Furthermore, this parameter choice implies that—using [Lemma F.17](#) to handle log factors—we have

$$\mathbb{E}[s_T] \leq \tilde{O}\left(\frac{1}{(\delta \mathbf{g}^*)^{8/7}} \log^{8/7} T + \text{poly}\left(V_{\mathcal{M}}, d_{\text{cov}}, \log C_{\text{cov}}, L_{\text{KL}}, \mathbf{g}^*, \frac{1}{\delta}, \frac{1}{\mathbf{g}^*} \right) + \frac{V_{\mathcal{M}}(s^*)^{3/2+\alpha_n/2}}{\delta \mathbf{g}^*} \right).$$

Plugging this into the regret bound given above, we have

$$\begin{aligned} \mathbb{E}\left[\sum_{s=s^*}^{s_T} \Delta^*(\pi^s) \right] &\leq (1 + \delta) \mathbf{g}^* \log T + \tilde{O}\left(\text{poly}\left(V_{\mathcal{M}}, d_{\text{cov}}, \log C_{\text{cov}}, L_{\text{KL}}, \mathbf{g}^*, \frac{1}{\delta}, \frac{1}{\mathbf{g}^*}, \log \log T \right) \cdot \log^{6/7} T \right. \\ &\quad \left. + \frac{(1 + 1/\mathbf{g}^*)(V_{\mathcal{M}} + L_{\text{KL}}) V_{\mathcal{M}}^2 d_{\text{cov}} \log(C_{\text{cov}}) + V_{\mathcal{M}}}{\delta} \cdot (s^*)^{3/2} \cdot \log^{3/2} \log T \right). \end{aligned}$$

Finally, by [Lemma F.5](#), we can bound s^* as

$$s^* \leq \left(\overline{\text{aec}}_{\varepsilon/2}^{\text{D}, \mathcal{M}}(\mathcal{M}^*) \right)^{\frac{1}{\alpha_{\mathcal{M}}}} + \left(n_{\varepsilon}^* + \frac{1}{\Delta_{\min}^*} + \frac{4\mathbf{g}^*}{\Delta_{\min}^*} + \frac{2n_{\varepsilon}^*}{\mathbf{g}^*} + \frac{4}{\Delta_{\min}^* \varepsilon} \right)^{\frac{2}{\alpha_n}},$$

and, by [Lemma E.3](#) and [Lemma F.13](#), we can lower bound $\mathbf{g}^* \geq \Delta_{\min}^*/2V_{\mathcal{M}}$. Plugging this in gives the final bound. \square

F.3. Estimation Guarantees

In this section, we analyze [Algorithm 7](#), which is a variant of the Tempered Aggregation algorithm ([Algorithm 4](#)) designed for infinite classes. This algorithm is used within [Algorithm 5](#) in order to prove [Theorem F.1](#).

[Algorithm 7](#) simply applies [Algorithm 4](#) to a sequence of covers for the class \mathcal{M} on a doubling epoch schedule. In particular, at every epoch ℓ , [Algorithm 7](#) restarts the Tempered Aggregation algorithm ([Algorithm 4](#)), clearing the Tempered Aggregation instance from the previous epoch from memory. We denote the ℓ th instantiation of Tempered Aggregation as `TemperedAggregationℓ`.

Algorithm 7 Estimation with Adaptive Covering

- 1: **input:** Class \mathcal{M} .
 - 2: $\ell \leftarrow 1, \mathfrak{D}^\ell \leftarrow \emptyset$.
 - 3: $\mathcal{M}_{\text{cov}}^\ell \leftarrow (\rho_\ell, \mu_\ell)$ -cover of \mathcal{M} for $\rho_\ell \leftarrow 2^{-\ell}, \mu_\ell \leftarrow 2^{-5\ell}$.
 - 4: Initialize `TemperedAggregationℓ` as an instance of [Algorithm 4](#) with $\mathcal{M}_{\text{cov}}^\ell$.
 - 5: **for** $s = 1, 2, 3, \dots$ **do**
 - 6: Receive $(\pi^s, r^s, o^s), \mathfrak{D}^\ell \leftarrow \mathfrak{D}^\ell \cup \{(\pi^s, r^s, o^s)\}$.
 - 7: $\xi^s \leftarrow \text{TemperedAggregation}^\ell(\mathfrak{D}^\ell), \widehat{M}^s \leftarrow \mathbb{E}_{M \sim \xi^s}[M]$
 - 8: **if** $s \geq 2^\ell$ **then.**
 - 9: $\ell \leftarrow \ell + 1, \mathfrak{D}^\ell \leftarrow \emptyset$.
 - 10: $\mathcal{M}_{\text{cov}}^\ell \leftarrow (\rho_\ell, \mu_\ell)$ -cover of \mathcal{M} for $\rho_\ell \leftarrow 2^{-\ell}, \mu_\ell \leftarrow 2^{-5\ell}$.
 - 11: Initialize `TemperedAggregationℓ` as an instance of [Algorithm 4](#) with $\mathcal{M}_{\text{cov}}^\ell$.
-

Lemma F.7. Let τ denote some stopping time with respect to the filtration $(\mathcal{F}^t)_{t=1}^T$ such that $\tau \leq T$ almost surely, and let $\alpha \in (0, 1)$. When running [Algorithm 7](#) under [Assumption A.3](#), we have

$$\mathbb{E} \left[\sum_{s=1}^{\tau} s^\alpha \cdot \mathbb{E}_{M \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{H}}^2(M^*(\pi), M(\pi))]] \right] \leq \mathbb{E} \left[(20d_{\text{cov}} \cdot \log(64C_{\text{cov}}\tau) + 1) \cdot \frac{2^{2\alpha}}{2^\alpha - 1} \tau^\alpha \right] + 4.$$

In addition, if [Assumption A.2](#) also holds, then

$$\begin{aligned} & \mathbb{E} \left[\sum_{s=1}^{\tau} s^\alpha \cdot \mathbb{E}_{M \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))]] \right] \\ & \leq \mathbb{E} \left[(2 + 6V_{\mathcal{M}})(20d_{\text{cov}} \cdot \log(64C_{\text{cov}}\tau) + 1) \cdot \frac{2^{2\alpha}}{2^\alpha - 1} \tau^\alpha \sqrt{\log(2\tau)} \right] \\ & \quad + \mathbb{E}[32(1 + V_{\mathcal{M}}) \log(\tau)] + 8 \end{aligned} \tag{63}$$

and

$$\begin{aligned} & \mathbb{E} \left[\sum_{s=1}^{\tau} s^\alpha \cdot \mathbb{E}_{M \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(M(\pi) \parallel M^*(\pi))]] \right] \\ & \leq \mathbb{E} \left[(2 + 6V_{\mathcal{M}})(20d_{\text{cov}} \cdot \log(64C_{\text{cov}}\tau) + 1) \cdot \frac{2^{2\alpha}}{2^\alpha - 1} \tau^\alpha \sqrt{\log(2\tau)} \right] \\ & \quad + \mathbb{E}[32(1 + V_{\mathcal{M}}) \log(\tau)] + 8. \end{aligned} \tag{64}$$

Proof of Lemma F.7. Let \mathcal{S}^k denote the set of s values for which $\ell = k$ and note that $\mathcal{S}^k = \{2^{k-1} + 1, 2^{k-1} + 2, \dots, 2^k\}$. We can bound

$$\begin{aligned} & \mathbb{E} \left[\sum_{s=1}^{\tau} s^\alpha \mathbb{E}_{M \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{H}}^2(M^*(\pi), M(\pi))]] \right] \\ & \leq \sum_{k=1}^{\lceil \log_2 T \rceil} \mathbb{E} \left[\sum_{s \in \mathcal{S}^k} s^\alpha \cdot \mathbb{E}_{M \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{H}}^2(M^*(\pi), M(\pi))]] \cdot \mathbb{I}\{s \leq \tau\} \right] \end{aligned} \tag{65}$$

$$\leq \sum_{k=1}^{\lceil \log_2 T \rceil} 2^{\alpha k} \cdot \mathbb{E} \left[\sum_{s \in \mathcal{S}^k} \mathbb{E}_{M \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\mathbb{H}}^2(M^*(\pi), M(\pi))] \cdot \mathbb{I}\{s \leq \tau\}] \right] \quad (66)$$

since we have that $\ell \leq \lceil \log_2 T \rceil$ by construction, and since $s \leq 2^k$ for $s \in \mathcal{S}^k$. Let A_k denote the event

$$A_k := \left\{ \forall s \in \mathcal{S}^k : \sum_{i=2^{k-1}}^s \mathbb{E}_{M \sim \xi^i} [\mathbb{E}_{\pi \sim p^i} [D_{\mathbb{H}}^2(M^*(\pi), M(\pi))] \leq 2 \log \frac{2^k \mathbf{N}_{\text{cov}}(\mathcal{M}, \rho_k, \mu_k)}{\delta_k} + 2^k \rho_k \right\}.$$

By [Proposition E.2](#) and a union bound, $\mathbb{P}[A_k] \geq 1 - \delta_k - 2^{2k} \mu_k$. Since the Hellinger distance is always bounded by 2 and $|\mathcal{S}^k| \leq 2^k$, we can upper bound

$$\text{Eq. (66)} \leq \sum_{k=1}^{\lceil \log_2 T \rceil} 2^{\alpha k} \left(\sum_{s \in \mathcal{S}^k} \mathbb{E} [\mathbb{E}_{M \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\mathbb{H}}^2(M^*(\pi), M(\pi))] \cdot \mathbb{I}\{s \leq \tau, A_k\}] + 2 \cdot 2^k \mathbb{E}[\mathbb{I}\{A_k^c\}] \right).$$

Choosing $\delta_k = 2^{-3k}$ and since $\mu_k = 2^{-5k}$, we have

$$\sum_{k=1}^{\lceil \log_2 T \rceil} 2^{\alpha k} \cdot 2^{k+1} \mathbb{E}[\mathbb{I}\{A_k^c\}] \leq \sum_{k=1}^{\lceil \log_2 T \rceil} 2^{2k+1} \cdot (\delta_k + 2^{2k} \mu_k) = \sum_{k=1}^{\lceil \log_2 T \rceil} 2^{2k+1} \cdot 2 \cdot 2^{-3k} \leq 4.$$

Note that for $s \in \mathcal{S}^k$, if $s \leq \tau$, this implies that $2^{k-1} \leq \tau$. On the event A_k , for $\rho_k = 2^{-k}$, when $\alpha > 0$ we can bound

$$\begin{aligned} & \sum_{k=1}^{\lceil \log_2 T \rceil} 2^{\alpha k} \sum_{s \in \mathcal{S}^k} \mathbb{E} [\mathbb{E}_{M \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\mathbb{H}}^2(M^*(\pi), M(\pi))] \cdot \mathbb{I}\{s \leq \tau, A_k\}] \\ & \leq \mathbb{E} \left[\sum_{k=1}^{\lceil \log_2 T \rceil} 2^{\alpha k} \left(2 \log \frac{2^k \mathbf{N}_{\text{cov}}(\mathcal{M}, 2^{-k}, 2^{-5k})}{2^{-3k}} + 1 \right) \cdot \mathbb{I}\{2^{k-1} \leq \tau\} \right] \\ & \leq \mathbb{E} \left[\left(2 \log \frac{\mathbf{N}_{\text{cov}}(\mathcal{M}, \tau^{-1}/2, \tau^{-5}/32)}{\tau^{-4}/16} + 1 \right) \cdot \max_k \left(\frac{2^\alpha}{2^\alpha - 1} 2^{\alpha k} \cdot \mathbb{I}\{2^{k-1} \leq \tau\} \right) \right] \\ & \leq \mathbb{E} \left[\left(2 \log \frac{\mathbf{N}_{\text{cov}}(\mathcal{M}, \tau^{-1}/2, \tau^{-5}/32)}{\tau^{-4}/16} + 1 \right) \cdot \frac{2^{2\alpha}}{2^\alpha - 1} \tau^\alpha \right] \end{aligned} \quad (67)$$

where the final two inequalities follow since $2^k \leq 2\tau$. Under [Assumption A.3](#) we have

$$\log \frac{\mathbf{N}_{\text{cov}}(\mathcal{M}, \tau^{-1}/2, \tau^{-5}/32)}{\tau^{-4}/16} \leq 10d_{\text{cov}} \cdot \log(64C_{\text{cov}}\tau),$$

which gives the first result.

Bound on KL Estimation Error. By [Lemma F.14](#), for any $x > 0$ we have

$$D_{\text{KL}}(M(\pi) \parallel \widetilde{M}(\pi)) \leq (2 + 2V_{\mathcal{M}} + x) \cdot D_{\mathbb{H}}^2(M(\pi), \widetilde{M}(\pi)) + 32(1 + V_{\mathcal{M}}^2/x + V_{\mathcal{M}}^3/x^2) \cdot \exp(-x^2/8V_{\mathcal{M}}^2).$$

In particular choosing $x = V_{\mathcal{M}}\sqrt{8\log s^2}$, we have

$$D_{\text{KL}}(M(\pi) \parallel \widetilde{M}(\pi)) \leq (2 + 2V_{\mathcal{M}} + V_{\mathcal{M}}\sqrt{8\log s}) \cdot D_{\text{H}}^2(M(\pi), \widetilde{M}(\pi)) + 32(1 + V_{\mathcal{M}})/s.$$

Repeating this for each step s , we can therefore bound

$$\begin{aligned} & \mathbb{E} \left[\sum_{s=1}^{\tau} s^{\alpha} \mathbb{E}_{M \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))]] \right] \\ & \leq (2 + 6V_{\mathcal{M}}) \cdot \mathbb{E} \left[\sum_{s=1}^{\tau} s^{\alpha} \sqrt{\log s} \mathbb{E}_{M \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{H}}^2(M^*(\pi), M(\pi))]] \right] \\ & \quad + 32(1 + V_{\mathcal{M}}) \cdot \mathbb{E}[\log \tau]. \end{aligned}$$

The result in (63) then follows from a calculation nearly identical to our above bound on Hellinger estimation error. Applying Lemma F.14 in a similar fashion with the arguments flipped gives (64). \square

In the following, we extend Lemma F.7 to the case when $\alpha = 0$.

Lemma F.8. *Let τ denote some stopping time with respect to the filtration $(\mathcal{F}^t)_{t=1}^T$ such that $\tau \leq T$ almost surely. When running Algorithm 7, under Assumption A.3, we have*

$$\mathbb{E} \left[\sum_{s=1}^{\tau} \mathbb{E}_{M \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{H}}^2(M^*(\pi), M(\pi))]] \right] \leq \mathbb{E}[(20d_{\text{cov}} \cdot \log(64C_{\text{cov}}\tau) + 1) \cdot (2\log \tau + 1)] + 4.$$

In addition, if Assumption A.2 also holds,

$$\begin{aligned} & \mathbb{E} \left[\sum_{s=1}^{\tau} \mathbb{E}_{M \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(M^*(\pi) \parallel M(\pi))]] \right] \\ & \leq \mathbb{E} \left[(2 + 5V_{\mathcal{M}})(20d_{\text{cov}} \cdot \log(64C_{\text{cov}}\tau) + 1)(2\log \tau + 1) \cdot \sqrt{\log(2\tau)} \right] \\ & \quad + \mathbb{E}[32(1 + V_{\mathcal{M}}) \log(\tau)] + 8 \end{aligned}$$

and

$$\begin{aligned} & \mathbb{E} \left[\sum_{s=1}^{\tau} \mathbb{E}_{M \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{KL}}(M(\pi) \parallel M^*(\pi))]] \right] \\ & \leq \mathbb{E} \left[(2 + 5V_{\mathcal{M}})(20d_{\text{cov}} \cdot \log(64C_{\text{cov}}\tau) + 1)(2\log \tau + 1) \cdot \sqrt{\log(2\tau)} \right] \\ & \quad + \mathbb{E}[32(1 + V_{\mathcal{M}}) \log(\tau)] + 8. \end{aligned}$$

Proof of Lemma F.8. This follows identically to Lemma F.7 but replacing Eq. (67) with

$$\begin{aligned} & \sum_{k=1}^{\lceil \log_2 T \rceil} 2^{\alpha k} \sum_{s \in \mathcal{S}^k} \mathbb{E} [\mathbb{E}_{M \sim \xi^s} [\mathbb{E}_{\pi \sim p^s} [D_{\text{H}}^2(M^*(\pi), M(\pi))]]] \cdot \mathbb{I}\{s \leq \tau, A_k\} \\ & \leq \mathbb{E} \left[\sum_{k=1}^{\lceil \log_2 T \rceil} \left(2\log \frac{2^k \mathbf{N}_{\text{cov}}(\mathcal{M}, 2^{-k}, 2^{-5k})}{2^{-3k}} + 1 \right) \cdot \mathbb{I}\{2^{k-1} \leq \tau\} \right] \end{aligned}$$

$$\leq \mathbb{E} \left[\left(2 \log \frac{\mathbf{N}_{\text{cov}}(\mathcal{M}, \tau^{-1}/2, \tau^{-5}/32)}{\tau^{-4}/16} + 1 \right) \cdot (2 \log \tau + 1) \right].$$

The bound on the KL estimation error also follows from the same reasoning as in [Lemma F.7](#). \square

F.4. Supporting Lemmas

Lemma F.9. Consider running either [Algorithm 5](#) or [Algorithm 6](#). Assume that $\pi_\star \notin \pi_{\widehat{M}^s}$ and that $\min_{M \in \mathcal{M} : \xi^s(M) > 0} \Delta_{\min}^M \geq \Delta$. Then $\mathbb{E}_{M \sim \xi^s} [\mathbb{I}\{M \in \mathcal{M}_{\text{alt}}^\star\}] \geq \frac{1}{2} \Delta$.

Proof of Lemma F.9. Recall that $\widehat{M}^s = \mathbb{E}_{M \sim \xi^s} [M]$, so $\pi \in \pi_{\widehat{M}^s}$ implies that $\pi \in \arg \max_{\pi' \in \Pi} \mathbb{E}_{M \sim \xi^s} [f^M(\pi')]$. If $\pi_\star \notin \pi_{\widehat{M}^s}$, then there exists some $\tilde{\pi}$ such that $\mathbb{E}_{M \sim \xi^s} [f^M(\tilde{\pi})] > \mathbb{E}_{M \sim \xi^s} [f^M(\pi_\star)]$. Since $f^M(\pi) \in [0, 1]$, this implies that

$$0 < \mathbb{E}_{M \sim \xi^s} [f^M(\tilde{\pi}) - f^M(\pi_\star)] \leq \mathbb{E}_{M \sim \xi^s} [\mathbb{I}\{M \in \mathcal{M}_{\text{alt}}^\star\}] - \mathbb{E}_{M \sim \xi^s} [(f^M(\pi_\star) - f^M(\tilde{\pi})) \cdot \mathbb{I}\{M \notin \mathcal{M}_{\text{alt}}^\star\}].$$

For $M \notin \mathcal{M}_{\text{alt}}^\star$, we have $f^M(\pi_\star) - f^M(\tilde{\pi}) \geq \Delta_{\min}^M \geq \Delta$. Thus, the above implies

$$\begin{aligned} 0 &< \mathbb{E}_{M \sim \xi^s} [\mathbb{I}\{M \in \mathcal{M}_{\text{alt}}^\star\}] - \Delta \cdot \mathbb{E}_{M \sim \xi^s} [\mathbb{I}\{M \notin \mathcal{M}_{\text{alt}}^\star\}] \\ \iff \Delta \cdot (1 - \mathbb{E}_{M \sim \xi^s} [\mathbb{I}\{M \in \mathcal{M}_{\text{alt}}^\star\}]) &< \mathbb{E}_{M \sim \xi^s} [\mathbb{I}\{M \in \mathcal{M}_{\text{alt}}^\star\}]. \end{aligned}$$

Rearranging gives $\mathbb{E}_{M \sim \xi^s} [\mathbb{I}\{M \in \mathcal{M}_{\text{alt}}^\star\}] \geq \frac{\Delta}{1+\Delta} \geq \frac{1}{2} \Delta$. \square

Lemma F.10. When running [Algorithm 5](#), on rounds s for which $M^\star \in \mathcal{M} \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda^s; \mathbf{n}_{\max})$, we have

$$\gamma^s \leq \frac{(1 + \delta)(4\mathbf{n}_{\max} + 2\delta \underline{\mathbf{g}}^{\mathcal{M}})}{\delta \underline{\mathbf{g}}^{\mathcal{M}}} \cdot \overline{\text{aec}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}),$$

for γ^s as defined in [Eq. \(50\)](#).

Proof of Lemma F.10. Recall that ω^s and λ^s are chosen to minimize

$$\sup_{M \in \mathcal{M} \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda^s; \mathbf{n}_{\max})} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_\omega [D(\widehat{M}(\pi) \| M(\pi))]]}.$$

Since $M^\star \in \mathcal{M} \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda^s; \mathbf{n}_{\max})$, we can therefore bound

$$\begin{aligned} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\omega^s} [D(\widehat{M}(\pi) \| M^\star(\pi))]]} &\leq \inf_{\omega \in \Delta_\Pi} \sup_{M \in \mathcal{M} \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda^s; \mathbf{n}_{\max})} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_\omega [D(\widehat{M}(\pi) \| M(\pi))]]} \\ &= \inf_{\lambda, \omega \in \Delta_\Pi} \sup_{M \in \mathcal{M} \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda; \mathbf{n}_{\max})} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_\omega [D(\widehat{M}(\pi) \| M(\pi))]]}. \end{aligned}$$

Recall that we set

$$\mathbf{n}_{\max} = \left(\frac{1}{\Delta_{\min} \varepsilon} + \frac{2V_{\mathcal{M}} \mathbf{n}_\varepsilon^{\mathcal{M}}}{\Delta_{\min}} \right) \cdot \max_{M \in \mathcal{M}} \frac{32\mathbf{g}^M}{\Delta_{\min}}.$$

By [Lemma F.13](#), under [Assumption A.2](#), we can bound $D_{\text{KL}}(M(\pi) \| M'(\pi)) \leq 2V_{\mathcal{M}}$ for all $M, M' \in \mathcal{M}$ and $\pi \in \Pi$. It follows from [Lemma E.3](#) that

$$\frac{2V_{\mathcal{M}}}{\Delta_{\min}} \geq \frac{1}{\min_{M \in \mathcal{M}: \mathbf{g}^M > 0} \mathbf{g}^M},$$

so

$$n_{\max} \geq \left(\frac{1}{\Delta_{\min} \varepsilon} + \frac{n_{\varepsilon}^{\mathcal{M}}}{\min_{M \in \mathcal{M}: \mathbf{g}^M > 0} \mathbf{g}^M} \right) \cdot \max_{M \in \mathcal{M}} \frac{32\mathbf{g}^M}{\Delta_{\min}}.$$

Given this, straightforward calculation shows that n_{\max} meets the condition required by [Lemma E.4](#), so [Lemma E.4](#) implies

$$\begin{aligned} \inf_{\lambda, \omega \in \Delta_{\Pi}} \sup_{M \in \mathcal{M} \setminus \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda; n_{\max})} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\omega} [D(\widehat{M}(\pi) \| M(\pi))]]} &\leq \inf_{\lambda, \omega \in \Delta_{\Pi}} \sup_{M \in \mathcal{M} \setminus \mathcal{M}_{\varepsilon/2}^{\text{gl}}(\lambda)} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\omega} [D(\widehat{M}(\pi) \| M(\pi))]]} \\ &= \overline{\text{aec}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}; \xi^s) \\ &\leq \overline{\text{aec}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}). \end{aligned}$$

By our choice of q we have $\frac{1}{1-q} = \frac{4n_{\max} + 2\delta \underline{\mathbf{g}}^{\mathcal{M}}}{\delta \underline{\mathbf{g}}^{\mathcal{M}}}$. We can then bound

$$\gamma^s = \frac{1 + \delta}{1 - q} \cdot \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\omega^s} [D(\widehat{M}(\pi) \| M^*(\pi))]]} \leq \frac{(1 + \delta)(4n_{\max} + 2\delta \underline{\mathbf{g}}^{\mathcal{M}})}{\delta \underline{\mathbf{g}}^{\mathcal{M}}} \cdot \overline{\text{aec}}_{\varepsilon/2}^{\text{D}}(\mathcal{M}).$$

□

Lemma F.11. Consider running [Algorithm 6](#). Then on rounds s for which $M^* \in \mathcal{M}^{\ell} \setminus \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda^s; n^s)$, we have

$$\gamma^s \leq (1 + \mathbf{g}^* \delta) \cdot s^{\alpha_q + \alpha_{\mathcal{M}}}$$

for γ^s as defined in [Eq. \(56\)](#), and $\alpha_q, \alpha_{\mathcal{M}}$ parameters of [Algorithm 6](#).

Proof of Lemma F.11. Assume we are at epoch ℓ . Recall that ω^s and λ^s minimize

$$\sup_{M \in \mathcal{M}^{\ell} \setminus \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda^s; n^s)} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\omega} [D(\widehat{M}(\pi) \| M(\pi))]]}.$$

Since $M^* \in \mathcal{M}^{\ell} \setminus \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda^s; n^s)$, we have can therefore bound

$$\begin{aligned} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\omega^s} [D(\widehat{M}(\pi) \| M^*(\pi))]]} &\leq \inf_{\omega \in \Delta_{\Pi}} \sup_{M \in \mathcal{M}^{\ell} \setminus \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda^s; n^s)} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\omega} [D(\widehat{M}(\pi) \| M(\pi))]]} \\ &= \inf_{\lambda, \omega \in \Delta_{\Pi}} \sup_{M \in \mathcal{M}^{\ell} \setminus \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda; n^s)} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\omega} [D(\widehat{M}(\pi) \| M(\pi))]]}. \end{aligned}$$

By construction, for every $M \in \mathcal{M}^\ell$, we have

$$n_\varepsilon^M + \frac{1}{\Delta_{\min}^M} + \frac{4\mathbf{g}^M}{\Delta_{\min}^M} + \frac{2n_\varepsilon^M}{\mathbf{g}^M} + \frac{4}{\Delta_{\min}^M \varepsilon} \leq \sqrt{n^s}.$$

This implies that

$$\zeta := \min_{M \in \mathcal{M}^\ell} \min \left\{ \frac{\mathbf{g}^M}{\mathbf{g}^M + 2n_\varepsilon^M}, \frac{\Delta_{\min}^M \varepsilon}{4} \right\} \geq \frac{1}{\sqrt{n^s}}$$

and

$$\max_{M \in \mathcal{M}^\ell} \max \left\{ n_\varepsilon^M, \frac{4\mathbf{g}^M}{\Delta_{\min}^M} \right\} \leq \sqrt{n^s},$$

which together imply that

$$\max_{M \in \mathcal{M}^\ell} \max \left\{ n_\varepsilon^M, \frac{4\mathbf{g}^M}{\Delta_{\min}^M}, \frac{2\mathbf{g}^M}{\zeta \Delta_{\min}^M} \right\} \leq n^s$$

By [Lemma E.4](#) and since \mathcal{M}^ℓ is constructed such that $\inf_{M \in \mathcal{M}^\ell} \Delta_{\min}^M > 0$, we can therefore bound

$$\begin{aligned} \inf_{\lambda, \omega \in \Delta_\Pi} \sup_{M \in \mathcal{M}^\ell \setminus \mathcal{M}_{\varepsilon/2}^{\text{gl}}(\lambda; n^s)} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_\omega [D(\widehat{M}(\pi) \| M(\pi))]]} &\leq \inf_{\lambda, \omega \in \Delta_\Pi} \sup_{M \in \mathcal{M}^\ell \setminus \mathcal{M}_{\varepsilon/2}^{\text{gl}}(\lambda)} \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_\omega [D(\widehat{M}(\pi) \| M(\pi))]]} \\ &\leq \overline{\text{aec}}_{\varepsilon/2}^{\text{D}, \mathcal{M}}(\mathcal{M}_{\Delta^\ell, \frac{1}{\Delta^\ell}}; \xi^s) \end{aligned}$$

where the last inequality holds by the definition of \mathcal{M}^ℓ and Δ^ℓ . Note that by construction, ξ^s is only supported on \mathcal{M}^ℓ , so we can bound $\overline{\text{aec}}_{\varepsilon/2}^{\text{D}, \mathcal{M}}(\mathcal{M}_{\Delta^\ell, \frac{1}{\Delta^\ell}}; \xi^s) \leq \overline{\text{aec}}_{\varepsilon/2}^{\text{D}, \mathcal{M}}(\mathcal{M}_{\Delta^\ell, \frac{1}{\Delta^\ell}})$. By construction, we also have $\overline{\text{aec}}_{\varepsilon/2}^{\text{D}, \mathcal{M}}(\mathcal{M}_{\Delta^\ell, \frac{1}{\Delta^\ell}}) \leq 2^{\ell \alpha_{\mathcal{M}}} \leq s^{\alpha_{\mathcal{M}}}$. Lastly, by our choice for q^s we have $\frac{1}{1-q^s} = s^{\alpha_q}$. We can then bound

$$\gamma^s = \frac{1 + \mathbf{g}^* \delta}{1 - q^s} \cdot \frac{1}{\mathbb{E}_{\widehat{M} \sim \xi^s} [\mathbb{E}_{\omega^s} [D(\widehat{M}(\pi) \| M^*(\pi))]]} \leq (1 + \delta) \cdot s^{\alpha_q + \alpha_{\mathcal{M}}}.$$

□

F.4.1. LIKELIHOOD RATIOS

Lemma F.12. *Consider running either [Algorithm 5](#) or [Algorithm 6](#). Under [Assumption A.2](#), with probability at least $1 - \delta$, we can bound, for any $M \in \mathcal{M}$, s , and $\beta_i > 0$,*

$$\begin{aligned} \sum_{i=1}^s \mathbb{E}_{\widehat{M} \sim \xi^i} [\mathbb{E}_{\pi \sim p^i} [D_{\text{KL}}(\widehat{M}(\pi) \| M(\pi))]] &\leq \sum_{i=1}^s \mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}_{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}_{M, \pi^i}(r^i, o^i)} \right] + V_{\mathcal{M}} \sqrt{56s \log 1/\delta} \\ &\quad + V_{\mathcal{M}} \cdot \left(\sum_{i=1}^s 1/\beta_i + \sum_{i=1}^s \beta_i \mathbb{E}_{\widehat{M} \sim \xi^i} [\mathbb{E}_{\pi \sim p^i} [D(\widehat{M}(\pi) \| M^*(\pi))]] \right). \end{aligned}$$

Proof of Lemma F.12. By Theorem 2 of [Shamir \(2011\)](#), under [Assumption A.2](#) we have that with probability at least $1 - \delta$,

$$\sum_{i=1}^s \mathbb{E}_{(r,o) \sim M^*} \left[\mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi}(r, o)}{\mathbb{P}_{M, \pi}(r, o)} \mid \mathcal{H}^{i-1} \right] \right] \leq \sum_{i=1}^s \mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}_{M, \pi^i}(r^i, o^i)} \right] + V_{\mathcal{M}} \sqrt{56s \log 1/\delta}. \quad (68)$$

Note that

$$\begin{aligned} \mathbb{E} \left[\mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}_{M, \pi^i}(r^i, o^i)} \mid \mathcal{H}^{i-1} \right] \right] &= \mathbb{E}_{\pi \sim p^i} \left[\mathbb{E}_{(r,o) \sim M^*(\pi)} \left[\mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi}(r, o)}{\mathbb{P}_{M, \pi}(r, o)} \right] \right] \right] \\ &= \mathbb{E}_{\pi \sim p^i} \left[\mathbb{E}_{\widehat{M} \sim \xi^i} \left[\mathbb{E}_{(r,o) \sim M^*(\pi)} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi}(r, o)}{\mathbb{P}_{M, \pi}(r, o)} \right] \right] \right] \end{aligned}$$

By Lemma B.4 of [Foster et al. \(2022b\)](#), we can bound, for any $\widehat{M} \in \mathcal{M}$,

$$\begin{aligned} &\left| \mathbb{E}_{(r,o) \sim M^*(\pi)} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi}(r, o)}{\mathbb{P}_{M, \pi}(r, o)} \right] - \mathbb{E}_{(r,o) \sim \widehat{M}(\pi)} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi}(r, o)}{\mathbb{P}_{M, \pi}(r, o)} \right] \right| \\ &\leq \sqrt{\frac{1}{2} \left(\mathbb{E}_{(r,o) \sim M^*(\pi)} \left[\log^2 \frac{\mathbb{P}^{\widehat{M}, \pi}(r, o)}{\mathbb{P}_{M, \pi}(r, o)} \right] + \mathbb{E}_{(r,o) \sim \widehat{M}(\pi)} \left[\log^2 \frac{\mathbb{P}^{\widehat{M}, \pi}(r, o)}{\mathbb{P}_{M, \pi}(r, o)} \right] \right)} \cdot D_{\mathbb{H}}^2(M^*(\pi), \widehat{M}(\pi)) \\ &\leq \sqrt{2V_{\mathcal{M}}^2 \cdot D_{\mathbb{H}}^2(M^*(\pi), \widehat{M}(\pi))} \end{aligned}$$

where the second inequality follows under the subgaussian assumption, [Assumption A.2](#). It follows that for any $\widehat{M} \in \mathcal{M}$,

$$\begin{aligned} &\mathbb{E}_{\pi \sim p^i} \left[\mathbb{E}_{(r,o) \sim M^*(\pi)} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi}(r, o)}{\mathbb{P}_{M, \pi}(r, o)} \right] \right] \\ &\geq \mathbb{E}_{\pi \sim p^i} \left[\mathbb{E}_{(r,o) \sim \widehat{M}(\pi)} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi}(r, o)}{\mathbb{P}_{M, \pi}(r, o)} \right] \right] - \sqrt{2V_{\mathcal{M}}^2 \cdot \mathbb{E}_{\pi \sim p^i} [D_{\mathbb{H}}^2(M^*(\pi), \widehat{M}(\pi))]} \\ &\stackrel{(a)}{=} \mathbb{E}_{\pi \sim p^i} [D_{\text{KL}}(\widehat{M}(\pi) \parallel M(\pi))] - \sqrt{2V_{\mathcal{M}}^2 \cdot \mathbb{E}_{\pi \sim p^i} [D_{\mathbb{H}}^2(M^*(\pi), \widehat{M}(\pi))]} \\ &\stackrel{(b)}{\geq} \mathbb{E}_{\pi \sim p^i} [D_{\text{KL}}(\widehat{M}(\pi) \parallel M(\pi))] - V_{\mathcal{M}}^2 \cdot (\beta \mathbb{E}_{\pi \sim p^i} [D_{\mathbb{H}}^2(M^*(\pi), \widehat{M}(\pi))] + 1/\beta) \end{aligned}$$

where (a) follows by the definition of KL divergence, and (b) from AM-GM for any $\beta > 0$. Since this bound holds uniformly for all $\widehat{M} \in \mathcal{M}$, this implies that

$$\begin{aligned} \mathbb{E} \left[\mathbb{E}_{\widehat{M} \sim \xi^i} \left[\log \frac{\mathbb{P}^{\widehat{M}, \pi^i}(r^i, o^i)}{\mathbb{P}_{M, \pi^i}(r^i, o^i)} \mid \mathcal{H}^{i-1} \right] \right] &\geq \mathbb{E}_{\widehat{M} \sim \xi^i} [\mathbb{E}_{\pi \sim p^i} [D_{\text{KL}}(\widehat{M}(\pi) \parallel M(\pi))]] \\ &\quad - V_{\mathcal{M}}^2 \cdot (\beta \mathbb{E}_{\widehat{M} \sim \xi^i} [\mathbb{E}_{\pi \sim p^i} [D_{\mathbb{H}}^2(M^*(\pi), \widehat{M}(\pi))]] + 1/\beta). \end{aligned}$$

Combining this with [Eq. \(68\)](#), and using [Assumption D.2](#) proves the result. \square

Lemma F.13. Under [Assumption A.2](#), we have, for any $M, M' \in \mathcal{M}$,

$$D_{\text{KL}}(M'(\pi) \parallel M(\pi)) \leq 2V_{\mathcal{M}}, \quad \forall \pi \in \Pi.$$

Proof of Lemma F.13. We have

$$D_{\text{KL}}(M'(\pi) \parallel M(\pi)) = \mathbb{E}_{(r,o) \sim M'(\pi)} \left[\log \frac{\mathbb{P}^{M',\pi}(r,o)}{\mathbb{P}^{M,\pi}(r,o)} \right] \leq \mathbb{E}_{(r,o) \sim M'(\pi)} \left[\left| \log \frac{\mathbb{P}^{M',\pi}(r,o)}{\mathbb{P}^{M,\pi}(r,o)} \right| \right] \leq 2V_{\mathcal{M}}$$

where the last inequality holds under [Assumption A.2](#). \square

Lemma F.14. Under [Assumption A.2](#), for any $x > 0$ and $M, \widetilde{M} \in \mathcal{M}$, we have

$$D_{\text{KL}}(M(\pi) \parallel \widetilde{M}(\pi)) \leq (2 + 2V_{\mathcal{M}} + x) \cdot D_{\text{H}}^2(M(\pi), \widetilde{M}(\pi)) + 32(1 + V_{\mathcal{M}}^2/x + V_{\mathcal{M}}^3/x^2) \cdot \exp(-x^2/8V_{\mathcal{M}}^2).$$

Proof of Lemma F.14. Fix π . Define

$$\mathcal{E} := \left\{ \left| \log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{\widetilde{M},\pi}(r,o)} - D_{\text{KL}}(M(\pi) \parallel \widetilde{M}(\pi)) \right| \leq x \right\}$$

and for $j \in \mathbb{N}$,

$$\mathcal{E}_j := \left\{ e^{j-1} \cdot x < \left| \log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{\widetilde{M},\pi}(r,o)} - D_{\text{KL}}(M(\pi) \parallel \widetilde{M}(\pi)) \right| \leq e^j \cdot x \right\}.$$

Note that $\mathcal{E}, (\mathcal{E}_j)_{j=1}^{\infty}$ form a partition of the probability space. Furthermore, since $D_{\text{KL}}(M(\pi) \parallel \widetilde{M}(\pi)) = \mathbb{E}_{o \sim M(\pi)} [\log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{\widetilde{M},\pi}(r,o)}]$, under [Assumption A.2](#) we have that $\mathbb{P}^{M,\pi}(\mathcal{E}_j) \leq 2 \exp(-x^2 e^{2(j-1)}/V_{\mathcal{M}}^2)$ and $\mathbb{P}^{M,\pi}(\mathcal{E}^c) \leq 2 \exp(-x^2/V_{\mathcal{M}}^2)$. Now,

$$\begin{aligned} D_{\text{KL}}(M(\pi) \parallel \widetilde{M}(\pi)) &= \int \log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{\widetilde{M},\pi}(r,o)} d\mathbb{P}^{M,\pi}(r,o) \\ &= \int_{\mathcal{E}} \log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{\widetilde{M},\pi}(r,o)} d\mathbb{P}^{M,\pi}(r,o) + \sum_{j=1}^{\infty} \int_{\mathcal{E}_j} \log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{\widetilde{M},\pi}(r,o)} d\mathbb{P}^{M,\pi}(r,o). \end{aligned}$$

Using that $D_{\text{KL}}(M(\pi) \parallel \widetilde{M}(\pi)) \leq 2V_{\mathcal{M}}$ by [Lemma F.13](#), we can bound

$$\begin{aligned} \sum_{j=1}^{\infty} \int_{\mathcal{E}_j} \log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{\widetilde{M},\pi}(r,o)} d\mathbb{P}^{M,\pi}(r,o) &\leq \sum_{j=1}^{\infty} (e^j x + 2V_{\mathcal{M}}) \cdot \int_{\mathcal{E}_j} d\mathbb{P}^{M,\pi}(r,o) \\ &= \sum_{j=1}^{\infty} (e^j x + 2V_{\mathcal{M}}) \cdot \mathbb{P}^{M,\pi}(\mathcal{E}_j) \\ &\leq \sum_{j=1}^{\infty} (e^j x + 2V_{\mathcal{M}}) \cdot 2 \exp(-x^2 e^{2(j-1)}/V_{\mathcal{M}}^2) \\ &\leq \int_0^{\infty} (e^j x + 2V_{\mathcal{M}}) \cdot 2 \exp(-x^2 e^{2(j-1)}/V_{\mathcal{M}}^2) dj \end{aligned}$$

$$\begin{aligned}
 &\leq 4(x + V_{\mathcal{M}}) \int_0^\infty e^j \exp(-e^j \cdot x^2 e^{-2}/V_{\mathcal{M}}^2) dj \\
 &= 4(x + V_{\mathcal{M}}) V_{\mathcal{M}}^2 e^2 \exp(-x^2 e^{-2}/V_{\mathcal{M}}^2)/x^2 \\
 &\leq 32(V_{\mathcal{M}}^2/x + V_{\mathcal{M}}^3/x^2) \cdot \exp(-x^2/8V_{\mathcal{M}}^2).
 \end{aligned}$$

We turn now to the first term. Note that we can write

$$\begin{aligned}
 \int_{\mathcal{E}} \log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{\widetilde{M},\pi}(r,o)} d\mathbb{P}^{M,\pi}(r,o) &= \int_{\mathcal{E}, \mathbb{P}^{\widetilde{M},\pi}(r,o) > \mathbb{P}^{M,\pi}(r,o)} \left(\log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{\widetilde{M},\pi}(r,o)} + \frac{\mathbb{P}^{\widetilde{M},\pi}(r,o)}{\mathbb{P}^{M,\pi}(r,o)} - 1 \right) d\mathbb{P}^{M,\pi}(r,o) - \mathbb{P}^{\widetilde{M},\pi}(\mathcal{E}^c) \\
 &\quad + \int_{\mathcal{E}, \mathbb{P}^{\widetilde{M},\pi}(r,o) \leq \mathbb{P}^{M,\pi}(r,o)} \left(\frac{\mathbb{P}^M(o|\pi)}{\mathbb{P}^{\widetilde{M},\pi}(r,o)} \log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{\widetilde{M},\pi}(r,o)} + 1 - \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{\widetilde{M},\pi}(r,o)} \right) d\mathbb{P}^{\widetilde{M},\pi}(r,o) + \mathbb{P}^{M,\pi}(\mathcal{E}^c)
 \end{aligned}$$

Following the proof of Lemma 4 of [Yang and Barron \(1998\)](#) and using that $\log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{\widetilde{M},\pi}(r,o)} \leq D_{\text{KL}}(M(\pi) \parallel \widetilde{M}(\pi)) + x \leq 2V_{\mathcal{M}} + x$ on \mathcal{E} , we can bound this as

$$\begin{aligned}
 &\leq (2 + 2V_{\mathcal{M}} + x) \int_{\mathcal{E}} (\sqrt{d\mathbb{P}^{M,\pi}(r,o)} - \sqrt{d\mathbb{P}^{\widetilde{M},\pi}(r,o)})^2 + \mathbb{P}^{M,\pi}(\mathcal{E}^c) \\
 &\leq (2 + 2V_{\mathcal{M}} + x) \int (\sqrt{d\mathbb{P}^{M,\pi}(r,o)} - \sqrt{d\mathbb{P}^{\widetilde{M},\pi}(r,o)})^2 + \mathbb{P}^{M,\pi}(\mathcal{E}^c) \\
 &\leq (2 + 2V_{\mathcal{M}} + x) \cdot D_{\text{H}}^2(M(\pi), \widetilde{M}(\pi)) + 2 \exp(-x^2/V_{\mathcal{M}}^2).
 \end{aligned}$$

□

F.4.2. COVERING NUMBERS

Lemma F.15. *For any subset $\mathcal{M}' \subseteq \mathcal{M}$, there exists some (ρ, μ) -cover $\mathcal{M}'_{\text{cov}} \subseteq \mathcal{M}'$ for \mathcal{M}' such that $|\mathcal{M}'_{\text{cov}}| \leq N_{\text{cov}}(\mathcal{M}, \rho/2, \mu)$.*

Proof of Lemma F.15. Let \mathcal{M}_{cov} denote a $(\rho/2, \mu)$ -cover of \mathcal{M} with event \mathcal{E} . Throughout the proof we use the shorthand $(r, o, \pi) \in \mathcal{E}$ to denote that there exists $M \in \mathcal{M}$ such that $\mathbb{P}^{M,\pi}(r, o | \mathcal{E}) > 0$. By definition, it follows that for any $M \in \mathcal{M}$, there exists $M' \in \mathcal{M}_{\text{cov}}$ such that

$$\sup_{r,o,\pi : (r,o,\pi) \in \mathcal{E}} \left| \log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{M',\pi}(r,o)} \right| \leq \rho/2. \quad (69)$$

Let $\mathcal{M}'_{\text{cov}} = \emptyset$ and consider running the following procedure for every $M' \in \mathcal{M}_{\text{cov}}$:

1. Choose a single $M \in \mathcal{M}'$ such that $\sup_{r,o,\pi : (r,o,\pi) \in \mathcal{E}} \left| \log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{M',\pi}(r,o)} \right| \leq \rho/2$ (if such an M exists).
2. If there exists an $M \in \mathcal{M}'$ in step 1, set $\mathcal{M}'_{\text{cov}} \leftarrow \mathcal{M}'_{\text{cov}} \cup \{M\}$. Otherwise $\mathcal{M}'_{\text{cov}}$ remains unchanged.

By construction $\mathcal{M}'_{\text{cov}} \subseteq \mathcal{M}'$, and $|\mathcal{M}'_{\text{cov}}| \leq |\mathcal{M}_{\text{cov}}|$. We claim that $\mathcal{M}'_{\text{cov}}$ is a (ρ, μ) -cover of \mathcal{M}' . To see why, take some $M \in \mathcal{M}'$. Let $M' \in \mathcal{M}_{\text{cov}}$ denote the point realizing [Eq. \(69\)](#) for M . Let

M'' denote the point chosen in the above procedure for M' . Note that there must exist some M'' chosen for this M' since Eq. (69) holds for M , so in particular $M \in \mathcal{M}'$ satisfies the condition of step 1 in the above procedure. Then,

$$\begin{aligned} \sup_{r,o,\pi : (r,o,\pi) \in \mathcal{E}} \left| \log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{M''}(o|\pi)} \right| &= \sup_{r,o,\pi : (r,o,\pi) \in \mathcal{E}} \left| \log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{M',\pi}(r,o)} + \log \frac{\mathbb{P}^{M',\pi}(r,o)}{\mathbb{P}^{M''}(o|\pi)} \right| \\ &\leq \sup_{r,o,\pi : (r,o,\pi) \in \mathcal{E}} \left| \log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{M',\pi}(r,o)} \right| + \sup_{r,o,\pi : (r,o,\pi) \in \mathcal{E}} \left| \log \frac{\mathbb{P}^{M''}(o|\pi)}{\mathbb{P}^{M',\pi}(r,o)} \right| \\ &\leq \rho/2 + \rho/2 = \rho \end{aligned}$$

where the last inequality follows by our choice of M' and the definition of M'' . Thus, it follows that $\mathcal{M}'_{\text{cov}}$ is a (ρ, μ) -cover of \mathcal{M}' . \square

F.4.3. FURTHER LEMMAS

Lemma F.16. *For $a > 0$, $\alpha \leq [0, 3/4]$, and $\beta > 0$, the function $x^\alpha \log^\beta(ax)$ is concave in x for $x \geq \frac{1}{a} \exp\left(\frac{4\beta}{\alpha}\right)$ when $\beta \leq 1$, and for $x \geq \max\left\{\frac{1}{a} \exp\left(\sqrt{\frac{8(\beta-1)\beta}{\alpha}}\right), \frac{1}{a} \exp\left(\frac{4\beta}{\alpha}\right)\right\}$ when $\beta > 1$.*

Proof of Lemma F.16. By some calculation, we have

$$\begin{aligned} \frac{d^2}{dx^2} \left(x^\alpha \log^\beta(ax) \right) &= (-1 + \beta)\beta x^{-2+\alpha} \log^{\beta-2}(ax) + (-\beta + 2\alpha\beta)x^{-2+\alpha} \log^{\beta-1}(ax) \\ &\quad + (-1 + \alpha)\alpha x^{-2+\alpha} \log^\beta(ax). \end{aligned}$$

If we restrict to $x \geq 1/a$, then to show that the function is concave it then suffices to show that

$$(-1 + \beta)\beta \log^{-2}(ax) + (-\beta + 2\alpha\beta) \log^{-1}(ax) + (-1 + \alpha)\alpha \leq 0$$

which, since $\alpha \leq 3/4$, is implied by

$$(-1 + \beta)\beta \log^{-2}(ax) + \frac{1}{2}\beta \log^{-1}(ax) \leq \frac{1}{4}\alpha$$

which is further implied by

$$(-1 + \beta)\beta \log^{-2}(ax) \leq \frac{1}{8}\alpha \quad \text{and} \quad \frac{1}{2}\beta \log^{-1}(ax) \leq \frac{1}{8}\alpha.$$

The former condition is met for $x > 1/a$ for all $\beta \in (0, 1]$. For $\beta > 1$, it is met as long as

$$\frac{8(\beta-1)\beta}{\alpha} \leq \log^2(ax) \iff x \geq \frac{1}{a} \exp\left(\sqrt{\frac{8(\beta-1)\beta}{\alpha}}\right).$$

The latter condition is met for

$$x \geq \frac{1}{a} \exp\left(\frac{4\beta}{\alpha}\right).$$

\square

Lemma F.17. For all $B, C, n \geq 1$, if $x \leq C \log^n(Bx)$, then $x \leq C(2n)^n \log^n(2nBC)$.

Proof of Lemma F.17. This is a direct consequence of Lemma A.4 of [Wagenmaker et al. \(2022a\)](#). \square

Appendix G. Proofs for Examples

In this section, we provide proofs for the examples given in [Appendix A](#). We begin in [Appendix G.1](#) by introducing a condition which implies $n_\varepsilon^{M^*}$ is bounded, and is easy to verify for many classes of interest. Next, in [Appendix G.2](#), we consider a variety of structured bandit settings, and in [Appendix G.3](#) extend this to contextual bandits with finitely many arms. In [Appendix G.4](#), we provide proofs for the informative arm example of [Section 1.3](#). Finally, in [Appendix G.5](#), we consider tabular MDPs.

G.1. Preliminaries: Regular Models

To bound the quantity $n_\varepsilon^{M^*} = n_\varepsilon^{M^*}$ for the examples we consider, it will be helpful to introduce the following notion of a *regular model*.

Definition G.1 (Regular Model). We say instance $M \in \mathcal{M}$ is a regular model if there exists some constant $L_{\mathcal{M}}^M > 0$ such that, for any $M' \in \mathcal{M}^{\text{alt}}(M)$ with $D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) > 0$, there exists $M'' \in \mathcal{M}^{\text{alt}}(M)$ such that $D_{\text{KL}}(M(\pi_M) \| M''(\pi_M)) = 0$ and, for all $\pi \in \Pi$,

$$|D_{\text{KL}}(M(\pi) \| M'(\pi)) - D_{\text{KL}}(M(\pi) \| M''(\pi))| \leq \sqrt{L_{\mathcal{M}}^M D_{\text{KL}}(M(\pi_M) \| M'(\pi_M))} + L_{\mathcal{M}}^M D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)). \quad (70)$$

Our definition of a regular model is a direct generalization of existing notions of class regularity found in the literature ([Degenne et al., 2020b](#)). As we will see, for a variety of standard bandit classes (including multi-armed bandits, linear bandits, and Lipschitz bandits), as well as tabular MDPs, one can show that M^* is a regular model with $L_{\mathcal{M}}^{M^*} = L_{\mathcal{M}}^M$ bounded by a polynomial function of problem parameters. Intuitively, M^* will be a regular model if, for any instance $M' \in \mathcal{M}^{\text{alt}}(M^*)$ for which it is sufficient to pull π_* in order to distinguish M^* and M' (thereby ruling out M' while incurring no regret), then there exists some other instance $M'' \in \mathcal{M}^{\text{alt}}(M^*)$ which is “close” to M' in a certain sense, and which cannot be distinguished from M^* by simply pulling π_* . As the following result shows, the quantity $n_\varepsilon^{M^*}$ can be bounded whenever M^* is a regular model.

Proposition G.1. If M is a regular model with $\Delta_{\min}^M > 0$, we can bound

$$n_\varepsilon^M \leq \frac{2g^M}{\Delta_{\min}^M} \cdot \left(1 + L_{\mathcal{M}}^M + \frac{2g^M}{\varepsilon \Delta_{\min}^M} \cdot L_{\mathcal{M}}^M \right).$$

Given [Proposition G.1](#), for many of the examples in this section, rather than bounding $n_\varepsilon^{M^*}$ directly, we first show that M^* is a regular model with $L_{\mathcal{M}}^{M^*}$ well-bounded, and then use [Proposition G.1](#) to obtain a bound on $n_\varepsilon^{M^*}$.

Proof of Proposition G.1. To prove this result, it suffices to show that, for every normalized allocation $\lambda \in \Lambda(M, \varepsilon)$ with normalization factor n , there exists some allocation $\eta \in \mathbb{R}_+^\Pi$ such that 1)

$\eta(\pi) = n\lambda(\pi)$ for $\pi \neq \pi_M$, and $\eta(\pi_M) \leq \bar{n}$ for some well-bounded \bar{n} , and 2) $\Delta^M(\eta) \leq (1 + 2\varepsilon)\mathbf{g}^M$ and $I^M(\eta) \geq 1 - 2\varepsilon$.

Fix some $\lambda \in \Lambda(M, \varepsilon)$ with normalization factor $n > 0$. Note that if M is a regular model, then $I^M(\mathbb{I}_{\pi_M}) = 0$. Since $\lambda \in \Lambda(M, \varepsilon)$, this implies that $\lambda(\pi_M) < 1$. Let λ' denote the allocation $\lambda'(\pi_M) = 1 - \zeta$ and $\lambda'(\pi) = \frac{\zeta}{1 - \lambda(\pi_M)}\lambda(\pi)$ for $\pi \neq \pi_M$, for some ζ to be chosen.

We have

$$\Delta^M(\lambda') = \frac{\zeta}{1 - \lambda(\pi_M)}\Delta^M(\lambda) \leq \frac{\zeta}{1 - \lambda(\pi_M)}(1 + \varepsilon)\mathbf{g}^M/n. \quad (71)$$

Take $M' \in \mathcal{M}^{\text{alt}}(M)$ such that $D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) > 0$ and let $M'' \in \mathcal{M}^{\text{alt}}(M)$ denote the instance guaranteed to exist under [Definition G.1](#). We then have that, for any π ,

$$\begin{aligned} D_{\text{KL}}(M(\pi) \| M'(\pi)) &\geq D_{\text{KL}}(M(\pi) \| M''(\pi)) - \sqrt{L_{\mathcal{M}}^M D_{\text{KL}}(M(\pi_M) \| M'(\pi_M))} - L_{\mathcal{M}}^M D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) \\ &\geq D_{\text{KL}}(M(\pi) \| M''(\pi)) - (1 + \alpha)L_{\mathcal{M}}^M D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) - \frac{1}{\alpha} \end{aligned}$$

where the last inequality follows for any $\alpha > 0$ by AM-GM. Then

$$\begin{aligned} &\sum_{\pi} \lambda'(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) \\ &= \sum_{\pi \neq \pi_M} \frac{\zeta}{1 - \lambda(\pi_M)} \lambda(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) + (1 - \zeta) D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) \\ &\geq \sum_{\pi \neq \pi_M} \frac{\zeta \lambda(\pi)}{1 - \lambda(\pi_M)} \left(D_{\text{KL}}(M(\pi) \| M''(\pi)) - (1 + \alpha)L_{\mathcal{M}}^M D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) - \frac{1}{\alpha} \right) \\ &\quad + (1 - \zeta) D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) \\ &= \sum_{\pi \neq \pi_M} \frac{\zeta}{1 - \lambda(\pi_M)} \lambda(\pi) D_{\text{KL}}(M(\pi) \| M''(\pi)) + \left((1 - \zeta) - (1 + \alpha)L_{\mathcal{M}}^M \zeta \right) D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) - \frac{\zeta}{\alpha}. \end{aligned}$$

We have $M'' \in \mathcal{M}^{\text{alt}}(M)$, so by definition

$$\sum_{\pi} \lambda(\pi) D_{\text{KL}}(M(\pi) \| M''(\pi)) \geq (1 - \varepsilon)/n.$$

However, $D_{\text{KL}}(M(\pi_M) \| M''(\pi_M)) = 0$ by assumption, so it follows that

$$\sum_{\pi \neq \pi_M} \frac{\zeta}{1 - \lambda(\pi_M)} \lambda(\pi) D_{\text{KL}}(M(\pi) \| M''(\pi)) \geq \frac{\zeta}{1 - \lambda(\pi_M)} \cdot \frac{1 - \varepsilon}{n}.$$

By assumption, $\Delta^M(\lambda) \leq (1 + \varepsilon)\mathbf{g}^M/n$, and we can also lower bound $\Delta^M(\lambda) \geq (1 - \lambda(\pi_M))\Delta_{\min}^M$. Rearranging these implies that

$$(1 - \lambda(\pi_M))n \leq (1 + \varepsilon)\mathbf{g}^M/\Delta_{\min}^M.$$

Set $\alpha = \frac{(1 + \varepsilon)\mathbf{g}^M}{\varepsilon\Delta_{\min}^M}$, then we have

$$\frac{\zeta}{1 - \lambda(\pi_M)} \cdot \frac{1 - \varepsilon}{n} - \frac{\zeta}{\alpha} = \frac{\zeta}{1 - \lambda(\pi_M)} \cdot \frac{1 - \varepsilon}{n} - \frac{\zeta\varepsilon}{(1 + \varepsilon)\mathbf{g}^M/\Delta_{\min}^M}$$

$$\begin{aligned} &\geq \frac{\zeta}{1 - \lambda(\pi_M)} \cdot \frac{1 - \varepsilon}{n} - \frac{\zeta \varepsilon}{(1 - \lambda(\pi_M))n} \\ &= \frac{\zeta}{1 - \lambda(\pi_M)} \cdot \frac{1 - 2\varepsilon}{n}. \end{aligned}$$

Furthermore, with this choice of α , if

$$\zeta \leq \frac{1}{1 + (1 + \alpha)L_{\mathcal{M}}^M} = \frac{1}{1 + (1 + (1 + \varepsilon)\mathbf{g}^M/\varepsilon\Delta_{\min}^M)L_{\mathcal{M}}^M},$$

we have

$$\left((1 - \zeta) - (1 + \alpha)L_{\mathcal{M}}^M\zeta \right) D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) \geq 0.$$

We therefore have that, for any $M' \in \mathcal{M}^{\text{alt}}(M)$ with $D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) > 0$, that

$$\sum_{\pi} \lambda'(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) \geq \frac{\zeta}{1 - \lambda(\pi_M)} \cdot \frac{1 - 2\varepsilon}{n}.$$

Now consider $M' \in \mathcal{M}^{\text{alt}}(M)$ with $D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) = 0$. In this case we have

$$\sum_{\pi} \lambda'(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) = \frac{\zeta}{1 - \lambda(\pi_M)} \cdot \sum_{\pi} \lambda(\pi) D_{\text{KL}}(M(\pi) \| M'(\pi)) \geq \frac{\zeta}{1 - \lambda(\pi_M)} \cdot \frac{1 - \varepsilon}{n}$$

where the inequality follows since $\lambda \in \Lambda(M, \varepsilon)$ with normalization factor n by assumption. Together these bounds then imply that:

$$I^M(\lambda') \geq \frac{\zeta}{1 - \lambda(\pi_M)} \cdot \frac{1 - 2\varepsilon}{n}.$$

Combining this with our bound on $\Delta^M(\lambda')$ in (71) implies that $\lambda' \in \Lambda(M; 2\varepsilon)$ with parameter $n' = \frac{1 - \lambda(\pi_M)}{\zeta} \cdot n$.

To conclude, define the allocation $\eta := n'\lambda'$. Then for $\pi \neq \pi_M$:

$$\eta(\pi) = \frac{1 - \lambda(\pi_M)}{\zeta} \cdot n\lambda'(\pi) = n\lambda(\pi)$$

and

$$\eta(\pi_M) = \frac{1 - \lambda(\pi_M)}{\zeta} \cdot n\lambda'(\pi_M) \leq \frac{1 - \lambda(\pi_M)}{\zeta} \cdot n$$

It follows then that $\Delta^M(\eta) = n\Delta^M(\lambda) \leq (1 + \varepsilon)\mathbf{g}^M$, and $I^M(\eta) \geq n'I^M(\lambda') \geq 1 - 2\varepsilon$. Therefore, η satisfies the desired condition. Since $\lambda \in \Lambda(M, \varepsilon)$, we have

$$\Delta^M(\lambda) \leq (1 + \varepsilon)\mathbf{g}^M/n \implies (1 - \lambda(\pi_M))n \leq (1 + \varepsilon)\mathbf{g}^M/\Delta_{\min}^M,$$

and thus

$$\eta(\pi_M) \leq \frac{(1 + \varepsilon)\mathbf{g}^M}{\Delta_{\min}^M \cdot \zeta} \leq \frac{2\mathbf{g}^M(1 + (1 + 2\mathbf{g}^M/\varepsilon\Delta_{\min}^M)L_{\mathcal{M}}^*)}{\Delta_{\min}^M}.$$

which proves the result. \square

G.2. Structured Bandits with Gaussian Noise

In this section, we consider the problem of structured bandits with Gaussian noise, in which $\mathcal{O} = \{\emptyset\}$, and the mean reward functions belong to a given function class \mathcal{F} . Concretely, we consider the model class

$$\mathcal{M} = \{M(\pi) = \mathcal{N}(f(\pi), \sigma^2) : f \in \mathcal{F}\}.$$

We set

$$D(M(\pi) \parallel \bar{M}(\pi)) \leftarrow D_{\text{KL}}(M(\pi) \parallel \bar{M}(\pi)) = \frac{1}{2\sigma^2}(f^M(\pi) - f^{\bar{M}}(\pi))^2 \quad (72)$$

for D the divergence used by AE_*^2 . In general in the following examples we take $\sigma = 1$ for simplicity.

We begin by verifying that the basic regularity conditions required by our results are satisfied for generic classes \mathcal{F} , then provide bounds on the AEC for specific classes of interest.

Lemma G.1. *For bandits with Gaussian noise:*

1. For $D \leftarrow D_{\text{KL}}$, [Assumptions A.2](#), [D.1](#) and [D.2](#) hold with parameters

$$L_{\text{KL}} = V_{\mathcal{M}} = \frac{2\sqrt{2}}{\sigma}.$$

2. We can bound $N_{\text{cov}}(\mathcal{M}, \rho, \mu)$ by the covering number of \mathcal{M} in the distance $d(M, M') := \sup_{\pi \in \Pi} |f^M(\pi) - f^{M'}(\pi)|$ at tolerance $\frac{\sigma^2 \cdot \rho}{2 + \sqrt{2\sigma^2 \log(2/\mu)}}$. Furthermore, it suffices to take $\mathcal{E} := \{|r| \leq 1 + \sqrt{2\sigma^2 \log(2/\mu)}\}$.

Proof of Lemma G.1. In this setting, we have that for any $M, M' \in \mathcal{M}$ and any $\pi \in \Pi$,

$$D_{\text{KL}}(M(\pi) \parallel M'(\pi)) = \frac{1}{2\sigma^2}(f^M(\pi) - f^{M'}(\pi))^2.$$

For $\bar{M} \in \mathcal{M}$, we therefore have

$$\begin{aligned} |D_{\text{KL}}(M(\pi) \parallel M'(\pi)) - D_{\text{KL}}(\bar{M}(\pi) \parallel M'(\pi))| &= \frac{1}{2\sigma^2} \left| (f^M(\pi) - f^{M'}(\pi))^2 - (f^{\bar{M}}(\pi) - f^{M'}(\pi))^2 \right| \\ &\leq \frac{2}{\sigma^2} |f^M(\pi) - f^{\bar{M}}(\pi)| \\ &= \frac{2\sqrt{2}}{\sigma} \sqrt{D(M(\pi) \parallel \bar{M}(\pi))}, \end{aligned}$$

where the inequality follows from the Mean Value Theorem and the assumption that $f^{\bar{M}}(\pi) \in [0, 1]$ for all $\pi \in \Pi$. This verifies that [Assumption D.1](#) holds with $L_{\text{KL}} = \frac{2\sqrt{2}}{\sigma}$.

To show that [Assumption A.2](#) is met, we note that for all $M, \bar{M}, \bar{M}' \in \mathcal{M}$,

$$\begin{aligned} \log \frac{\mathbb{P}^{\bar{M}, \pi}(r, \mathcal{O})}{\mathbb{P}^{M, \pi}(r, \mathcal{O})} &= \frac{1}{2\sigma^2}(r - f^M(\pi))^2 - \frac{1}{2\sigma^2}(r - f^{\bar{M}}(\pi))^2 \\ &= \frac{1}{2\sigma^2} \left[f^M(\pi)^2 + f^{\bar{M}}(\pi)^2 - 2r(f^M(\pi) - f^{\bar{M}}(\pi)) \right] \end{aligned}$$

$$\begin{aligned}
 &= \frac{1}{2\sigma^2} \left[f^M(\pi)^2 + f^{\bar{M}}(\pi)^2 - 2f^{\bar{M}'}(\pi)(f^M(\pi) - f^{\bar{M}}(\pi)) + 2(f^{\bar{M}'}(\pi) - r)(f^M(\pi) - f^{\bar{M}}(\pi)) \right] \\
 &= \mathbb{E}_{(r,o) \sim \bar{M}'(\pi)} \left[\log \frac{\mathbb{P}^{\bar{M},\pi}(r,o)}{\mathbb{P}^{M,\pi}(r,o)} \right] + \frac{1}{\sigma^2} (f^{\bar{M}'}(\pi) - r)(f^M(\pi) - f^{\bar{M}}(\pi)).
 \end{aligned}$$

It follows that

$$\log \frac{\mathbb{P}^{\bar{M},\pi}(r,o)}{\mathbb{P}^{M,\pi}(r,o)} - \mathbb{E}_{(r,o) \sim \bar{M}'(\pi)} \left[\log \frac{\mathbb{P}^{\bar{M},\pi}(r,o)}{\mathbb{P}^{M,\pi}(r,o)} \right]$$

is sub-gaussian with parameter $\mathbb{E}_{r \sim \bar{M}'(\pi)} [(\frac{1}{\sigma^2} (f^{\bar{M}'}(\pi) - r)(f^M(\pi) - f^{\bar{M}}(\pi)))^2] \leq \frac{4}{\sigma^2}$, which verifies [Assumption A.2](#) with $V_{\mathcal{M}}^2 = \frac{8}{\sigma^2}$.

Finally, we bound the covering number. Let $\mathcal{E} := \{|r| \leq 1 + \sqrt{2\sigma^2 \log(2/\mu)}\}$. Elementary manipulations show that $\mathbb{P}^{\bar{M},\pi}(\mathcal{E}^c) \leq \mu$ for any $\bar{M} \in \mathcal{M}$ and π . Using the same calculation as above, we have

$$\log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{M',\pi}(r,o)} = \frac{1}{2\sigma^2} (r - f^M(\pi))^2 - \frac{1}{2\sigma^2} (r - f^{M'}(\pi))^2 \leq \frac{1 + |r|}{\sigma^2} \cdot |f^M(\pi) - f^{M'}(\pi)|$$

where the inequality follows from the Mean Value Theorem. We therefore have that for any $M, M' \in \mathcal{M}$,

$$\sup_{r,o,\pi : |r| \leq 1 + \sqrt{2\sigma^2 \log(2/\mu)}} \left| \log \frac{\mathbb{P}^{M,\pi}(r,o)}{\mathbb{P}^{M',\pi}(r,o)} \right| \leq \frac{2 + \sqrt{2\sigma^2 \log(2/\mu)}}{\sigma^2} \cdot \sup_{\pi} |f^M(\pi) - f^{M'}(\pi)|.$$

It follows that if we can form a $\frac{\sigma^2 \cdot \rho}{2 + \sqrt{2\sigma^2 \log(2/\mu)}}$ -cover of \mathcal{M} in the distance $d(M, M') = \sup_{\pi \in \Pi} |f^M(\pi) - f^{M'}(\pi)|$, this will serve as an (ρ, μ) cover of \mathcal{M} . \square

G.2.1. DISCRETE STRUCTURED BANDITS

As a first example of bandits with Gaussian noise, we present an additional class that satisfies the uniformly regular assumption.

Example G.1 (Discrete Structured Bandits). Fix $\Delta_{\min} > 0$, and consider a discrete reward space $\mathcal{R} \subseteq [0, 1]$ satisfying $\min_{r,r' \in \mathcal{R}} |r - r'| \geq \Delta_{\min}$. Consider any function class $\mathcal{F} \subseteq (\Pi \rightarrow \mathcal{R})$ defined such that each $f \in \mathcal{F}$ has a unique optimal decision. Let our model class be defined as

$$\mathcal{M} = \{M(\pi) = \mathcal{N}(f(\pi), 1) \mid f \in \mathcal{F}\}.$$

It is straightforward to show that [Assumption A.1](#) and [Assumption A.2](#) are met with $L_{\text{KL}}, V_{\mathcal{M}} \leq 4$, and that [Assumption A.3](#) is satisfied with d_{cov} scaling with the log-covering number of \mathcal{F} in the distance $d(f, f') = \sup_{\pi \in \Pi} |f(\pi) - f'(\pi)|$, and $C_{\text{cov}} = O(1)$. Furthermore, [Assumption A.4](#) is satisfied by construction of \mathcal{R} and \mathcal{F} , and we can bound $\eta_{\varepsilon}^{\mathcal{M}} \leq \frac{2}{\Delta_{\min}^2}$.⁸ We thus have the following corollary to [Theorem A.1](#).

8. It is not difficult to see that, given the construction of \mathcal{R} , once the optimal arm has been played $\frac{2}{\Delta_{\min}^2}$ times, no additional information can be extracted from playing it.

Corollary G.1. *In the discrete structured bandits setting considered above, if $\mathbf{g}^* > 0$, the regret of AE^2 is bounded as*

$$\mathbb{E}^{M^*}[\mathbf{Reg}(T)] \leq (1 + \varepsilon)\mathbf{g}^* \cdot \log(T) + \overline{\text{aec}}_{\varepsilon/12}(\mathcal{M}) \cdot \text{poly}\left(\max_{M \in \mathcal{M}} \mathbf{g}^M, d_{\text{cov}}, \frac{1}{\varepsilon}, \frac{1}{\Delta_{\min}}, \log \log T\right) \cdot \log^{1/2}(T)$$

As we show in [Example A.3](#), the AEC in this setting can be bounded in terms of the eluder dimension of \mathcal{F} . \triangleleft

Proof for Example G.1. We provide calculations for the discrete structured bandits setting of [Example G.1](#). First, note that [Assumptions A.1](#) to [A.3](#) all hold by [Lemma G.1](#). To bound $n^{\mathcal{M}}$, consider $M \in \mathcal{M}$ and $M' \in \mathcal{M}^{\text{alt}}(M)$. Note that either $D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) = 0$, in which case there is no advantage to playing π_M , or, due to the discretization of the means, $D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) \geq \frac{1}{2}\Delta_{\min}^2$. Thus for any allocation $\eta \in \mathbb{R}_+^{\Pi}$, as long as $\eta(\pi_M) \geq \frac{2}{\Delta_{\min}^2}$, we have $\eta(\pi_M)D_{\text{KL}}(M(\pi_M) \| M'(\pi_M)) \geq 1$. It follows that there is no advantage to choosing $\eta(\pi_M)$ larger than $\frac{2}{\Delta_{\min}^2}$, so we can bound $n^{\mathcal{M}} \leq \frac{2}{\Delta_{\min}^2}$. \square

G.2.2. MULTI-ARMED BANDITS (EXAMPLE A.2)

In this section we prove the result in [Example A.2](#). First, note that for any $M \in \mathcal{M}$, we have $\mathbf{g}^M \leq A/\Delta_{\min}^M$. [Assumptions A.2](#), [D.1](#) and [D.2](#) are met due to [Lemma G.1](#), and with constants $L_{\text{KL}}, V_{\mathcal{M}} \leq 4$, $d_{\text{cov}} = O(A)$, and $C_{\text{cov}} = O(1)$. By [Lemma G.2](#)—stated and proven below— M^* is a regular model with $L_{\mathcal{M}}^* = \sqrt{2}$ as long as $f^{M^*}(\pi_*) < 1$. It remains to bound the AEC.

Proof of Proposition A.2. It is immediate to see that $C_{\text{exp}}(\mathcal{M}, \varepsilon) \leq O(\frac{A}{\varepsilon})$ by choosing the exploration distribution to be uniform over A . By [Lemma E.6](#), we then

$$\overline{\text{aec}}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}^*) \leq c_1 \cdot \frac{A^3}{\varepsilon^2 \Delta_{\star}^6}$$

for a universal constant c_1 . By [Proposition G.1](#), we have

$$n_{\varepsilon/36}^* \leq c_2 \cdot \frac{\mathbf{g}^*}{\Delta_{\min}^*} \cdot \left(1 + \frac{\mathbf{g}^*}{\varepsilon \Delta_{\min}^*}\right) \leq c_2 \cdot \frac{A^2}{\varepsilon (\Delta_{\min}^*)^4},$$

for a universal constant c_2 , so we can lower bound $\Delta_{\star} \geq c_3 \varepsilon (\Delta_{\min}^*)^4 / A^2$, giving

$$\overline{\text{aec}}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}^*) \leq c_4 \cdot \frac{A^{15}}{\varepsilon^8 (\Delta_{\min}^*)^{24}}.$$

\square

Lemma G.2. *In the multi-armed bandit setting of [Example A.2](#), any model $M^* \in \mathcal{M}$ is a regular model with $L_{\mathcal{M}}^* = \sqrt{3}$ as long as $f^{M^*}(\pi_*) < 1$.*

Proof of Lemma G.2. Let $M' \in \mathcal{M}^{\text{alt}}(M^*)$ and assume that $D_{\text{KL}}(M^*(\pi_*) \| M'(\pi_*)) > 0$.

Case 1: $f^{M'}(\pi_{M'}) + f^{M^*}(\pi_*) - f^{M'}(\pi_*) \leq 1$. Let M'' denote the Gaussian bandit instance given by $f^{M''}(\pi) = f^{M'}(\pi) + f^{M^*}(\pi_*) - f^{M'}(\pi_*)$; by our assumption, $M'' \in \mathcal{M}$, and $M'' \in \mathcal{M}^{\text{alt}}(M^*)$. Furthermore, $f^{M''}(\pi_*) = f^{M^*}(\pi_*)$ which implies $D_{\text{KL}}(M^*(\pi_*) \| M''(\pi_*)) = 0$ as desired. Finally, for any π , we have

$$f^{M'}(\pi) - f^{M''}(\pi) = f^{M^*}(\pi_*) - f^{M'}(\pi_*).$$

This implies that for all π , since all models in \mathcal{M} have unit Gaussian rewards, using the expression for the KL divergence given in Eq. (72):

$$|D_{\text{KL}}(M^*(\pi) \| M'(\pi)) - D_{\text{KL}}(M^*(\pi) \| M''(\pi))| \leq |f^{M^*}(\pi_*) - f^{M'}(\pi_*)| = \sqrt{2D_{\text{KL}}(M^*(\pi_*) \| M'(\pi_*))}$$

which implies the condition of Definition G.1 is met with $L_{\mathcal{M}}^* = \sqrt{2}$.

Case 2: $f^{M'}(\pi_{M'}) + f^{M^*}(\pi_*) - f^{M'}(\pi_*) > 1$. For this case the model M'' constructed in Case 1 will not be in \mathcal{M} . Assume first that $f^{M^*}(\pi_*) \geq f^{M'}(\pi_{M'})$ and in this case define M'' to be the instance

$$f^{M''}(\pi) = \begin{cases} \min\{f^{M^*}(\pi_*), f^{M'}(\pi) + f^{M^*}(\pi_*) - f^{M'}(\pi_*)\} & \pi \neq \pi_{M'} \\ f^{M^*}(\pi_*) + \delta & \pi = \pi_{M'} \end{cases}$$

or some $\delta > 0$ such that $f^{M^*}(\pi_*) + \delta < 1$ (note that such a δ exists since we have assumed $f^{M^*}(\pi_*) < 1$). Note that we now have $M'' \in \mathcal{M}$, and $f^{M''}(\pi) < f^{M''}(\pi_{M'})$ for all $\pi \neq \pi_{M'}$, so $M'' \in \mathcal{M}^{\text{alt}}(M^*)$. Furthermore, we have $f^{M''}(\pi_*) = f^{M^*}(\pi_*)$, so $D_{\text{KL}}(M^*(\pi_*) \| M''(\pi_*)) = 0$. For $\pi \neq \pi_{M'}$, if $\min\{f^{M^*}(\pi_*), f^{M'}(\pi) + f^{M^*}(\pi_*) - f^{M'}(\pi_*)\} = f^{M^*}(\pi_*)$, this implies that

$$f^{M^*}(\pi_*) \leq f^{M'}(\pi) + f^{M^*}(\pi_*) - f^{M'}(\pi_*) \implies f^{M^*}(\pi_*) - f^{M'}(\pi) \leq f^{M^*}(\pi_*) - f^{M'}(\pi_*).$$

Since we have assumed $f^{M^*}(\pi_*) \geq f^{M'}(\pi_{M'})$, this implies that

$$|f^{M''}(\pi) - f^{M'}(\pi)| = |f^{M^*}(\pi_*) - f^{M'}(\pi)| \leq |f^{M^*}(\pi_*) - f^{M'}(\pi_*)|$$

So by the expression given for the KL divergence in Eq. (72), we have

$$|D_{\text{KL}}(M^*(\pi) \| M'(\pi)) - D_{\text{KL}}(M^*(\pi) \| M''(\pi))| \leq \sqrt{3D_{\text{KL}}(M^*(\pi_*) \| M'(\pi_*))}. \quad (73)$$

For $\pi \neq \pi_{M'}$ with $\min\{f^{M^*}(\pi_*), f^{M'}(\pi) + f^{M^*}(\pi_*) - f^{M'}(\pi_*)\} = f^{M'}(\pi) + f^{M^*}(\pi_*) - f^{M'}(\pi_*)$, the bound on $|D_{\text{KL}}(M^*(\pi) \| M'(\pi)) - D_{\text{KL}}(M^*(\pi) \| M''(\pi))|$ follows identically to Case 1. For $\pi = \pi_{M'}$, since we have assumed that $f^{M^*}(\pi_*) \geq f^{M'}(\pi_{M'})$ we have

$$|f^{M''}(\pi_{M'}) - f^{M'}(\pi_{M'})| = f^{M^*}(\pi_*) - f^{M'}(\pi_{M'}) + \delta \leq f^{M^*}(\pi_*) - f^{M'}(\pi_*) + \delta.$$

For small enough δ , this implies that Eq. (73) is satisfied for $\pi = \pi_{M'}$ as well.

Consider now the case where $f^{M^*}(\pi_*) < f^{M'}(\pi_{M'})$. In this case define M'' by

$$f^{M''}(\pi) = \begin{cases} \min\{f^{M'}(\pi_{M'}) - \delta, f^{M'}(\pi) + f^{M^*}(\pi_*) - f^{M'}(\pi_*)\} & \pi \neq \pi_{M'} \\ f^{M'}(\pi_{M'}) & \pi = \pi_{M'} \end{cases}$$

for $\delta > 0$ small enough that $f^{M'}(\pi_{M'}) - \delta > f^{M^*}(\pi_*)$. Note that M'' and $M''' \in \mathcal{M}^{\text{alt}}(M^*)$ by construction, and that $f^{M''}(\pi_*) = f^{M^*}(\pi_*)$ by our choice of δ , so $D_{\text{KL}}(M^*(\pi_*) \| M''(\pi_*)) = 0$. For $\pi \neq \pi_{M'}$, if $\min\{f^{M'}(\pi_{M'}) - \delta, f^{M'}(\pi) + f^{M^*}(\pi_*) - f^{M'}(\pi_*)\} = f^{M'}(\pi_{M'}) - \delta$, then we have

$$\begin{aligned} |f^{M''}(\pi) - f^{M'}(\pi)| &= |f^{M'}(\pi_{M'}) - \delta - f^{M'}(\pi)| \\ &\leq f^{M'}(\pi_{M'}) - f^{M'}(\pi) + \delta \\ &\leq |f^{M^*}(\pi_*) - f^{M'}(\pi_*)| + \delta \end{aligned}$$

where for the final inequality we have used that $\min\{f^{M'}(\pi_{M'}) - \delta, f^{M'}(\pi) + f^{M^*}(\pi_*) - f^{M'}(\pi_*) = f^{M'}(\pi_{M'}) - \delta$. It follows that [Eq. \(73\)](#) is satisfied for this π for sufficiently small δ . If we instead have $\min\{f^{M'}(\pi_{M'}) - \delta, f^{M'}(\pi) + f^{M^*}(\pi_*) - f^{M'}(\pi_*)\} = f^{M'}(\pi) + f^{M^*}(\pi_*) - f^{M'}(\pi_*)$, then the bound on $|D_{\text{KL}}(M^*(\pi) \| M'(\pi)) - D_{\text{KL}}(M^*(\pi) \| M''(\pi))|$ follows identically to Case 1.

For $\pi = \pi_{M'}$, we have $|f^{M''}(\pi_{M'}) - f^{M'}(\pi_{M'})| = 0$. This proves the result. \square

G.2.3. STRUCTURED BANDITS WITH BOUNDED ELUDER DIMENSION (EXAMPLE A.3)

In this section, we give generic bounds on the uniform exploration coefficient and Allocation-Estimation Coefficient for structured bandit classes with bounded eluder dimension (cf. [Definition A.5](#)). These results are used by subsequent examples, including linear bandits.

Lemma G.3. *Let $\mathcal{M}\{M(\pi) = \mathcal{N}(f(\pi), 1) \mid f \in \mathcal{F}\}$. Then for all $\varepsilon > 0$, we have*

$$C_{\text{exp}}(\mathcal{M}, \varepsilon) \leq \frac{16d_{\text{E}}(\mathcal{F}, \sqrt{\varepsilon}/2)}{\varepsilon}.$$

Proof of Lemma G.3. Let $\xi \in \Delta(\mathcal{M})$. Recall the expression for KL divergence between Gaussians of unit variance:

$$\mathbb{E}_{\bar{M} \sim \xi}[\mathbb{E}_p[D_{\text{KL}}(\bar{M}(\pi) \| M(\pi))] = \frac{1}{2} \mathbb{E}_{\bar{M} \sim \xi}[\mathbb{E}_p[(f^M(\pi) - f^{\bar{M}}(\pi))^2]].$$

Abbreviate $d_{\text{E}} := d_{\text{E}}(\mathcal{F}, \sqrt{\varepsilon}/2)$ and let $\{\pi_1, \dots, \pi_{d_{\text{E}}}\}$ denote a maximal sequence of ε -independent points. By the definition of the eluder dimension, for any $\pi \in \Pi$ and any $M, \bar{M} \in \mathcal{M}$, we have:

$$\sqrt{\sum_{i=1}^{d_{\text{E}}} (f^M(\pi_i) - f^{\bar{M}}(\pi_i))^2} \leq \sqrt{\varepsilon/2} \implies |f^M(\pi) - f^{\bar{M}}(\pi)| \leq \sqrt{\varepsilon/2}.$$

Now, set p to be the uniform distribution over $\{\pi_1, \dots, \pi_{d_{\text{E}}}\}$. Assume that $M, M' \in \mathcal{M}$ are such that

$$\max_{M'' \in \{M, M'\}} \mathbb{E}_{\bar{M} \sim \xi} \mathbb{E}_p[D_{\text{KL}}(\bar{M}(\pi) \| M''(\pi))] = \frac{1}{2} \max_{M'' \in \{M, M'\}} \mathbb{E}_{\bar{M} \sim \xi} \mathbb{E}_p[(f^{M''}(\pi) - f^{\bar{M}}(\pi))^2] \leq \varepsilon/(16d_{\text{E}}),$$

Markov's inequality implies that for each $M'' \in \{M, M'\}$, with probability at least 3/4 over the draw of $\bar{M} \sim \xi$,

$$\mathbb{E}_p[(f^{M''}(\pi) - f^{\bar{M}}(\pi))^2] \leq \varepsilon/(2d_{\text{E}}).$$

Taking a union bound, we conclude that with probability at least $1/2$ over the draw of $\bar{M} \sim \xi$,

$$\max_{M'' \in \{M, M'\}} \mathbb{E}_p[(f^{M''}(\pi) - f^{\bar{M}}(\pi))^2] \leq \varepsilon/(2d_E). \quad (74)$$

Going forward, let $\bar{M} \in \mathcal{M}$ be any model such that [Eq. \(74\)](#) holds; we have just proven that such a model exists. It follows from the maximality of π_1, \dots, π_{d_E} and the definition of p that for all $\tilde{\pi} \in \Pi$, and $M'' \in \mathcal{M}$,

$$\mathbb{E}_p[(f^{M''}(\pi) - f^{\bar{M}}(\pi))^2] \leq \varepsilon/(2d_E) \implies (f^{M''}(\tilde{\pi}) - f^{\bar{M}}(\tilde{\pi}))^2 \leq \varepsilon/2.$$

In particular, since this holds for both $M \in \{M', M''\}$, and since [Eq. \(74\)](#) holds, we have that for all π ,

$$\begin{aligned} D_{\text{KL}}(M'(\pi) \| M''(\pi)) &= \frac{1}{2}(f^{M'}(\pi) - f^{M''}(\pi))^2 \\ &\leq (f^{M'}(\pi) - f^{\bar{M}}(\pi))^2 + (f^{\bar{M}}(\pi) - f^{M''}(\pi))^2 \\ &\leq \varepsilon. \end{aligned}$$

As this is the condition required by [Definition A.4](#), it suffices to take $C_{\text{exp}}^\xi(\varepsilon) = 16d_E(\mathcal{F}, \sqrt{\varepsilon}/2)/\varepsilon$. Since this bound holds uniformly for all choices of ξ , the result follows. \square

Proof of [Proposition A.3](#). The bound

$$C_{\text{exp}}(\mathcal{M}^*, \delta) \leq \frac{16d_E(\mathcal{F}, \sqrt{\delta}/2)}{\delta}.$$

follows from [Lemma G.3](#), since $d_E(\mathcal{F}', \delta) \leq d_E(\mathcal{F}, \delta)$ for all $\mathcal{F}' \subseteq \mathcal{F}$.

By [Lemma E.6](#) we can bound $\overline{\text{aec}}_\varepsilon^{\mathcal{M}}(\mathcal{M}^*)$:

$$\overline{\text{aec}}_\varepsilon^{\mathcal{M}}(\mathcal{M}^*) \leq C_{\text{exp}}(\mathcal{M}^*, \delta) \quad \text{for} \quad \delta = \min_{M \in \mathcal{M}^*} \min \left\{ \min \left\{ \frac{1}{81L_{\text{KL}}}, \frac{\Delta_{\min}^M}{34V_{\mathcal{M}}} \right\}, \frac{\varepsilon}{2\mathbf{g}^M/\Delta_{\min}^M + \mathbf{n}_{\varepsilon/36}^M}, \frac{\Delta_{\min}^M}{3} \right\}^2.$$

[Lemma E.2](#) implies that for all $M \in \mathcal{M}^*$,

$$\mathbf{g}^M \leq C_{\text{exp}}(\mathcal{M}^*, \frac{1}{4}(\Delta_{\min}^M)^2) \leq \frac{64d_E(\mathcal{F}, \frac{1}{2}\Delta_{\min}^M)}{(\Delta_{\min}^M)^2} \leq \frac{64d_E(\mathcal{F}, \frac{1}{2}\Delta_*)}{\Delta_*^2}$$

where we have used that the eluder dimension increases as its scale ε decreases; by [Lemma G.1](#), it suffices to take $L_{\text{KL}} = V_{\mathcal{M}} = 2$. A sufficient value for δ is therefore

$$\delta = c \cdot \frac{\varepsilon^2 \Delta_*^8}{d_E(\mathcal{F}, \frac{1}{2}\Delta_*)^2}$$

The result follows. \square

G.2.4. LINEAR BANDITS (EXAMPLE A.4)

Proof of Proposition A.4. The result follows directly from [Proposition A.3](#), since it is known that linear bandits have eluder dimension which scales as $d_{\mathbb{E}}(\mathcal{F}, \varepsilon) = O(d \cdot \log 1/\varepsilon)$ ([Russo and Van Roy, 2013](#)). \square

In what follows, we prove [Proposition A.5](#), providing sufficient conditions under which it is possible to bound the regularity constant $L_{\mathcal{M}}^*$ (and hence n_{ε}^*) for linear bandits.

We begin with an geometric assumption on Θ , \mathcal{X} , and θ^* , which we will show ensures that $M^*(\pi) = \mathcal{N}(\langle x_{\pi}, \theta^* \rangle, 1)$ is a regular model. To state our condition, we denote, for any vectors x and y , we define x_y and $x_{\bar{y}}$ to be unique vectors satisfying $x = x_y + x_{\bar{y}}$ for $x_y \parallel y$ and $x_{\bar{y}} \perp y$.

Assumption G.1 (Regular Linear Bandits). *The sets Θ , \mathcal{X} and model parameter θ^* satisfy:*

1. Θ is a convex polytope.
2. For all $\theta \in \Theta$, we have that there exists some $\delta_{\theta} > 0$ such that $\{\theta' \in \mathbb{R}^d : \|\theta' - \theta\|_2 \leq \delta_{\theta}\} \subseteq \Theta$.
3. Letting $x^* \in \mathcal{X}$ denote the optimal action for θ^* , we have

$$\left\{ \theta \in \mathbb{R}^d : \|\theta - \theta^*\|_2 \leq \max_{x \in \mathcal{X}, x \neq x^*} \Delta^*(x) / \|x_{\bar{x}^*}\|_2 \right\} \subseteq \Theta.$$

The first two points above are quite mild. The primary restriction of [Assumption G.1](#) is Point 3, which requires that θ^* is located sufficiently far within the interior of Θ . Using [Assumption G.1](#) we can state the full version of [Proposition A.5](#).

Proposition G.2 (Full Version of [Proposition A.5](#)). *If Θ , \mathcal{X} , and θ^* satisfy [Assumption G.1](#), then n_{ε}^* is bounded by a polynomial function of d , $1/\Delta_{\min}^*$, $1/\varepsilon$, \mathbf{g}^* , and a geometry-dependent term scaling with the structure of \mathcal{X} and Θ .*

Remark G.1 (Comparison to Existing Work). *We remark that [Assumption G.1](#) is similar to the conditions required by existing works which achieve instance-optimality in linear bandits with polynomial lower-order terms ([Tirinzoni et al., 2020](#); [Kirschner et al., 2021](#)). Though neither of these works explicitly states such a condition, closer inspection of their analysis reveals it is indeed required. In particular, the proof of Lemma 1 of [Tirinzoni et al. \(2020\)](#) relies on a result from [Degenne et al. \(2020b\)](#) which shows that a condition analogous to [Definition G.1](#) is met for linear bandits. However, the proof given in [Degenne et al. \(2020b\)](#) appears to only hold when Θ is unbounded, or a condition such as [Assumption G.1](#) holds. As [Tirinzoni et al. \(2020\)](#) assumes that Θ is bounded, their results therefore only appear to hold if a condition similar to [Assumption G.1](#) also holds. Similarly, in the proof of Lemma 10 of [Kirschner et al. \(2021\)](#), it is assumed that for every arm $x \neq x_{\pi^*}$, there exists some instance in the alternate set with optimal arm x . To satisfy this condition, it appears that an assumption similar to [Assumption G.1](#) is required.*

Thus, while not stated explicitly in the existing literature, it therefore seems that all existing results which obtain reasonable lower-order terms require an assumption similar to [Assumption G.1](#). Removing this assumption (or showing it is necessary) is an interesting direction for future work.

Proof of Proposition G.2. Under Assumption G.1, this follows directly from Lemma G.4 and Proposition G.1. \square

Lemma G.4. Under Assumption G.1, the linear bandit model M^* is regular for some $L_{\mathcal{M}}^* < \infty$ whose value depends on the geometry of Θ and \mathcal{X} .

Proof of Lemma G.4. Fix some $\theta^* \in \Theta$ and let x^* denote its optimal arm. Let $\Theta^{\text{alt}}(\theta^*) \subseteq \Theta$ denote parameters with optimal arm $x \neq x^*$. Assume there exists some $\theta \in \Theta^{\text{alt}}(\theta^*)$ such that $\langle \theta - \theta^*, x^* \rangle \neq 0$ (if this is not the case, M^* immediately satisfies Definition G.1 and we are done). Let $\Theta_x = \{\theta \in \Theta : \langle x, \theta \rangle \geq \langle x^*, \theta \rangle\}$, $\Theta^* = \{\theta \in \mathbb{R}^d : \langle \theta - \theta^*, x^* \rangle = 0\}$, and $\Theta_x^* = \Theta_x \cap \Theta^*$. We first show that, under Assumption G.1, $\Theta_x^* \neq \emptyset$ for all $x \in \mathcal{X}$, $x \neq x^*$. We then use this to show that $M^*(\pi) = \mathcal{N}(\langle x_\pi, \theta^* \rangle, 1)$ is a regular model.

Part 1: $\Theta_x^* \neq \emptyset$. Fix some $x \in \mathcal{X}$ with $x \neq x^*$. Consider $\theta = \theta^* + ax_{\bar{x}^*}$ for some $a \in \mathbb{R}$ to be chosen. By construction we have $\langle \theta, x^* \rangle = \langle \theta^*, x^* \rangle$, which implies that $\theta \in \Theta^*$ for all $a \in \mathbb{R}$. We wish to choose a large enough that $\langle \theta, x \rangle \geq \langle \theta, x^* \rangle$. Note that

$$\langle \theta, x \rangle = \langle \theta^*, x \rangle + a \langle x_{\bar{x}^*}, x_{\bar{x}^*} + x_{x^*} \rangle = \langle \theta^*, x \rangle + a \|x_{\bar{x}^*}\|_2^2$$

and $\langle \theta, x^* \rangle = \langle \theta^*, x^* \rangle$. Thus, to satisfy $\langle \theta, x \rangle \geq \langle \theta, x^* \rangle$, we need

$$a \|x_{\bar{x}^*}\|_2^2 \geq \langle \theta^*, x^* - x \rangle \iff a \geq \Delta^*(x) / \|x_{\bar{x}^*}\|_2^2.$$

Let $a = \Delta^*(x) / \|x_{\bar{x}^*}\|_2^2$, then it follows that $\langle \theta, x \rangle \geq \langle \theta, x^* \rangle$. Furthermore, we can bound

$$\|\theta - \theta^*\|_2 \leq a \|x_{\bar{x}^*}\|_2 = \Delta^*(x) / \|x_{\bar{x}^*}\|_2.$$

Under Assumption G.1, it follows that $\theta \in \Theta$.

Part 2: M^* is a Regular Model. Let $\bar{\Theta}_x = \{\theta - \theta^* : \theta \in \Theta_x\}$. Note that, since Θ is a convex polytope, and Θ_x simply adds a linear inequality constraint, $\bar{\Theta}_x$ is also convex. Let $\bar{\Theta}_x^* = \{\phi \in \bar{\Theta}_x : \langle \phi, x^* \rangle = 0\}$. From Part 1, we have $\bar{\Theta}_x^* \neq \emptyset$. Lemma 23 of Kirschner et al. (2021) then gives that there exists a geometry-dependent constant $C(\Theta, \mathcal{X})$ such that, for all $\phi \in \bar{\Theta}_x$:

$$\min_{\phi' \in \bar{\Theta}_x^*} \|\phi - \phi'\|_2 \leq C(\Theta, \mathcal{X}) \cdot |\langle \phi, x^* \rangle|.$$

This implies that for all $\theta \in \Theta_x$, we have:

$$\min_{\theta' \in \bar{\Theta}_x^*} \|\theta - \theta'\|_2 \leq C(\Theta, \mathcal{X}) \cdot |\langle \theta - \theta^*, x^* \rangle|.$$

Now consider some $\theta \in \Theta^{\text{alt}}(\theta^*)$, and assume that $\langle \theta - \theta^*, x^* \rangle \neq 0$ (by assumption such a θ exists). Assume that θ has optimal arm x , which implies that $\theta \in \Theta_x$. By what we have just shown, we know that there exists some $\theta' \in \Theta$ with $\langle \theta', x \rangle > \langle \theta', x^* \rangle$ so that $\theta' \in \Theta^{\text{alt}}(\theta^*)$, $\langle \theta' - \theta^*, x^* \rangle = 0$, and

$$\|\theta - \theta'\|_2 \leq C(\Theta, \mathcal{X}) \cdot |\langle \theta - \theta^*, x^* \rangle|.$$

Note that, for any $x' \in \mathcal{X}$, we have

$$|D_{\text{KL}}(\theta^*(x') \parallel \theta(x')) - D_{\text{KL}}(\theta^*(x') \parallel \theta'(x'))| = \frac{1}{2} |\langle \theta^* - \theta, x' \rangle^2 - \langle \theta^* - \theta', x' \rangle^2|$$

$$\begin{aligned} &\leq 2 \max_{\theta'' \in \Theta} |\langle \theta'', x' \rangle| \cdot |\langle \theta - \theta', x' \rangle| \\ &\leq 2 \max_{\theta'' \in \Theta} \|\theta''\|_2 \|x'\|_2^2 \cdot \|\theta - \theta'\|_2. \end{aligned}$$

Furthermore, note that

$$\sqrt{D_{\text{KL}}(\theta^*(x^*) \parallel \theta(x^*))} = \frac{1}{\sqrt{2}} |\langle \theta^* - \theta, x^* \rangle|,$$

so

$$\begin{aligned} &|D_{\text{KL}}(\theta^*(x') \parallel \theta(x')) - D_{\text{KL}}(\theta^*(x') \parallel \theta'(x'))| \\ &\leq \left(2\sqrt{2}C(\Theta, \mathcal{X}) \max_{\theta'' \in \Theta, x'' \in \mathcal{X}} \|\theta''\|_2 \|x''\|_2^2 \right) \cdot \sqrt{D_{\text{KL}}(\theta^*(x^*) \parallel \theta(x^*))}. \end{aligned}$$

As $\theta \in \Theta^*(\theta^*)$ was arbitrary, we have therefore shown that M^* is a regular model with

$$L_{\mathcal{M}}^* = \left(2\sqrt{2} \cdot C(\Theta, \mathcal{X}) \cdot \max_{\theta'' \in \Theta, x'' \in \mathcal{X}} \|\theta''\|_2 \|x''\|_2^2 \right)^2.$$

□

G.2.5. GENERALIZED LINEAR MODELS (EXAMPLE A.5)

Proof Sketch for Example A.5. The bound on the AEC follows as in [Proposition A.4](#), using that the eluder dimension for generalized linear models is bounded as $O(d \cdot (\frac{g_{\max}}{g_{\min}})^2 \cdot \log \frac{g_{\max}}{\varepsilon})$ ([Russo and Van Roy, 2013](#)). For the other regularity assumptions, note that by the Mean Value Theorem, we have

$$\begin{aligned} |D_{\text{KL}}(M(\pi) \parallel M'(\pi)) - D_{\text{KL}}(M(\pi) \parallel M''(\pi))| &= \frac{1}{2} |(g(\langle \theta, x \rangle) - g(\langle \theta', x \rangle))^2 - (g(\langle \theta, x \rangle) - g(\langle \theta'', x \rangle))^2| \\ &\leq 2g_{\max} |\langle \theta' - \theta'', x \rangle| \end{aligned}$$

and

$$\sqrt{D_{\text{KL}}(M(\pi) \parallel M'(\pi))} = \frac{1}{\sqrt{2}} |g(\langle \theta, x \rangle) - g(\langle \theta', x \rangle)| \geq \frac{g_{\min}}{\sqrt{2}} |\langle \theta - \theta', x \rangle|.$$

In light of these inequalities, bounds on all relevant regularity parameters for generalized linear bandits follow from similar reasoning to the proofs for linear bandits. In particular, the conclusion of [Lemma G.4](#) holds for generalized linear bandits under [Assumption G.1](#), with $L_{\mathcal{M}}^*$ as in [Lemma G.4](#), but scaled by $(\frac{g_{\max}}{g_{\min}})^2$.

□

G.3. Contextual Bandits with Finitely Many Actions (Example A.6)

In this setting we take $D(M(\pi) \parallel \bar{M}(\pi)) \leftarrow D_{\text{KL}}(M(\pi) \parallel \bar{M}(\pi))$ for D the divergence employed by AE_*^2 . Note that we have

$$D_{\text{KL}}(M(\pi) \parallel \bar{M}(\pi)) = \frac{1}{2} \mathbb{E}_{x \sim p_{\mathcal{X}}} [\mathbb{E}_{a \sim \pi(x)} [(f^M(x, a) - f^{\bar{M}}(x, a))^2]].$$

Lemma G.5. *In the contextual bandits setting of Example A.6:*

1. *Assumptions A.2, D.1 and D.2 hold with parameters*

$$L_{\text{KL}} = V_{\mathcal{M}} = 2\sqrt{2}$$

and $D(\cdot \parallel \cdot) = D_{\text{KL}}(\cdot \parallel \cdot)$.

2. *We can bound $N_{\text{cov}}(\mathcal{M}, \rho, \mu)$ by the covering number of \mathcal{M} in the distance $d(M, M') = \sup_{x \in \mathcal{X}, a \in \mathcal{A}} |f^M(x, a) - f^{M'}(x, a)|$ at tolerance $\frac{\sigma^2 \cdot \rho}{2 + \sqrt{2 \log(2/\mu)}}$. Furthermore, it suffices to take $\mathcal{E} := \{|r| \leq 1 + \sqrt{2 \log(2/\mu)}\}$.*

Proof of Lemma G.5. Using the expression for the KL divergence given above, for any $M, M', \bar{M} \in \mathcal{M}$ and $\pi \in \Pi$, we have

$$\begin{aligned} & |D_{\text{KL}}(M(\pi) \parallel M'(\pi)) - D_{\text{KL}}(\bar{M}(\pi) \parallel M'(\pi))| \\ &= \frac{1}{2} \left| \mathbb{E}_{x \sim p_{\mathcal{X}}} [\mathbb{E}_{a \sim \pi(x)} [(f^M(x, a) - f^{M'}(x, a))^2]] - \mathbb{E}_{x \sim p_{\mathcal{X}}} [\mathbb{E}_{a \sim \pi(x)} [(f^{\bar{M}}(x, a) - f^{M'}(x, a))^2]] \right| \\ &\leq \frac{1}{2} \mathbb{E}_{x \sim p_{\mathcal{X}}} \left[\mathbb{E}_{a \sim \pi(x)} \left[\left| (f^M(x, a) - f^{M'}(x, a))^2 - (f^{\bar{M}}(x, a) - f^{M'}(x, a))^2 \right| \right] \right] \\ &\leq 2 \mathbb{E}_{x \sim p_{\mathcal{X}}} \left[\mathbb{E}_{a \sim \pi(x)} \left[\left| f^M(x, a) - f^{\bar{M}}(x, a) \right| \right] \right] \\ &\leq 2 \sqrt{\mathbb{E}_{x \sim p_{\mathcal{X}}} \left[\mathbb{E}_{a \sim \pi(x)} \left[(f^M(x, a) - f^{\bar{M}}(x, a))^2 \right] \right]} \\ &= 2\sqrt{2} \sqrt{D_{\text{KL}}(M(\pi) \parallel \bar{M}(\pi))}. \end{aligned}$$

This verifies Assumption D.1 holds with $L_{\text{KL}} = 2\sqrt{2}$. Assumption D.2 is immediate.

To show that Assumption A.2 is met, we note that

$$\log \frac{\mathbb{P}^{\bar{M}, \pi}(r, o)}{\mathbb{P}^{M, \pi}(r, o)} = \log \frac{\mathbb{P}^{\bar{M}, \pi}(r \mid o) \mathbb{P}^{\bar{M}, \pi}(o)}{\mathbb{P}^{M, \pi}(r \mid o) \mathbb{P}^{M, \pi}(o)} = \log \frac{\mathbb{P}^{\bar{M}, \pi}(r \mid o)}{\mathbb{P}^{M, \pi}(r \mid o)},$$

where the second equality holds because the context distribution is identical for all models. As the reward likelihoods conditioned on the context are Gaussian, a calculation similar to Lemma G.1 shows that Assumption A.2 with $V_{\mathcal{M}} = 2\sqrt{2}$. The covering number bound also follows from the same reasoning as Lemma G.1. \square

Lemma G.6. *For the contextual bandit setting described above, we can bound*

$$C_{\text{exp}}(\mathcal{M}, \varepsilon) \leq \frac{4A}{\varepsilon}.$$

Proof of Lemma G.6. Fix some $\xi \in \Delta_{\mathcal{M}}$ and let π_{exp} be uniform over \mathcal{A} for each context. Then, for any $p' \in \Delta_{\Pi}$, we have

$$\begin{aligned}
 & \mathbb{E}_{\pi \sim p'} [D_{\text{KL}}(M(\pi) \| M'(\pi))] \\
 &= \frac{1}{2} \mathbb{E}_{\pi \sim p'} [\mathbb{E}_{x \sim p_{\mathcal{X}}} [\mathbb{E}_{a \sim \pi(x)} [(f^M(x, a) - f^{M'}(x, a))^2]]] \\
 &\leq \mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{\pi \sim p'} [\mathbb{E}_{x \sim p_{\mathcal{X}}} [\mathbb{E}_{a \sim \pi(x)} [(f^M(x, a) - f^{\bar{M}}(x, a))^2 + (f^{\bar{M}}(x, a) - f^{M'}(x, a))^2]]]] \\
 &\leq \sum_{a \in \mathcal{A}} \mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{x \sim p_{\mathcal{X}}} [(f^M(x, a) - f^{\bar{M}}(x, a))^2 + (f^{\bar{M}}(x, a) - f^{M'}(x, a))^2]] \\
 &= A \mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{x \sim p_{\mathcal{X}}} [\mathbb{E}_{a \sim p_{\text{exp}}(x)} [(f^M(x, a) - f^{\bar{M}}(x, a))^2 + (f^{\bar{M}}(x, a) - f^{M'}(x, a))^2]]] \\
 &\leq 2A \mathbb{E}_{\bar{M} \sim \xi} [D_{\text{KL}}(\bar{M}(\pi_{\text{exp}}) \| M(\pi_{\text{exp}})) + D_{\text{KL}}(\bar{M}(\pi_{\text{exp}}) \| M'(\pi_{\text{exp}}))].
 \end{aligned}$$

It follows that if

$$\mathbb{E}_{\bar{M} \sim \xi} [D_{\text{KL}}(\bar{M}(\pi_{\text{exp}}) \| M(\pi_{\text{exp}}))] \leq \frac{\varepsilon}{4A} \quad \text{and} \quad \mathbb{E}_{\bar{M} \sim \xi} [D_{\text{KL}}(\bar{M}(\pi_{\text{exp}}) \| M'(\pi_{\text{exp}}))] \leq \frac{\varepsilon}{4A},$$

then we can bound $\mathbb{E}_{\pi \sim p'} [D_{\text{KL}}(M(\pi) \| M'(\pi))] \leq \varepsilon$. Thus, choosing $p_{\text{exp}} \in \Delta(\Pi)$ to place probability mass 1 on π_{exp} , a sufficient bound on $C_{\text{exp}}(\mathcal{M}, \varepsilon)$ is $4A/\varepsilon$. \square

Proof of Proposition A.7. The bound on $C_{\text{exp}}(\mathcal{M}^*, \varepsilon)$ follows from Lemma G.6. Hence, by Lemma E.6 we can bound $\overline{\text{aec}}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}^*)$ as:

$$\overline{\text{aec}}_{\varepsilon}^{\mathcal{M}}(\mathcal{M}^*) \leq \frac{4A}{\delta} \quad \text{for} \quad \delta = \min_{M \in \mathcal{M}^*} \min \left\{ \min \left\{ \frac{1}{81L_{\text{KL}}}, \frac{\Delta_{\min}^M}{34V_{\mathcal{M}}} \right\} \cdot \frac{\varepsilon}{2\mathbf{g}^M/\Delta_{\min}^M + \mathbf{n}_{\varepsilon/36}^M}, \frac{\Delta_{\min}^M}{3} \right\}^2.$$

By Lemma E.2, we have that for all $M \in \mathcal{M}^*$,

$$\mathbf{g}^M \leq C_{\text{exp}}(\mathcal{M}^*, \frac{1}{4}(\Delta_{\min}^M)^2) \leq \frac{16A}{\Delta_{\min}^M}.$$

By Lemma G.5, we can take $L_{\text{KL}} = V_{\mathcal{M}} = 2\sqrt{2}$. A sufficient choice for δ is therefore

$$\delta = c \cdot \frac{\varepsilon^2 \Delta_{\min}^8}{A^2}.$$

The result follows. \square

G.4. Informative Arms (Example A.1)

In this section, we provide calculations for the bandits with informative arms setting in Example A.1. We first show that Assumptions A.1 to A.3 are satisfied.

- **Lemma G.1** If $\pi \in [A]$, then the response is simply Gaussian, so by Lemma G.1, the condition of Assumption A.1 is met with $L_{\text{KL}} = 2$. If $\pi = \pi_i^{\circ}$, then by the Mean Value Theorem we have

$$|D_{\text{KL}}(M(\pi) \| M'(\pi)) - D_{\text{KL}}(\bar{M}(\pi) \| M'(\pi))|$$

$$\begin{aligned}
 &= \left| \sum_{a \in [A]} \mathbb{P}^{M, \pi}(a) \log \frac{\mathbb{P}^{M, \pi}(a)}{\mathbb{P}^{M', \pi}(a)} - \sum_{a \in [A]} \mathbb{P}^{\bar{M}, \pi}(a) \log \frac{\mathbb{P}^{\bar{M}, \pi}(a)}{\mathbb{P}^{M', \pi}(a)} \right| \\
 &\leq \left(1 + \max_{a \in [A]} \max \left\{ \left| \log \frac{\mathbb{P}^{M, \pi}(a)}{\mathbb{P}^{M', \pi}(a)} \right|, \left| \log \frac{\mathbb{P}^{\bar{M}, \pi}(a)}{\mathbb{P}^{M', \pi}(a)} \right| \right\} \right) \cdot \sum_{a \in [A]} |\mathbb{P}^{M, \pi}(a) - \mathbb{P}^{\bar{M}, \pi}(a)| \\
 &= \left(1 + \max_{a \in [A]} \max \left\{ \left| \log \frac{\mathbb{P}^{M, \pi}(a)}{\mathbb{P}^{M', \pi}(a)} \right|, \left| \log \frac{\mathbb{P}^{\bar{M}, \pi}(a)}{\mathbb{P}^{M', \pi}(a)} \right| \right\} \right) \cdot D_{\text{TV}}(\mathbb{P}^{M, \pi}, \mathbb{P}^{\bar{M}, \pi}) \\
 &\leq \left(1 + \max_{a \in [A]} \max \left\{ \left| \log \frac{\mathbb{P}^{M, \pi}(a)}{\mathbb{P}^{M', \pi}(a)} \right|, \left| \log \frac{\mathbb{P}^{\bar{M}, \pi}(a)}{\mathbb{P}^{M', \pi}(a)} \right| \right\} \right) \cdot \sqrt{\frac{1}{2} D_{\text{KL}}(\mathbb{P}^{M, \pi} \parallel \mathbb{P}^{\bar{M}, \pi})}.
 \end{aligned}$$

Using the bound on the log-likelihood ratio given above, this verifies that [Assumption A.1](#) holds with $L_{\text{KL}} = \max\{2, 1 + \log \frac{A}{1-\beta}\}$.

- If $\pi \in [A]$, then since the response is Gaussian, by [Lemma G.1](#), the condition of [Assumption A.2](#) is met with $V_{\mathcal{M}} = 2$. If $\pi = \pi_i^\circ$, then for $M \in \mathcal{M}$, either the observation is distributed as $1/A$, so $\mathbb{P}^{M, \pi}(r, o) = 1/A$ for all $o \in [A]$, or i is the informative arm for instance M , in which case $\mathbb{P}^{M, \pi}(r, o) = (1 - \beta)/A$ for $o \neq \pi_M$, and $\mathbb{P}^{M, \pi}(r, o) = \beta + (1 - \beta)/A$ for $o = \pi_M$ (note that we can disregard $o = \perp$ since it occurs with probability 0 if an informative arm is pulled). The log-likelihood ratio is then at most

$$\log \frac{\beta + (1 - \beta)/A}{(1 - \beta)/A} \leq \log \frac{A}{1 - \beta}.$$

Thus, [Assumption A.2](#) is satisfied with $V_{\mathcal{M}} = \max\{2, \log \frac{A}{1-\beta}\}$.

- Using [Lemma G.1](#), it is easy to see that [Assumption A.3](#) is met with $d_{\text{cov}} = O(A)$ and $C_{\text{cov}} = O(1)$.

To bound the parameter $n^{\mathcal{M}}$, consider $M \in \mathcal{M}$ and $M' \in \mathcal{M}^{\text{alt}}(M)$. Note that either $D_{\text{KL}}(M(\pi_M) \parallel M'(\pi_M)) = 0$, in which case there is no advantage to playing π_M , or, due to the discretization of the means, $D_{\text{KL}}(M(\pi_M) \parallel M'(\pi_M)) \geq \frac{1}{2} \Delta_{\min}^2$. Thus, for any allocation $\eta \in \mathbb{R}_+^{\Pi}$, as long as $\eta(\pi_M) \geq \frac{2}{\Delta_{\min}^2}$, we have $\eta(\pi_M) D_{\text{KL}}(M(\pi_M) \parallel M'(\pi_M)) \geq 1$. It follows that there is no advantage to choosing $\eta(\pi_M)$ larger than $\frac{2}{\Delta_{\min}^2}$, so we can bound $n^{\mathcal{M}} \leq \frac{2}{\Delta_{\min}^2}$.

G.4.1. BOUNDING THE ALLOCATION-ESTIMATION COEFFICIENT

We begin with some basic observations. First, since we restrict \mathcal{M} to only contain instances with a single optimal decision, if $f^M(\pi_M) = \lfloor \frac{1}{\Delta_{\min}} \rfloor \Delta_{\min}$ for some $M \in \mathcal{M}$, this implies that $f^M(\pi) < \lfloor \frac{1}{\Delta_{\min}} \rfloor \Delta_{\min}$ for all $\pi \neq \pi_M$. Fix some $M \in \mathcal{M}$ satisfying $f^M(\pi_M) = \lfloor \frac{1}{\Delta_{\min}} \rfloor \Delta_{\min}$. It follows that, for every $M' \in \mathcal{M}^{\text{alt}}(M)$, it must be the case that $f^{M'}(\pi_M) < \lfloor \frac{1}{\Delta_{\min}} \rfloor \Delta_{\min}$. Therefore, since M and M' have different reward means at π_M , and since this holds for all $M' \in \mathcal{M}^{\text{alt}}(M)$, M can be distinguished from every $M' \in \mathcal{M}^{\text{alt}}(M)$ by playing π_M . In this case, then, $\mathbf{g}^M = 0$, so any ε -optimal Graves-Lai allocation must put all its mass on π_M , implying $\Lambda(M, \varepsilon) = \{\mathbb{I}_{\pi_M}\}$. Denote such instances M with $f^M(\pi_M) = \lfloor \frac{1}{\Delta_{\min}} \rfloor \Delta_{\min}$ as $\bar{\mathcal{M}}$. Note that for M with $f^M(\pi_M) < \lfloor \frac{1}{\Delta_{\min}} \rfloor \Delta_{\min}$, we have $\mathbb{I}_{\pi_M^\circ} \in \Lambda(M, \varepsilon)$.

We proceed to bound the value of the $\text{Fix } \bar{M} \in \text{co}(\mathcal{M})$. For a first case, assume that $f^{\bar{M}}(\pi_{\bar{M}}) \leq \lfloor \frac{1}{\Delta_{\min}} \rfloor \Delta_{\min} - \frac{1}{2} \Delta_{\min}$ and let $k = \arg \max_{i \in [N]} \mathbb{P}^{\bar{M}, \pi_i^\circ}(o = \pi_{\bar{M}})$ denote the index of the most informative arm for \bar{M} . Let $\lambda = \mathbb{I}_{\pi_k^\circ}$, and note that $\mathcal{M}_\varepsilon^{\text{gl}}(\lambda)$ contains only instances in \mathcal{M} that have informative arm k . Let $\mathcal{M}' = \{M \in \mathcal{M} : \pi_M^\circ \neq \pi_k^\circ\} \cup \bar{M}$. Then $\mathcal{M} \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda) \subseteq \mathcal{M}'$. Let $\omega = \frac{1}{2} \text{Unif}(\{\pi_i^\circ\}_{i \in [N]}) + \frac{1}{2} \text{Unif}([A])$. Then,

$$\begin{aligned} \text{aec}_\varepsilon(\mathcal{M}, \bar{M}) &\leq \sup_{M \in \mathcal{M}'} \frac{1}{\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \| M(\pi))]} \\ &\leq \sup_{M \in \mathcal{M}, \pi_M^\circ \neq \pi_k^\circ} \frac{2N}{\sum_{i=1}^N D_{\text{KL}}(\bar{M}(\pi_i^\circ) \| M(\pi_i^\circ))} + \sup_{M \in \bar{\mathcal{M}}} \frac{2A}{\sum_{\pi \in [A]} D_{\text{KL}}(\bar{M}(\pi) \| M(\pi))} \end{aligned}$$

If $\pi_M^\circ \neq \pi_k^\circ$, this implies that $o \sim M(\pi_i^\circ)$ is uniform on $[A]$ for $\pi_i^\circ \neq \pi_M^\circ$, and $o \sim M(\pi_i^\circ)$ is distributed as $\beta \mathbb{I}_{\pi_M} + (1 - \beta) \text{Unif}([A])$ for $\pi_i^\circ = \pi_M^\circ$. Note that since $k = \arg \max_{i \in [N]} \mathbb{P}^{\bar{M}, \pi_i^\circ}(o = \pi_{\bar{M}})$ and $\bar{M} \in \text{co}(\mathcal{M})$, we can have at most $\mathbb{P}^{\bar{M}, \pi_i^\circ}(o) \leq 1/A + \beta/2$ for all o if $i \neq k$, since if this were not the case, then i must be k . It follows from Pinsker's inequality that for M with $\pi_M^\circ \neq \pi_k^\circ$:

$$\begin{aligned} D_{\text{KL}}(\bar{M}(\pi_M^\circ) \| M(\pi_M^\circ)) &\geq 2D_{\text{TV}}(M(\pi_M^\circ), \bar{M}(\pi_M^\circ))^2 \\ &\geq 2|\mathbb{P}^{M, \pi_M^\circ}(o = \pi_M) - \mathbb{P}^{\bar{M}, \pi_M^\circ}(o = \pi_M)|^2 \\ &= 2|\beta - 1/A - \beta/2|^2 \\ &\geq 2|\beta/4|^2 \end{aligned}$$

where the last inequality uses our assumption that $\beta \geq 4/A$. For $M \in \bar{\mathcal{M}}$, since $f^{\bar{M}}(\pi_{\bar{M}}) \leq \lfloor \frac{1}{\Delta_{\min}} \rfloor \Delta_{\min} - \frac{1}{2} \Delta_{\min}$, we have that

$$D_{\text{KL}}(\bar{M}(\pi) \| M(\pi)) \geq \frac{1}{8} \Delta_{\min}^2.$$

Thus, we can bound

$$\text{aec}_\varepsilon(\mathcal{M}, \bar{M}) \leq \frac{64N}{\beta^2} + \frac{16A}{\Delta_{\min}^2}.$$

Now, consider the second case where \bar{M} has $f^{\bar{M}}(\pi_{\bar{M}}) > \lfloor \frac{1}{\Delta_{\min}} \rfloor \Delta_{\min} - \frac{1}{2} \Delta_{\min}$. Note that in this case we must have $|\pi_{\bar{M}}| = 1$. Set $\lambda = \mathbb{I}_{\pi_{\bar{M}}}$. Then we have that $\mathcal{M}_\varepsilon^{\text{gl}}(\lambda)$ contains every instance except the single instance with $f^M(\pi_{\bar{M}}) = \lfloor \frac{1}{\Delta_{\min}} \rfloor \Delta_{\min}$. Let $\omega = \mathbb{I}_{\pi_{\bar{M}}}$. Note that for any instance with $f^M(\pi_{\bar{M}}) < \lfloor \frac{1}{\Delta_{\min}} \rfloor \Delta_{\min}$, i.e. every instance in $\mathcal{M} \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)$, we have

$$D_{\text{KL}}(\bar{M}(\pi_{\bar{M}}) \| M(\pi_{\bar{M}})) \geq \frac{1}{8} \Delta_{\min}^2.$$

It follows that with such an \bar{M} , we can bound

$$\text{aec}_\varepsilon(\mathcal{M}, \bar{M}) \leq \frac{8}{\Delta_{\min}^2}.$$

G.5. Tabular Reinforcement Learning (Appendix A.7)

In this section, we prove all of the claims in [Appendix A.7](#) concerning tabular reinforcement learning.

Throughout this section we let $M_{sh}(a)$ denote the joint distribution of the next state and reward if we play action a in state s at step h on MDP $M \in \mathcal{M}$. We also define

$$w_h^{M,\pi}(s, a) = \mathbb{P}^{M,\pi}(s_h = s, a_h = a)$$

as the state-action visitation probabilities on MDP M under policy π (and define $w_h^{M,\pi}(s)$ analogously). We let $r_h^M(s, a) = \mathbb{E}_{r \sim R_h^M(s,a)}[r]$ denote the mean reward on MDP M at (s, a, h) , and let $r_h(s_h, a_h)$ denote the realized (random reward) at step h . We let $r := (r_1(s_1, a_1), \dots, r_H(s_H, a_H))$ denote the vector of all random rewards in a given episode. $\tau = (s_1, \dots, s_H)$ denotes a trajectory of states, and $\tau_h = s_h$ the h th state in the trajectory. We denote the Q -value function on M for policy π by

$$Q_h^{M,\pi}(s, a) = \mathbb{E}^{M,\pi} \left[\sum_{h'=h}^H r_{h'}(s_{h'}, a_{h'}) \mid s_h = s, a_h = a \right]$$

and the value function by $V_h^{M,\pi}(s) = \mathbb{E}_{a \sim \pi_h(s)}[Q_h^{M,\pi}(s, a)]$. We denote the value of a policy by $V_1^{M,\pi} := V_1^{M,\pi}(s_1)$. For any function $V : \mathcal{S} \rightarrow \mathbb{R}$ we denote

$$\mathbb{P}_h^M[V](s, a) = \mathbb{E}_{s' \sim P_h^M(\cdot | s, a)}[V(s')].$$

For all results in this section concerning general divergences, we take $D(\cdot \| \cdot) \leftarrow D_{\text{KL}}(\cdot \| \cdot)$.

Proof of [Proposition A.8](#). To bound the AEC, we first move from KL divergence to Hellinger distance. Since we always have $D_{\text{KL}}(\bar{M}(\pi) \| M(\pi)) \geq D_{\text{H}}^2(\bar{M}(\pi), M(\pi))$, we upper bound

$$\overline{\text{aEC}}_\varepsilon^{\mathcal{M}}(\mathcal{M}^*) \leq \sup_{\xi \in \Delta_{\mathcal{M}}} \inf_{\lambda, \omega \in \Delta_{\Pi}} \sup_{M \in \mathcal{M}^* \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)} \frac{1}{\mathbb{E}_{\bar{M} \sim \xi}[\mathbb{E}_\omega[D_{\text{H}}^2(\bar{M}(\pi), M(\pi))]]}.$$

We then apply [Lemma E.6](#) to bound this by $C_{\text{exp}}^{\text{D}}(\mathcal{M}^*, \delta)$, with $D(\cdot \| \cdot) \leftarrow D_{\text{H}}^2(\cdot, \cdot)$. The bound on $C_{\text{exp}}^{\text{D}}(\mathcal{M}^*, \varepsilon)$ then follows directly from [Lemma G.10](#), and gives

$$\overline{\text{aEC}}_\varepsilon^{\mathcal{M}}(\mathcal{M}^*) \leq \frac{SAH^2 \cdot \log^2 H}{\delta^2} \quad \text{for} \quad \delta = \min_{M \in \mathcal{M}^*} \min \left\{ \min \left\{ \frac{1}{81L_{\text{KL}}}, \frac{\Delta_{\min}^M}{34V_{\mathcal{M}}} \right\} \cdot \frac{\varepsilon}{2g^M/\Delta_{\min}^M + n_{\varepsilon/36}^M}, \frac{\Delta_{\min}^M}{3} \right\}^2.$$

By [Lemma G.8](#), we have that [Assumptions A.2](#) and [D.1](#) hold with

$$L_{\text{KL}} = V_{\mathcal{M}} = 4H + \max_{\bar{M}, \bar{M}' \in \mathcal{M}} \max_{\pi \in \Pi} \max_{\tau \in \mathcal{T}} \left| \log \frac{\mathbb{P}^{\bar{M}', \pi}(\tau)}{\mathbb{P}^{\bar{M}, \pi}(\tau)} \right|.$$

As we assume every transition has probability at least P_{\min} , we can lower bound $\mathbb{P}^{\bar{M}, \pi}(\tau) \geq P_{\min}^H$, so it suffices to take

$$L_{\text{KL}} = V_{\mathcal{M}} = H(4 + \log 1/P_{\min}).$$

By [Lemma E.2](#), for $M \in \mathcal{M}^*$ we can bound

$$\mathbf{g}^M \leq C_{\text{exp}}(\mathcal{M}^*, \frac{1}{4}(\Delta_{\min}^M)^2) \leq c \cdot \frac{SAH^2 \cdot \log^2 H}{\Delta_{\star}^4}.$$

A sufficient choice of δ is therefore

$$\delta = c \cdot \frac{\varepsilon^2 \Delta_{\star}^{12}}{S^2 A^2 H^6 \cdot \log^4 H \cdot \log^2 1/P_{\min}}.$$

□

G.5.1. TABULAR MDPs ARE REGULAR CLASSES

Lemma G.7. *If M^* is such that $r_h^{M^*}(s, a) \in [0, 1/H^2]$ for all (s, a, h) , then in the setting of $\mathcal{M} \leftarrow \mathcal{M}_{\text{tab}}(P_{\min})$, M^* is a regular model with constant*

$$L_{\mathcal{M}}^{\star} = \frac{96}{\Delta_{\min}^{\star}}$$

Proof of Lemma G.7. Take some MDP $M' \in \mathcal{M}^{\text{alt}}(M^*)$ such that $D_{\text{KL}}(M^*(\pi_{\star}) \| M'(\pi_{\star})) > 0$. Let M'' be such that

$$D_{\text{KL}}(M_{sh}^*(\pi_{\star}(s, h)) \| M_{sh}''(\pi_{\star}(s, h))) = 0, \quad \forall s, h$$

and

$$D_{\text{KL}}(M_{sh}'(a) \| M_{sh}''(a)) = 0, \quad \forall s, h, a \neq \pi_{\star}(s, h),$$

so that M'' is the MDP which is identical to M^* on optimal actions, and identical to M' on suboptimal actions (recall that optimal actions for M^* are unique). By construction, we have that $D_{\text{KL}}(M^*(\pi_{\star}) \| M''(\pi_{\star})) = 0$. Furthermore, it is not difficult to see that $M'' \in \mathcal{M}$. In particular, to verify that $P_h^{M''}(s' | s, a) \geq P_{\min}$ for each (s, a, h, s') , we note that since $M^*, M' \in \mathcal{M}$, for every (s, a, h, s') , we have $P_h^{M^*}(s' | s, a) \geq P_{\min}$ and $P_h^{M'}(s' | s, a) \geq P_{\min}$. By construction, we have that $P_h^{M''}(\cdot | s, a)$ is identical to either $P_h^{M^*}(\cdot | s, a)$ or $P_h^{M'}(\cdot | s, a)$ for each (s, a, h) , so it follows that $P_h^{M''}(s' | s, a) \geq P_{\min}$. The remaining conditions for inclusion in \mathcal{M} are immediate. We consider two cases.

Case 1: $M'' \in \mathcal{M}^{\text{alt}}(M^*)$. For $\pi \in \Pi$, by [Lemma G.15](#) we have

$$D_{\text{KL}}(M^*(\pi) \| M'(\pi)) = \sum_{s, a, h} w_h^{M^*, \pi}(s, a) D_{\text{KL}}(M_{sh}^*(a) \| M_{sh}'(a)),$$

and

$$\begin{aligned} D_{\text{KL}}(M^*(\pi) \| M''(\pi)) &= \sum_{s, a, h} w_h^{M^*, \pi}(s, a) D_{\text{KL}}(M_{sh}^*(a) \| M_{sh}''(a)) \\ &= \sum_{s, a, h} w_h^{M^*, \pi}(s, a) D_{\text{KL}}(M_{sh}^*(a) \| M_{sh}'(a)) \cdot \mathbb{I}\{a \neq \pi_{\star}(s, h)\} \end{aligned}$$

so that

$$\begin{aligned}
 & |D_{\text{KL}}(M^*(\pi) \parallel M'(\pi)) - D_{\text{KL}}(M^*(\pi) \parallel M''(\pi))| \\
 &= \sum_{s,h} w_h^{M^*,\pi}(s, \pi_*(s, h)) D_{\text{KL}}(M_{sh}^*(\pi_*(s, h)) \parallel M'_{sh}(\pi_*(s, h))) \\
 &\leq \sup_{s,h} \frac{w_h^{M^*,\pi}(s, \pi_*(s, h))}{w_h^{M^*,\pi_*}(s, \pi_*(s, h))} \cdot D_{\text{KL}}(M^*(\pi_*) \parallel M'(\pi_*)) \\
 &\leq \sup_{s,h} \frac{1}{w_h^{M^*,\pi_*}(s)} \cdot D_{\text{KL}}(M^*(\pi_*) \parallel M'(\pi_*)) \\
 &\leq \frac{1}{\Delta_{\min}^*} \cdot D_{\text{KL}}(M^*(\pi_*) \parallel M'(\pi_*))
 \end{aligned}$$

where the last inequality follows from [Lemma G.13](#). Thus, in this case, M^* is a regular model with

$$L_{\mathcal{M}}^* = \frac{1}{\Delta_{\min}^*}.$$

Case 2: $M'' \notin \mathcal{M}^{\text{alt}}(M^*)$. Let $(\tilde{s}, \tilde{a}, \tilde{h})$ be such that $Q_{\tilde{h}}^{M',\pi_*}(\tilde{s}, \tilde{a}) > Q_{\tilde{h}}^{M',\pi_*}(\tilde{s}, \pi_*(\tilde{s}, \tilde{h}))$, and note that such a tuple is guaranteed to exist by [Lemma G.11](#) since $M' \in \mathcal{M}^{\text{alt}}(M^*)$. Let \tilde{M}'' denote an MDP that is identical to M'' everywhere except for at $(\tilde{s}, \tilde{h}, \tilde{a})$, where we set $r_{\tilde{h}}^{\tilde{M}''}(\tilde{s}, \tilde{a})$ so that

$$Q_{\tilde{h}}^{\tilde{M}'',\pi_*}(\tilde{s}, \tilde{a}) = Q_{\tilde{h}}^{\tilde{M}'',\pi_*}(\tilde{s}, \pi_*(\tilde{s}, \tilde{h})) + \delta \quad (75)$$

for some $\delta > 0$ to be chosen. This will ensure \tilde{a} is the optimal action in (\tilde{s}, \tilde{h}) , so $\pi_{\tilde{M}''} \neq \pi_*$. By construction we have that M^* and \tilde{M}'' behave identically on π_* , which implies that $Q_{\tilde{h}}^{\tilde{M}'',\pi_*}(\tilde{s}, \pi_*(\tilde{s}, \tilde{h})) = Q_{\tilde{h}}^{M^*,\pi_*}(\tilde{s}, \pi_*(\tilde{s}, \tilde{h}))$. Furthermore, by assumption we have $r_{\tilde{h}}^{M^*}(s, a) < 1/H^2$ for all (s, a, h) , which implies $Q_{\tilde{h}}^{M^*,\pi_*}(\tilde{s}, \pi_*(\tilde{s}, \tilde{h})) < 1/H$. As $Q_{\tilde{h}}^{\tilde{M}'',\pi_*}(\tilde{s}, \tilde{a}) = r_{\tilde{h}}^{\tilde{M}''}(\tilde{s}, \tilde{a}) + \mathbb{P}_{\tilde{h}}^{\tilde{M}''}[V_{\tilde{h}+1}^{\tilde{M}'',\pi_*}](\tilde{s}, \tilde{a}) \geq r_{\tilde{h}}^{\tilde{M}''}(\tilde{s}, \tilde{a})$, it follows that for small enough δ , we can ensure [Eq. \(75\)](#) is met with $r_{\tilde{h}}^{\tilde{M}''}(\tilde{s}, \tilde{a}) < 1/H$, so that $\tilde{M}'' \in \mathcal{M}$.

If we can show that, for all π , $|D_{\text{KL}}(M^*(\pi) \parallel M'(\pi)) - D_{\text{KL}}(M^*(\pi) \parallel \tilde{M}''(\pi))|$ is bounded by some function of $D_{\text{KL}}(M^*(\pi_*) \parallel M'(\pi_*))$, we are then done. We proceed to show this. First, note that, similar to Case 1:

$$\begin{aligned}
 & |D_{\text{KL}}(M^*(\pi) \parallel M'(\pi)) - D_{\text{KL}}(M^*(\pi) \parallel \tilde{M}''(\pi))| \\
 &= \left| \sum_{s,h} w_h^{M^*,\pi}(s, \pi_*(s, h)) D_{\text{KL}}(M_{sh}^*(\pi_*(s, h)) \parallel M'_{sh}(\pi_*(s, h))) \right| \\
 &\quad + w_{\tilde{h}}^{M^*,\pi}(\tilde{s}, \tilde{a}) \left| D_{\text{KL}}(M_{\tilde{s}\tilde{h}}^*(\tilde{a}) \parallel M'_{\tilde{s}\tilde{h}}(\tilde{a})) - D_{\text{KL}}(M_{\tilde{s}\tilde{h}}^*(\tilde{a}) \parallel \tilde{M}''_{\tilde{s}\tilde{h}}(\tilde{a})) \right| \\
 &\leq \sup_{s,h} \frac{1}{w_h^{M^*,\pi_*}(s)} \cdot D_{\text{KL}}(M^*(\pi_*) \parallel M'(\pi_*)) \\
 &\quad + \frac{1}{2} w_{\tilde{h}}^{M^*,\pi}(\tilde{s}, \tilde{a}) \left| (r_{\tilde{h}}^{M^*}(\tilde{s}, \tilde{a}) - r_{\tilde{h}}^{M'}(\tilde{s}, \tilde{a}))^2 - (r_{\tilde{h}}^{M^*}(\tilde{s}, \tilde{a}) - r_{\tilde{h}}^{\tilde{M}''}(\tilde{s}, \tilde{a}))^2 \right|
 \end{aligned} \quad (76)$$

where the inequality follows by what we showed in Case 1, and since M' and \widetilde{M}'' have identical transitions at $(\widetilde{s}, \widetilde{a}, \widetilde{h})$, so the contribution to the KL divergence from the transitions cancels, leaving only the KL divergence between unit-variance Gaussians. By the Mean Value Theorem and since rewards are in $[0, 1/H]$, we have

$$\frac{1}{2} \left| (r_{\widetilde{h}}^{M^*}(\widetilde{s}, \widetilde{a}) - r_{\widetilde{h}}^{M'}(\widetilde{s}, \widetilde{a}))^2 - (r_{\widetilde{h}}^{M^*}(\widetilde{s}, \widetilde{a}) - r_{\widetilde{h}}^{\widetilde{M}''}(\widetilde{s}, \widetilde{a}))^2 \right| \leq \frac{1}{H} |r_{\widetilde{h}}^{M'}(\widetilde{s}, \widetilde{a}) - r_{\widetilde{h}}^{\widetilde{M}''}(\widetilde{s}, \widetilde{a})|.$$

Thus, it suffices to bound $|r_{\widetilde{h}}^{M'}(\widetilde{s}, \widetilde{a}) - r_{\widetilde{h}}^{\widetilde{M}''}(\widetilde{s}, \widetilde{a})|$.

By assumption $Q_{\widetilde{h}}^{M', \pi_*}(\widetilde{s}, \widetilde{a}) > Q_{\widetilde{h}}^{M', \pi_*}(\widetilde{s}, \pi_*(\widetilde{s}, \widetilde{h}))$. We can then ensure

$$Q_{\widetilde{h}}^{\widetilde{M}'', \pi_*}(\widetilde{s}, \widetilde{a}) - Q_{\widetilde{h}}^{\widetilde{M}'', \pi_*}(\widetilde{s}, \pi_*(\widetilde{s}, \widetilde{h})) \leq Q_{\widetilde{h}}^{M', \pi_*}(\widetilde{s}, \widetilde{a}) - Q_{\widetilde{h}}^{M', \pi_*}(\widetilde{s}, \pi_*(\widetilde{s}, \widetilde{h}))$$

for δ sufficiently small. This is equivalent to, abbreviating $\widetilde{a}_{M^*} := \pi_*(\widetilde{s}, \widetilde{h})$:

$$\begin{aligned} & r_{\widetilde{h}}^{\widetilde{M}''}(\widetilde{s}, \widetilde{a}) + \mathbb{P}_{\widetilde{h}}^{\widetilde{M}''} [V_{\widetilde{h}+1}^{\widetilde{M}'', \pi_*}](\widetilde{s}, \widetilde{a}) - r_{\widetilde{h}}^{\widetilde{M}''}(\widetilde{s}, \widetilde{a}_{M^*}) - \mathbb{P}_{\widetilde{h}}^{\widetilde{M}''} [V_{\widetilde{h}+1}^{\widetilde{M}'', \pi_*}](\widetilde{s}, \widetilde{a}_{M^*}) \\ & \leq r_{\widetilde{h}}^{M'}(\widetilde{s}, \widetilde{a}) + \mathbb{P}_{\widetilde{h}}^{M'} [V_{\widetilde{h}+1}^{M', \pi_*}](\widetilde{s}, \widetilde{a}) - r_{\widetilde{h}}^{M'}(\widetilde{s}, \widetilde{a}_{M^*}) - \mathbb{P}_{\widetilde{h}}^{M'} [V_{\widetilde{h}+1}^{M', \pi_*}](\widetilde{s}, \widetilde{a}_{M^*}). \end{aligned}$$

By construction we have that $V_{\widetilde{h}}^{\widetilde{M}'', \pi_*}(s) = V_{\widetilde{h}}^{M^*, \pi_*}(s)$ for all (s, h) , $r_{\widetilde{h}}^{\widetilde{M}''}(\widetilde{s}, \widetilde{a}_{M^*}) = r_{\widetilde{h}}^{M^*}(\widetilde{s}, \widetilde{a}_{M^*})$, and $\mathbb{P}_{\widetilde{h}}^{\widetilde{M}''} [V_{\widetilde{h}+1}^{\widetilde{M}'', \pi_*}](\widetilde{s}, \widetilde{a}_{M^*}) = \mathbb{P}_{\widetilde{h}}^{M^*} [V_{\widetilde{h}+1}^{M^*, \pi_*}](\widetilde{s}, \widetilde{a}_{M^*})$, since \widetilde{M}'' behaves identically to M on actions taken by π_* . Furthermore, we have $\mathbb{P}_{\widetilde{h}}^{\widetilde{M}''} [V_{\widetilde{h}+1}^{\widetilde{M}'', \pi_*}](\widetilde{s}, \widetilde{a}) = \mathbb{P}_{\widetilde{h}}^{M'} [V_{\widetilde{h}+1}^{M', \pi_*}](\widetilde{s}, \widetilde{a})$ since \widetilde{M}'' behaves identically to M' on actions not taken by π_* , other than the reward at $(\widetilde{s}, \widetilde{a}, \widetilde{h})$. Using these simplifications and rearranging, we get

$$\begin{aligned} |r_{\widetilde{h}}^{\widetilde{M}''}(\widetilde{s}, \widetilde{a}) - r_{\widetilde{h}}^{M'}(\widetilde{s}, \widetilde{a})| & \leq |r_{\widetilde{h}}^{\widetilde{M}''}(\widetilde{s}, \widetilde{a}_{M^*}) - r_{\widetilde{h}}^{M^*}(\widetilde{s}, \widetilde{a}_{M^*})| + |\mathbb{P}_{\widetilde{h}}^{\widetilde{M}''} [V_{\widetilde{h}+1}^{\widetilde{M}'', \pi_*}](\widetilde{s}, \widetilde{a}_{M^*}) - \mathbb{P}_{\widetilde{h}}^{M'} [V_{\widetilde{h}+1}^{M', \pi_*}](\widetilde{s}, \widetilde{a}_{M^*})| \\ & \quad + |\mathbb{P}_{\widetilde{h}}^{M'} [V_{\widetilde{h}+1}^{M', \pi_*}](\widetilde{s}, \widetilde{a}) - \mathbb{P}_{\widetilde{h}}^{M'} [V_{\widetilde{h}+1}^{M', \pi_*}](\widetilde{s}, \widetilde{a})| \\ & \leq |r_{\widetilde{h}}^{M'}(\widetilde{s}, \widetilde{a}_{M^*}) - r_{\widetilde{h}}^{M^*}(\widetilde{s}, \widetilde{a}_{M^*})| + |\mathbb{P}_{\widetilde{h}}^{M^*} [V_{\widetilde{h}+1}^{M^*, \pi_*}](\widetilde{s}, \widetilde{a}_{M^*}) - \mathbb{P}_{\widetilde{h}}^{M'} [V_{\widetilde{h}+1}^{M', \pi_*}](\widetilde{s}, \widetilde{a}_{M^*})| \\ & \quad + |\mathbb{P}_{\widetilde{h}}^{M'} [V_{\widetilde{h}+1}^{M', \pi_*}](\widetilde{s}, \widetilde{a}_{M^*}) - \mathbb{P}_{\widetilde{h}}^{M'} [V_{\widetilde{h}+1}^{M', \pi_*}](\widetilde{s}, \widetilde{a}_{M^*})| \\ & \quad + |\mathbb{P}_{\widetilde{h}}^{M'} [V_{\widetilde{h}+1}^{M', \pi_*}](\widetilde{s}, \widetilde{a}) - \mathbb{P}_{\widetilde{h}}^{M'} [V_{\widetilde{h}+1}^{M', \pi_*}](\widetilde{s}, \widetilde{a})|. \end{aligned}$$

Since rewards are unit-variance Gaussian, we have

$$|r_{\widetilde{h}}^{M'}(\widetilde{s}, \widetilde{a}_{M^*}) - r_{\widetilde{h}}^{M^*}(\widetilde{s}, \widetilde{a}_{M^*})| \leq \sqrt{2D_{\text{KL}}(M_{\widetilde{h}, \widetilde{s}}^{M^*}(\widetilde{a}_{M^*}) \| M_{\widetilde{h}, \widetilde{s}}^{M'}(\widetilde{a}_{M^*}))} \leq \sqrt{\frac{2}{w_{\widetilde{h}}^{M^*, \pi_*}(\widetilde{s})} D_{\text{KL}}(M^*(\pi_*) \| M'(\pi_*))}.$$

Since $V_{\widetilde{h}+1}^{M^*, \pi_*} \in [0, 1]$, we have

$$\begin{aligned} |\mathbb{P}_{\widetilde{h}}^{M^*} [V_{\widetilde{h}+1}^{M^*, \pi_*}](\widetilde{s}, \widetilde{a}_{M^*}) - \mathbb{P}_{\widetilde{h}}^{M'} [V_{\widetilde{h}+1}^{M', \pi_*}](\widetilde{s}, \widetilde{a}_{M^*})| & \leq \sum_{s'} |P_{\widetilde{h}}^{M^*}(s' | \widetilde{s}, \widetilde{a}_{M^*}) - P_{\widetilde{h}}^{M'}(s' | \widetilde{s}, \widetilde{a}_{M^*})| \\ & \leq 2D_{\text{TV}}(P_{\widetilde{h}}^{M^*}(\cdot | \widetilde{s}, \widetilde{a}_{M^*}), P_{\widetilde{h}}^{M'}(\cdot | \widetilde{s}, \widetilde{a}_{M^*})) \\ & \leq \sqrt{2D_{\text{KL}}(P_{\widetilde{h}}^{M^*}(\cdot | \widetilde{s}, \widetilde{a}_{M^*}) \| P_{\widetilde{h}}^{M'}(\cdot | \widetilde{s}, \widetilde{a}_{M^*}))} \end{aligned}$$

$$\begin{aligned}
 &\leq \sqrt{\frac{2}{w_{\tilde{h}}^{M^*, \pi^*}(\tilde{s})} D_{\text{KL}}(P_{\tilde{h}}^M(\cdot | \tilde{s}, \tilde{a}_{M^*}) \| P_{\tilde{h}}^{M'}(\cdot | \tilde{s}, \tilde{a}_{M^*}))} \\
 &\leq \sqrt{\frac{2}{w_{\tilde{h}}^{M^*, \pi^*}(\tilde{s})} D_{\text{KL}}(M^*(\pi_*) \| M'(\pi_*))}.
 \end{aligned}$$

By [Lemma G.12](#) we have

$$\begin{aligned}
 |\mathbb{P}_{\tilde{h}}^{M'}[V_{\tilde{h}+1}^{M^*, \pi^*}](\tilde{s}, \tilde{a}_{M^*}) - \mathbb{P}_{\tilde{h}}^{M'}[V_{\tilde{h}+1}^{M', \pi^*}](\tilde{s}, \tilde{a}_{M^*})| &\leq \sum_{s'} P_{\tilde{h}}^{M'}(s' | \tilde{s}, \tilde{a}_{M^*}) |V_{\tilde{h}+1}^{M^*, \pi^*}(s') - V_{\tilde{h}+1}^{M', \pi^*}(s')| \\
 &\leq \sum_{s'} P_{\tilde{h}}^{M'}(s' | \tilde{s}, \tilde{a}_{M^*}) \cdot \sqrt{\frac{8H}{w_{\tilde{h}+1}^{M^*, \pi^*}(s')} \cdot D_{\text{KL}}(M^*(\pi_*) \| M'(\pi_*))} \\
 &\leq \sup_s \sqrt{\frac{8H}{w_{\tilde{h}+1}^{M^*, \pi^*}(s)} \cdot D_{\text{KL}}(M^*(\pi_*) \| M'(\pi_*))}
 \end{aligned}$$

and similarly

$$|\mathbb{P}_{\tilde{h}}^{M'}[V_{\tilde{h}+1}^{M^*, \pi^*}](\tilde{s}, \tilde{a}) - \mathbb{P}_{\tilde{h}}^{M'}[V_{\tilde{h}+1}^{M', \pi^*}](\tilde{s}, \tilde{a})| \leq \sup_s \sqrt{\frac{8H}{w_{\tilde{h}+1}^{M^*, \pi^*}(s)} \cdot D_{\text{KL}}(M^*(\pi_*) \| M'(\pi_*))}.$$

Altogether then:

$$\begin{aligned}
 |r_{\tilde{h}}^{\tilde{M}''}(\tilde{s}, \tilde{a}) - r_{\tilde{h}}^{M'}(\tilde{s}, \tilde{a})| &\leq \left(\sqrt{\frac{8}{w_{\tilde{h}}^{M^*, \pi^*}(\tilde{s})}} + \sup_s \sqrt{\frac{32H}{w_{\tilde{h}+1}^{M^*, \pi^*}(s)}} \right) \cdot \sqrt{D_{\text{KL}}(M^*(\pi_*) \| M'(\pi_*))} \\
 &\leq \sup_{s, h} \sqrt{\frac{96H}{w_h^{M^*, \pi^*}(s)}} \cdot \sqrt{D_{\text{KL}}(M^*(\pi_*) \| M'(\pi_*))}.
 \end{aligned}$$

Combining this with [Eq. \(76\)](#), we have

$$\begin{aligned}
 &|D_{\text{KL}}(M^*(\pi) \| M'(\pi)) - D_{\text{KL}}(M^*(\pi) \| \tilde{M}''(\pi))| \\
 &\leq \sup_{s, h} \frac{1}{w_h^{M^*, \pi^*}(s)} \cdot D_{\text{KL}}(M^*(\pi_*) \| M'(\pi_*)) + \sup_{s, h} \sqrt{\frac{96}{H w_h^{M^*, \pi^*}(s)}} \cdot \sqrt{D_{\text{KL}}(M^*(\pi_*) \| M'(\pi_*))} \\
 &\leq \frac{1}{\Delta_{\min}^*} \cdot D_{\text{KL}}(M^*(\pi_*) \| M'(\pi_*)) + \sqrt{\frac{96}{H \Delta_{\min}^*}} \cdot \sqrt{D_{\text{KL}}(M^*(\pi_*) \| M'(\pi_*))}
 \end{aligned}$$

where the second inequality uses [Lemma G.13](#). Thus, in this case M^* is a regular model with

$$L_{\mathcal{M}}^* = \frac{96}{\Delta_{\min}^*}.$$

□

G.5.2. TABULAR MDPs SATISFY BASIC ASSUMPTIONS

Lemma G.8. *Tabular MDPs with unit-variance Gaussian rewards satisfy [Assumptions A.2, D.1](#) and [D.2](#) with*

$$L_{\text{KL}} = V_{\mathcal{M}} = 8H + \max_{\bar{M}, \bar{M}' \in \mathcal{M}} \max_{\pi \in \Pi} \max_{\tau \in \mathcal{T}} \left| \log \frac{\mathbb{P}^{\bar{M}', \pi}(\tau)}{\mathbb{P}^{\bar{M}, \pi}(\tau)} \right|,$$

and $D(\cdot \parallel \cdot) = D_{\text{KL}}(\cdot \parallel \cdot)$, where $\mathcal{T} := \mathcal{S}^H$ and $\mathbb{P}^{\bar{M}, \pi}(\tau)$ denotes the probability of observing state sequence $\tau \in \mathcal{T}$ on \bar{M} when playing policy π .

Proof of Lemma G.8. We verify each assumption separately.

Verifying Assumption D.1. Fix some $M, M', \bar{M} \in \mathcal{M}$. Our goal is to bound

$$\left| D_{\text{KL}}(M(\pi) \parallel M'(\pi)) - D_{\text{KL}}(\bar{M}(\pi) \parallel M'(\pi)) \right|.$$

Let \widetilde{M} denote the MDP with transitions identical to M but rewards identical to \bar{M} . Then

$$\begin{aligned} \left| D_{\text{KL}}(M(\pi) \parallel M'(\pi)) - D_{\text{KL}}(\bar{M}(\pi) \parallel M'(\pi)) \right| &\leq \left| D_{\text{KL}}(M(\pi) \parallel M'(\pi)) - D_{\text{KL}}(\widetilde{M}(\pi) \parallel M'(\pi)) \right| \\ &\quad + \left| D_{\text{KL}}(\widetilde{M}(\pi) \parallel M'(\pi)) - D_{\text{KL}}(\bar{M}(\pi) \parallel M'(\pi)) \right| \end{aligned}$$

We bound these terms separately. First, by [Lemma G.15](#) we have

$$\begin{aligned} &D_{\text{KL}}(M(\pi) \parallel M'(\pi)) \\ &= \sum_{s, a, h} w_h^{M, \pi}(s, a) D_{\text{KL}}(M_{sh}(\pi(s, h)) \parallel M'_{sh}(\pi(s, h))) \\ &= \sum_{s, a, h} w_h^{M, \pi}(s, a) \left[\frac{1}{2} (r_h^M(s, a) - r_h^{M'}(s, a))^2 + D_{\text{KL}}(P_h^M(\cdot \mid s, \pi(s, h)) \parallel P_h^{M'}(\cdot \mid s, \pi(s, h))) \right] \end{aligned}$$

and, given our definition of \widetilde{M} ,

$$\begin{aligned} &D_{\text{KL}}(\widetilde{M}(\pi) \parallel M'(\pi)) \\ &= \sum_{s, a, h} w_h^{M, \pi}(s, a) \left[\frac{1}{2} (r_h^{\bar{M}}(s, a) - r_h^{M'}(s, a))^2 + D_{\text{KL}}(P_h^M(\cdot \mid s, \pi(s, h)) \parallel P_h^{M'}(\cdot \mid s, \pi(s, h))) \right]. \end{aligned}$$

Thus,

$$\begin{aligned} &\left| D_{\text{KL}}(M(\pi) \parallel M'(\pi)) - D_{\text{KL}}(\widetilde{M}(\pi) \parallel M'(\pi)) \right| \\ &= \left| \frac{1}{2} \sum_{s, a, h} w_h^{M, \pi}(s, a) \left[(r_h^M(s, a) - r_h^{M'}(s, a))^2 - (r_h^{\bar{M}}(s, a) - r_h^{M'}(s, a))^2 \right] \right| \\ &\stackrel{(a)}{\leq} \sum_{s, a, h} w_h^{M, \pi}(s, a) |r_h^M(s, a) - r_h^{\bar{M}}(s, a)| \end{aligned}$$

$$\begin{aligned}
 &\leq \sum_{s,a,h} w_h^{M,\pi}(s,a) \sqrt{2D_{\text{KL}}(M_{sh}(a) \parallel \bar{M}_{sh}(a))} \\
 &\leq \sqrt{2H \cdot \sum_{s,a,h} w_h^{M,\pi}(s,a) D_{\text{KL}}(M_{sh}(a) \parallel \bar{M}_{sh}(a))} \\
 &= \sqrt{2H \cdot D_{\text{KL}}(M(\pi) \parallel \bar{M}(\pi))},
 \end{aligned}$$

where (a) holds by the Mean Value Theorem and the assumption that reward means are in $[0, 1]$, and the final equality holds by [Lemma G.15](#).

We turn now to bounding the second term. Let $\mathcal{T} = \mathcal{S}^H$ denote the space of all possible state trajectories. Let $\mathbb{P}^{M,\pi}(\tau = \cdot)$ denote the probability of observing $\tau \in \mathcal{T}$ when playing policy π on M . We then have

$$\begin{aligned}
 D_{\text{KL}}(\widetilde{M}(\pi) \parallel M'(\pi)) &= \int \log \frac{\mathbb{P}^{\widetilde{M},\pi}(r, \tau)}{\mathbb{P}^{M',\pi}(r, \tau)} d\mathbb{P}^{\widetilde{M},\pi}(r, \tau) \\
 &= \int \int \log \frac{\mathbb{P}^{\widetilde{M},\pi}(r \mid \tau) \mathbb{P}^{M,\pi}(\tau)}{\mathbb{P}^{M',\pi}(r \mid \tau) \mathbb{P}^{M',\pi}(\tau)} d\mathbb{P}^{\widetilde{M},\pi}(r \mid \tau) d\mathbb{P}^{M,\pi}(\tau) \\
 &= \int \log \frac{\mathbb{P}^{M,\pi}(\tau)}{\mathbb{P}^{M',\pi}(\tau)} d\mathbb{P}^{M,\pi}(\tau) + \int \left(\int \log \frac{\mathbb{P}^{\widetilde{M},\pi}(r \mid \tau)}{\mathbb{P}^{M',\pi}(r \mid \tau)} d\mathbb{P}^{\widetilde{M},\pi}(r \mid \tau) \right) d\mathbb{P}^{M,\pi}(\tau) \\
 &= \sum_{\tau \in \mathcal{T}} \mathbb{P}^{M,\pi}(\tau) \log \frac{\mathbb{P}^{M,\pi}(\tau)}{\mathbb{P}^{M',\pi}(\tau)} + \sum_{\tau \in \mathcal{T}} \mathbb{P}^{M,\pi}(\tau) D_{\text{KL}}\left(\mathbb{P}^{\widetilde{M},\pi}(r \mid \tau) \parallel \mathbb{P}^{M',\pi}(r \mid \tau)\right).
 \end{aligned}$$

It follows that

$$\begin{aligned}
 &\left| D_{\text{KL}}(\widetilde{M}(\pi) \parallel M'(\pi)) - D_{\text{KL}}(\bar{M}(\pi) \parallel M'(\pi)) \right| \\
 &\leq \left| \sum_{\tau \in \mathcal{T}} \mathbb{P}^{M,\pi}(\tau) \log \frac{\mathbb{P}^{M,\pi}(\tau)}{\mathbb{P}^{M',\pi}(\tau)} - \sum_{\tau \in \mathcal{T}} \mathbb{P}^{\bar{M},\pi}(\tau) \log \frac{\mathbb{P}^{\bar{M},\pi}(\tau)}{\mathbb{P}^{M',\pi}(\tau)} \right| \\
 &\quad + \sum_{\tau \in \mathcal{T}} |\mathbb{P}^{M,\pi}(\tau) - \mathbb{P}^{\bar{M},\pi}(\tau)| D_{\text{KL}}\left(\mathbb{P}^{M,\pi}(r \mid \tau) \parallel \mathbb{P}^{M',\pi}(r \mid \tau)\right).
 \end{aligned}$$

Note that, since rewards at each state are independent,

$$D_{\text{KL}}\left(\mathbb{P}^{M,\pi}(r \mid \tau) \parallel \mathbb{P}^{M',\pi}(r \mid \tau)\right) = \sum_{h=1}^H D_{\text{KL}}\left(\mathbb{P}^{M,\pi}(r_h \mid \tau) \parallel \mathbb{P}^{M',\pi}(r_h \mid \tau)\right) \leq H,$$

since rewards have means in $[0, 1]$ and are unit Gaussian. This implies

$$\begin{aligned}
 \sum_{\tau \in \mathcal{T}} |\mathbb{P}^{M,\pi}(\tau) - \mathbb{P}^{\bar{M},\pi}(\tau)| D_{\text{KL}}\left(\mathbb{P}^{M,\pi}(r \mid \tau) \parallel \mathbb{P}^{M',\pi}(r \mid \tau)\right) &\leq H \sum_{\tau \in \mathcal{T}} |\mathbb{P}^{M,\pi}(\tau) - \mathbb{P}^{\bar{M},\pi}(\tau)| \\
 &= H D_{\text{TV}}(M(\pi), \bar{M}(\pi)) \\
 &\leq H \sqrt{\frac{1}{2} D_{\text{KL}}(M(\pi) \parallel \bar{M}(\pi))}.
 \end{aligned}$$

Note that $\frac{d}{dx} x \log \frac{x}{y} = 1 + \log \frac{x}{y}$, so by the Mean Value Theorem we have

$$\begin{aligned} & \left| \mathbb{P}^{M,\pi}(\tau) \log \frac{\mathbb{P}^{M,\pi}(\tau)}{\mathbb{P}^{M',\pi}(\tau)} - \mathbb{P}^{\bar{M},\pi}(\tau) \log \frac{\mathbb{P}^{\bar{M},\pi}(\tau)}{\mathbb{P}^{M',\pi}(\tau)} \right| \\ & \leq \left(1 + \max \left\{ \left| \log \frac{\mathbb{P}^{M,\pi}(\tau)}{\mathbb{P}^{M',\pi}(\tau)} \right|, \left| \log \frac{\mathbb{P}^{\bar{M},\pi}(\tau)}{\mathbb{P}^{M',\pi}(\tau)} \right| \right\} \right) \cdot |\mathbb{P}^{M,\pi}(\tau) - \mathbb{P}^{\bar{M},\pi}(\tau)| \\ & \leq \left(1 + \max_{\bar{M}' \in \mathcal{M}} \max_{\tau' \in \mathcal{T}} \left| \log \frac{\mathbb{P}^{\bar{M}',\pi}(\tau')}{\mathbb{P}^{M',\pi}(\tau')} \right| \right) \cdot |\mathbb{P}^{M,\pi}(\tau) - \mathbb{P}^{\bar{M},\pi}(\tau)|. \end{aligned}$$

It follows that

$$\begin{aligned} & \left| \sum_{\tau \in \mathcal{T}} \mathbb{P}^{M,\pi}(\tau) \log \frac{\mathbb{P}^{M,\pi}(\tau)}{\mathbb{P}^{M',\pi}(\tau)} - \sum_{\tau \in \mathcal{T}} \mathbb{P}^{\bar{M},\pi}(\tau) \log \frac{\mathbb{P}^{\bar{M},\pi}(\tau)}{\mathbb{P}^{M',\pi}(\tau)} \right| \\ & \leq \left(1 + \max_{\bar{M}' \in \mathcal{M}} \max_{\tau' \in \mathcal{T}} \left| \log \frac{\mathbb{P}^{\bar{M}',\pi}(\tau')}{\mathbb{P}^{M',\pi}(\tau')} \right| \right) \cdot \sum_{\tau \in \mathcal{T}} |\mathbb{P}^{M,\pi}(\tau) - \mathbb{P}^{\bar{M},\pi}(\tau)| \\ & = \left(1 + \max_{\bar{M}' \in \mathcal{M}} \max_{\tau' \in \mathcal{T}} \left| \log \frac{\mathbb{P}^{\bar{M}',\pi}(\tau')}{\mathbb{P}^{M',\pi}(\tau')} \right| \right) \cdot D_{\text{TV}}(M(\pi), \bar{M}(\pi)) \\ & \leq \left(1 + \max_{\bar{M}' \in \mathcal{M}} \max_{\tau' \in \mathcal{T}} \left| \log \frac{\mathbb{P}^{\bar{M}',\pi}(\tau')}{\mathbb{P}^{M',\pi}(\tau')} \right| \right) \cdot \sqrt{\frac{1}{2} D_{\text{KL}}(M(\pi) \| \bar{M}(\pi))}. \end{aligned}$$

This verifies [Assumption D.1](#) with

$$L_{\text{KL}} = 1 + \sqrt{2H} + H + \max_{M' \in \mathcal{M}, \bar{M}' \in \mathcal{M}} \max_{\pi \in \Pi} \max_{\tau \in \mathcal{T}} \left| \log \frac{\mathbb{P}^{\bar{M}',\pi}(\tau)}{\mathbb{P}^{M',\pi}(\tau)} \right|.$$

Verifying [Assumption D.2](#). That $D_{\text{H}}^2(\bar{M}(\pi), \bar{M}'(\pi)) \leq D(\bar{M}(\pi) \| \bar{M}'(\pi))$ is immediate, since KL always upper bounds Hellinger squared.

Verifying [Assumption A.2](#). We have

$$\log \frac{\mathbb{P}^{\bar{M},\pi}(r, o)}{\mathbb{P}^{M,\pi}(r, o)} = \log \frac{\mathbb{P}^{\bar{M},\pi}(\tau)}{\mathbb{P}^{M,\pi}(\tau)} + \log \frac{\mathbb{P}^{\bar{M},\pi}(r | \tau)}{\mathbb{P}^{M,\pi}(r | \tau)} = \log \frac{\mathbb{P}^{\bar{M},\pi}(\tau)}{\mathbb{P}^{M,\pi}(\tau)} + \sum_{h=1}^H \log \frac{\mathbb{P}^{\bar{M},\pi}(r_h | \tau)}{\mathbb{P}^{M,\pi}(r_h | \tau)}.$$

Using the same calculation as in [Lemma G.1](#), we have that $\log \frac{\mathbb{P}^{\bar{M},\pi}(r_h | \tau)}{\mathbb{P}^{M,\pi}(r_h | \tau)}$ is 8-subgaussian, since rewards are unit-variance Gaussian. As $\log \frac{\mathbb{P}^{\bar{M},\pi}(r_h | \tau)}{\mathbb{P}^{M,\pi}(r_h | \tau)}$ and $\log \frac{\mathbb{P}^{\bar{M},\pi}(r_{h'} | \tau)}{\mathbb{P}^{M,\pi}(r_{h'} | \tau)}$ are independent for $h \neq h'$, it follows that $\log \frac{\mathbb{P}^{\bar{M},\pi}(r | \tau)}{\mathbb{P}^{M,\pi}(r | \tau)}$ is $8H$ -subgaussian.

Furthermore, bounding

$$\log \frac{\mathbb{P}^{\bar{M},\pi}(\tau)}{\mathbb{P}^{M,\pi}(\tau)} \leq \sup_{\bar{M}, M \in \mathcal{M}} \sup_{\pi \in \Pi} \sup_{\tau \in \mathcal{T}} \left| \log \frac{\mathbb{P}^{\bar{M},\pi}(\tau)}{\mathbb{P}^{M,\pi}(\tau)} \right| =: V_{\mathcal{T}},$$

we have that $\log \frac{\mathbb{P}^{\bar{M},\pi}(\tau)}{\mathbb{P}^{M,\pi}(\tau)}$ is $V_{\mathcal{T}}^2$ -subgaussian. Since the sum of subgaussian random variables is subgaussian, it follows that $\log \frac{\mathbb{P}^{\bar{M},\pi}(r, o)}{\mathbb{P}^{M,\pi}(r, o)} = \log \frac{\mathbb{P}^{\bar{M},\pi}(\tau)}{\mathbb{P}^{M,\pi}(\tau)} + \log \frac{\mathbb{P}^{\bar{M},\pi}(r | \tau)}{\mathbb{P}^{M,\pi}(r | \tau)}$ is $(V_{\mathcal{T}}^2 + 8H)$ -subgaussian, which verifies [Assumption A.2](#).

□

Lemma G.9. Let $P_{\min} := \inf_{M \in \mathcal{M}} \inf_{h, s', s, a} P_h^M(s' | s, a)$ and assume \mathcal{M} is such that $P_{\min} > 0$. We can construct a (ρ, μ) -cover of \mathcal{M} with respect to $\mathcal{E} := \{|r_h| \leq 1 + \sqrt{2 \log(2H/\mu)}, \forall h \in [H]\}$, with

$$N_{\text{cov}}(\mathcal{M}, \rho, \mu) \leq \frac{1}{\min\{\frac{\rho P_{\min}}{4H}, 2P_{\min}\} S^{2AH}} \cdot \frac{(2H(2 + \sqrt{\log(H/\mu)}))^{SAH}}{\rho^{SAH}}.$$

Proof of Lemma G.9. Throughout this proof, we use $M = \{(P_h^M)_{h=1}^H, (r_h^M)_{h=1}^H\} \in \mathcal{M}$ to denote the MDP in \mathcal{M} with $R_h^M(s, a) = \mathcal{N}(r_h^M(s, a), 1)$; for brevity, we \mathcal{S} , \mathcal{A} , and s_1 to be fixed and the dependence on them.

Observe that for any models $M, M' \in \mathcal{M}$, we have

$$\begin{aligned} \left| \log \frac{\mathbb{P}^{M, \pi}(r, o)}{\mathbb{P}^{M', \pi}(r, o)} \right| &\leq \left| \log \frac{\mathbb{P}^{M, \pi}(\tau)}{\mathbb{P}^{M', \pi}(\tau)} \right| + \left| \log \frac{\mathbb{P}^{M, \pi}(r | \tau)}{\mathbb{P}^{M', \pi}(r | \tau)} \right| \\ &= \sum_{h=1}^H \left| \log \frac{P_h^M(\tau_{h+1} | \tau_h, \pi(\tau_h, h))}{P_h^{M'}(\tau_{h+1} | \tau_h, \pi(\tau_h, h))} \right| + \sum_{h=1}^H \left| \log \frac{\mathbb{P}^{M, \pi}(r_h | \tau)}{\mathbb{P}^{M', \pi}(r_h | \tau)} \right|. \end{aligned}$$

Let $\mathcal{I}_\varepsilon = \{\varepsilon, 2\varepsilon, \dots, \lfloor 1/\varepsilon \rfloor \varepsilon\}$, so that $|\mathcal{I}_\varepsilon| \leq 1/\varepsilon$. Let \mathcal{P}_ε denote an ε cover of $\Delta_{\mathcal{S}}$ in the ℓ_∞ -norm, so that for any $P \in \Delta_{\mathcal{S}}$, there exists some $P' \in \mathcal{P}_\varepsilon$ such that $\sup_{s \in \mathcal{S}} |P_s - P'_s| \leq \varepsilon$. It suffices to choose $\mathcal{P}_\varepsilon = \mathcal{I}_\varepsilon^{\mathcal{S}} \cap \Delta_{\mathcal{S}}$, so we can bound $|\mathcal{P}_\varepsilon| \leq 1/\varepsilon^S$. Let

$$\mathcal{M}_{\text{cov}} = \{M = \{(P_h^M)_{h=1}^H, (r_h^M)_{h=1}^H\} : P_h^M(\cdot | s, a) \in \mathcal{P}_{\varepsilon_1}, r_h^M(s, a) \in \mathcal{I}_{\varepsilon_2}, \quad \forall s, a, h\}$$

for parameters $\varepsilon_1, \varepsilon_2 > 0$ to be chosen. Then

$$|\mathcal{M}_{\text{cov}}| = (|\mathcal{P}_{\varepsilon_1}| |\mathcal{I}_{\varepsilon_2}|)^{SAH} \leq \frac{1}{\varepsilon_1^{S^2 AH}} \cdot \frac{1}{\varepsilon_2^{SAH}}.$$

We will show that \mathcal{M}_{cov} is a (ρ, μ) -cover of \mathcal{M} for appropriately chosen \mathcal{E} and $\varepsilon_1, \varepsilon_2 > 0$.

Let $\mathcal{E} := \{|r_h| \leq 1 + \sqrt{2 \log(2H/\mu)}, \forall h \in [H]\}$. As we assume rewards are unit-variance Gaussian and have means in $[0, 1]$, it is straightforward to see that $\mathbb{P}[\mathcal{E}^c] \leq \mu$. Fix M and let $M' \in \mathcal{M}_{\text{cov}}$ denote the instance such that

$$|r_h^M(s, a) - r_h^{M'}(s, a)| \leq \varepsilon_2 \quad \text{and} \quad \sup_{s' \in \mathcal{S}} |P_h^M(s' | s, a) - P_h^{M'}(s' | s, a)| \leq \varepsilon_1, \quad \forall s, a, h.$$

Note that such an instance is guaranteed to exist by definition of \mathcal{M}_{cov} .

By a similar argument as in Lemma G.1, we can bound, on \mathcal{E} ,

$$\begin{aligned} \sum_{h=1}^H \left| \log \frac{\mathbb{P}^{M, \pi}(r_h | \tau)}{\mathbb{P}^{M', \pi}(r_h | \tau)} \right| &\leq \sum_{h=1}^H (1 + |r_h|) \cdot \sup_{s, a} |r_h^M(s, a) - r_h^{M'}(s, a)| \\ &\leq H(2 + \sqrt{2 \log(2H/\mu)}) \cdot \sup_{s, a, h} |r_h^M(s, a) - r_h^{M'}(s, a)| \\ &\leq H(2 + \sqrt{2 \log(2H/\mu)}) \cdot \varepsilon_2. \end{aligned}$$

We also have

$$\begin{aligned}
 \sum_{h=1}^H \left| \log \frac{P_h^M(\tau_{h+1} \mid \tau_h, \pi(\tau_h, h))}{P_h^{M'}(\tau_{h+1} \mid \tau_h, \pi(\tau_h, h))} \right| &\leq H \cdot \sup_{h, s', s, a} \left| \log \frac{P_h^M(s' \mid s, a)}{P_h^{M'}(s' \mid s, a)} \right| \\
 &\leq H \cdot \sup_{|x| \leq \varepsilon_1} \sup_{h, s', s, a} \left| \log \frac{P_h^M(s' \mid s, a)}{P_h^M(s' \mid s, a) - x} \right| \\
 &\leq H \cdot \sup_{h, s', s, a} \frac{\varepsilon_1}{P_h^M(s' \mid s, a) - \varepsilon_1}
 \end{aligned}$$

where the last inequality holds as long as $\inf_{h, s', s, a} P_h^M(s' \mid s, a) - \varepsilon_1 > 0$. Denoting $P_{\min} := \inf_{M \in \mathcal{M}} \inf_{h, s', s, a} P_h^M(s' \mid s, a)$, for \mathcal{M}_{cov} to be a (ρ, μ) -cover, it therefore suffices that

$$H(2 + \sqrt{2 \log(2H/\mu)}) \cdot \varepsilon_2 \leq \rho/2, \quad \frac{2H\varepsilon_1}{P_{\min}} \leq \rho/2, \quad \text{and} \quad P_{\min} \geq \varepsilon_1/2$$

so it suffices to take

$$\varepsilon_1 = \min\left\{\frac{\rho P_{\min}}{4H}, 2P_{\min}\right\} \quad \text{and} \quad \varepsilon_2 = \frac{\rho}{2H(2 + \sqrt{\log(H/\mu)})}.$$

The result now follows from our bound on $|\mathcal{M}_{\text{cov}}|$. □

G.5.3. TABULAR MDPs HAVE BOUNDED UNIFORM EXPLORATION COEFFICIENT

Lemma G.10. *For the tabular MDP class \mathcal{M} in (27), we can bound, for all $\varepsilon > 0$,*

$$C_{\text{exp}}^{\text{D}}(\mathcal{M}, \varepsilon) \leq \frac{320000SAH^2 \cdot \log^2 H}{\varepsilon^2}.$$

for $D(\cdot \parallel \cdot) \leftarrow D_{\text{H}}^2(\cdot, \cdot)$.

Proof of Lemma G.10. Let $\xi \in \Delta(\mathcal{M})$ be given. Define

$$p_{\text{exp}} = \arg \min_{p \in \Delta_{\Pi}} \max_{q \in \Delta_{\Pi}} \sum_{s, a, h} \mathbb{E}_{\pi \sim q} \left[\mathbb{E}_{\bar{M} \sim \xi} \left[\frac{(w_h^{\bar{M}, \pi}(s, a))^2}{\mathbb{E}_{\pi' \sim p}[w_h^{\bar{M}, \pi'}(s, a)]} \right] \right].$$

We first show that, for any $M \in \mathcal{M}$ and any π ,

$$\mathbb{E}_{\bar{M} \sim \xi}[D_{\text{H}}^2(\bar{M}(\pi), M(\pi))] \leq \sqrt{SAH^2 \cdot \mathbb{E}_{\bar{M} \sim \xi}[\mathbb{E}_{\pi \sim p_{\text{exp}}}[D_{\text{H}}^2(\bar{M}(\pi), M(\pi))]]}.$$

Consider any policy π . We can bound

$$\begin{aligned}
 &\mathbb{E}_{\bar{M} \sim \xi}[D_{\text{H}}^2(\bar{M}(\pi), M(\pi))] \\
 &\stackrel{(a)}{\leq} 100 \log(H) \cdot \sum_{s, a, h} \mathbb{E}_{\bar{M} \sim \xi} \left[w_h^{\bar{M}, \pi}(s, a) D_{\text{H}}^2(\bar{M}_{sh}(a), M_{sh}(a)) \right]
 \end{aligned}$$

$$\begin{aligned}
 &\stackrel{(b)}{\leq} 100 \log(H) \cdot \sqrt{\sum_{s,a,h} \mathbb{E}_{\bar{M} \sim \xi} \left[\frac{w_h^{\bar{M},\pi}(s,a)^2}{\mathbb{E}_{\pi' \sim p_{\text{exp}}} [w_h^{\bar{M},\pi'}(s,a)]} \right]} \\
 &\quad \cdot \sqrt{\sum_{s,a,h} \mathbb{E}_{\pi' \sim p_{\text{exp}}} \left[\mathbb{E}_{\bar{M} \sim \xi} \left[w_h^{\bar{M},\pi'}(s,a) D_{\mathbb{H}}^4(\bar{M}_{sh}(a), M_{sh}(a)) \right] \right]} \\
 &\stackrel{(c)}{\leq} 200 \log(H) \cdot \sqrt{\sum_{s,a,h} \mathbb{E}_{\bar{M} \sim \xi} \left[\frac{w_h^{\bar{M},\pi}(s,a)^2}{\mathbb{E}_{\pi' \sim p_{\text{exp}}} [w_h^{\bar{M},\pi'}(s,a)]} \right]} \\
 &\quad \cdot \sqrt{\sum_{s,a,h} \mathbb{E}_{\pi' \sim p_{\text{exp}}} \left[\mathbb{E}_{\bar{M} \sim \xi} \left[w_h^{\bar{M},\pi'}(s,a) D_{\mathbb{H}}^2(\bar{M}_{sh}(a), M_{sh}(a)) \right] \right]}
 \end{aligned}$$

where (a) follows from Lemma A.13 of Foster et al. (2021), (b) follows from Cauchy-Schwarz and Jensen's inequality, and (c) follows because the Hellinger distance is always bounded by 2. Now note that, by definition of p_{exp} , we have

$$\sum_{s,a,h} \mathbb{E}_{\bar{M} \sim \xi} \left[\frac{w_h^{\bar{M},\pi}(s,a)^2}{\mathbb{E}_{\pi' \sim p_{\text{exp}}} [w_h^{\bar{M},\pi'}(s,a)]} \right] \leq \min_{p \in \Delta_{\Pi}} \max_{q \in \Delta_{\Pi}} \sum_{s,a,h} \mathbb{E}_{\pi \sim q} \left[\mathbb{E}_{\bar{M} \sim \xi} \left[\frac{(w_h^{\bar{M},\pi}(s,a))^2}{\mathbb{E}_{\pi' \sim p} [w_h^{\bar{M},\pi'}(s,a)]} \right] \right]$$

and by the minimax theorem, we can bound

$$\begin{aligned}
 \min_{p \in \Delta_{\Pi}} \max_{q \in \Delta_{\Pi}} \sum_{s,a,h} \mathbb{E}_{\pi \sim q} \left[\mathbb{E}_{\bar{M} \sim \xi} \left[\frac{(w_h^{\bar{M},\pi}(s,a))^2}{\mathbb{E}_{\pi' \sim p} [w_h^{\bar{M},\pi'}(s,a)]} \right] \right] &= \max_{q \in \Delta_{\Pi}} \min_{p \in \Delta_{\Pi}} \sum_{s,a,h} \mathbb{E}_{\pi \sim q} \left[\mathbb{E}_{\bar{M} \sim \xi} \left[\frac{(w_h^{\bar{M},\pi}(s,a))^2}{\mathbb{E}_{\pi' \sim p} [w_h^{\bar{M},\pi'}(s,a)]} \right] \right] \\
 &\leq \max_{q \in \Delta_{\Pi}} \sum_{s,a,h} \mathbb{E}_{\pi \sim q} \left[\mathbb{E}_{\bar{M} \sim \xi} \left[\frac{(w_h^{\bar{M},\pi}(s,a))^2}{\mathbb{E}_{\pi' \sim q} [w_h^{\bar{M},\pi'}(s,a)]} \right] \right] \\
 &\leq \max_{q \in \Delta_{\Pi}} \sum_{s,a,h} \mathbb{E}_{\pi \sim q} \left[\mathbb{E}_{\bar{M} \sim \xi} \left[\frac{w_h^{\bar{M},\pi}(s,a)}{\mathbb{E}_{\pi' \sim q} [w_h^{\bar{M},\pi'}(s,a)]} \right] \right] \\
 &= SAH.
 \end{aligned}$$

By Lemma A.9 of Foster et al. (2021), since $\mu^{\bar{M}}(s,a,h) := \frac{1}{H} w_h^{\bar{M},\pi'}(s,a)$ forms a valid distribution on $\mathcal{S} \times \mathcal{A} \times [H]$, we can upper bound

$$\sum_{s,a,h} w_h^{\bar{M},\pi'}(s,a) D_{\mathbb{H}}^2(\bar{M}_{sh}(a), M_{sh}(a)) \leq H D_{\mathbb{H}}^2(\bar{M}(\pi'), M(\pi')).$$

Altogether then, we have shown that for all $\pi \in \Pi$,

$$\mathbb{E}_{\bar{M} \sim \xi} [D_{\mathbb{H}}^2(\bar{M}(\pi), M(\pi))] \leq 200 \log(H) \cdot \sqrt{SAH^2 \cdot \mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{\pi \sim p_{\text{exp}}} [D_{\mathbb{H}}^2(\bar{M}(\pi), M(\pi))]]}$$

as desired. Since the Hellinger distance is a metric and satisfies the triangle inequality, this in particular implies that, for any M, M' ,

$$D_{\mathbb{H}}^2(M(\pi), M'(\pi)) \leq 2 \mathbb{E}_{\bar{M} \sim \xi} [D_{\mathbb{H}}^2(\bar{M}(\pi), M(\pi))] + 2 \mathbb{E}_{\bar{M} \sim \xi} [D_{\mathbb{H}}^2(\bar{M}(\pi), M'(\pi))]$$

$$\begin{aligned} &\leq 400 \log H \cdot \sqrt{SAH^2 \cdot \mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{\pi \sim p_{\text{exp}}} [D_{\mathbb{H}}^2(\bar{M}(\pi), M(\pi))]]} \\ &\quad + 400 \log H \sqrt{SAH^2 \cdot \mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{\pi \sim p_{\text{exp}}} [D_{\mathbb{H}}^2(\bar{M}(\pi), M'(\pi))]]}. \end{aligned}$$

Thus, if

$$\mathbb{E}_{\bar{M} \sim \xi} [\mathbb{E}_{\pi \sim p_{\text{exp}}} [D_{\mathbb{H}}^2(\bar{M}(\pi), M''(\pi))]] \leq \frac{\varepsilon^2}{320000SAH^2 \cdot \log^2 H}$$

for both $M'' \in \{M, M'\}$, then $D_{\mathbb{H}}^2(M(\pi), M'(\pi)) \leq \varepsilon$. It follows that a sufficient choice for $C_{\text{exp}}^{\mathbb{D}}$ is $320000SAH^2 \log^2 H / \varepsilon^2$. \square

G.5.4. SUPPORTING LEMMAS

Lemma G.11. *If M has a unique optimal policy π_M and $M' \in \mathcal{M}^{\text{alt}}(M)$, then there exists some $(\tilde{s}, \tilde{a}, \tilde{h})$ such that*

$$Q_{\tilde{h}}^{M', \pi_M}(\tilde{s}, \tilde{a}) > V_{\tilde{h}}^{M', \pi_M}(\tilde{s}).$$

Proof of Lemma G.11. Assume that this is not the case, i.e. that for all (s, a, h) ,

$$Q_h^{M', \pi_M}(s, a) \leq V_h^{M', \pi_M}(s) = Q_h^{M', \pi_M}(s, \pi_M(s, h)).$$

Our goal is to show that in this case $\pi_{M'} = \pi_M$, which contradicts the fact that $M' \in \mathcal{M}^{\text{alt}}(M)$. We proceed by induction.

Base Case. Let $h = H$ and assume that for all (s, a) ,

$$Q_H^{M', \pi_M}(s, a) \leq Q_H^{M', \pi_M}(s, \pi_M(s, h)).$$

This contradicts the assumption that π_M is unique.

Inductive Case. Assume that $\pi_{M'}(s, h') = \pi_M(s, h')$ for all s and $h' > h$. This then implies that $V_{h+1}^{M', \pi_M}(s) = V_{h+1}^{M', \pi_{M'}}(s)$ for all s . It follows that for all a

$$Q_h^{M', \pi_M}(s, a) = r_h^{M'}(s, a) + \mathbb{P}_h^{M'} [V_{h+1}^{M', \pi_M}](s, a) = r_h^{M'}(s, a) + \mathbb{P}_h^{M'} [V_{h+1}^{M', \pi_{M'}}](s, a) = Q_h^{M', \pi_{M'}}(s, a)$$

so in particular $Q_h^{M', \pi_M}(s, \pi_M(s, h)) = Q_h^{M', \pi_{M'}}(s, \pi_M(s, h))$. Since we have assumed that for all (s, a)

$$Q_h^{M', \pi_M}(s, a) \leq Q_h^{M', \pi_M}(s, \pi_M(s, h))$$

we have

$$Q_h^{M', \pi_{M'}}(s, \pi_{M'}(s, h)) \leq Q_h^{M', \pi_{M'}}(s, \pi_M(s, h)).$$

However, since each $M \in \mathcal{M}$ has a unique optimal action at each state, this is a contradiction unless $\pi_{M'}(s, h) = \pi_M(s, h)$, which proves the inductive hypothesis. The result follows. \square

Lemma G.12. For MDPs M, M' with unit variance Gaussian rewards, we have

$$V_h^{M',\pi}(s) - V_h^{M,\pi}(s) \leq \sqrt{\frac{8H}{w_h^{M,\pi}(s)} \cdot D_{\text{KL}}(M(\pi) \| M'(\pi))}.$$

Proof of Lemma G.12. In the Gaussian reward setting, we have

$$r_{h'}^{M'}(s', a') - r_{h'}^M(s', a') \leq \sqrt{(r_{h'}^{M'}(s', a') - r_{h'}^M(s', a'))^2} \leq \sqrt{2D_{\text{KL}}(M_{h',s'}(a') \| M'_{h',s'}(a'))}.$$

Furthermore, since $V_{h'+1}^{M',\pi}(s') \in [0, 1]$, we have

$$\begin{aligned} \mathbb{P}_{h'}^{M'}[V_{h'+1}^{M',\pi}(s', a')] - \mathbb{P}_{h'}^M[V_{h'+1}^{M',\pi}(s', a')] &\leq \sum_{s''} |P_{h'}^{M'}(s'' | s', a') - P_{h'}^M(s'' | s', a')| \\ &= 2D_{\text{TV}}(P_{h'}^{M'}(\cdot | s', a'), P_{h'}^M(\cdot | s', a')) \\ &\leq \sqrt{2D_{\text{KL}}(P_{h'}^{M'}(\cdot | s', a') \| P_{h'}^M(\cdot | s', a'))} \\ &\leq \sqrt{2D_{\text{KL}}(M_{h',s'}(a') \| M'_{h',s'}(a'))}. \end{aligned}$$

By Lemma G.14, Jensen's inequality, and Lemma G.15, it follows that

$$\begin{aligned} V_h^{M',\pi}(s) - V_h^{M,\pi}(s) &\leq \sum_{h'=h}^H \sum_{s', a'} w_{h'}^{M,\pi}(s', a' | s_h = s) \cdot 2\sqrt{2D_{\text{KL}}(M_{h',s'}(a') \| M'_{h',s'}(a'))} \\ &\leq 2\sqrt{2H \sum_{h'=h}^H \sum_{s', a'} w_{h'}^{M,\pi}(s', a' | s_h = s) D_{\text{KL}}(M_{h',s'}(a') \| M'_{h',s'}(a'))} \\ &\leq 2\sqrt{\frac{2H}{w_h^{M,\pi}(s)} \cdot \sum_{h'=h}^H \sum_{s', a'} w_{h'}^{M,\pi}(s', a') D_{\text{KL}}(M_{h',s'}(a') \| M'_{h',s'}(a'))} \\ &\leq 2\sqrt{\frac{2H}{w_h^{M,\pi}(s)} D_{\text{KL}}(M(\pi) \| M'(\pi))} \end{aligned}$$

where we have used that, for $h < h'$,

$$w_{h'}^{M,\pi}(s', a') = \sum_{s''} w_{h'}^{M,\pi}(s', a' | s_h = s'') w_{\pi}^M(s'', h) \geq w_{h'}^{M,\pi}(s', a' | s_h = s) w_h^{M,\pi}(s).$$

□

Lemma G.13. For any $M \in \mathcal{M}$ for which π_M is unique, we have

$$\Delta_{\min}^M \leq \min_{s,h} w_h^{M,\pi_M}(s).$$

Proof of Lemma G.13. Let $\tilde{\pi}$ denote the policy that differs from policy π_M only at the state \tilde{s} and layer \tilde{h} given by $(\tilde{s}, \tilde{h}) = \arg \min_{s, h} w_h^{M, \pi_M}(s)$. Note that this implies, in particular, that $w_{\tilde{h}}^{M, \pi_M}(\tilde{s}) = w_{\tilde{h}}^{M, \tilde{\pi}}(\tilde{s})$ since $\tilde{\pi}$ and π_M take identical actions up to step \tilde{h} . By the Performance-Difference Lemma (Kakade, 2003), we have

$$\begin{aligned} V_1^{M, \pi_M} - V_1^{M, \tilde{\pi}} &= \sum_{h=1}^H \sum_{s, a} w_h^{M, \tilde{\pi}}(s, a) (V_h^{M, \pi_M}(s) - Q_h^{M, \pi_M}(s, a)) \\ &\stackrel{(a)}{=} w_{\tilde{h}}^{M, \tilde{\pi}}(\tilde{s}, \tilde{\pi}(\tilde{s}, \tilde{h})) (V_{\tilde{h}}^{M, \pi_M}(\tilde{s}) - Q_{\tilde{h}}^{M, \pi_M}(\tilde{s}, \tilde{\pi}(\tilde{s}, \tilde{h}))) \\ &= w_{\tilde{h}}^{M, \pi_M}(\tilde{s}) (V_{\tilde{h}}^{M, \pi_M}(\tilde{s}) - Q_{\tilde{h}}^{M, \pi_M}(\tilde{s}, \tilde{\pi}(\tilde{s}, \tilde{h}))) \\ &\leq w_{\tilde{h}}^{M, \pi_M}(\tilde{s}) \end{aligned}$$

where (a) holds since $V_h^{M, \pi_M}(s) = Q_h^{M, \pi_M}(s, a)$ for all (s, a, h) with $w_h^{M, \tilde{\pi}}(s, a, h) > 0$ other than at (\tilde{s}, \tilde{h}) . By assumption, the optimal policy is unique, so $V_1^{M, \pi_M} - V_1^{M, \tilde{\pi}} > 0$, and thus

$$\Delta_{\min}^M = \min_{\pi \in \Pi : V_1^{M, \pi_M} - V_1^{M, \pi} > 0} V_1^{M, \pi_M} - V_1^{M, \pi} \leq V_1^{M, \pi_M} - V_1^{M, \tilde{\pi}} \leq w_{\tilde{h}}^{M, \pi_M}(\tilde{s}) = \min_{s, h} w_h^{M, \pi_M}(s).$$

□

Lemma G.14 (Lemma E.15 of Dann et al. (2017)). *For MDPs M, M' and policy π , we have*

$$\begin{aligned} V_h^{M', \pi}(s) - V_h^{M, \pi}(s) &= \sum_{h'=h}^H \sum_{s', a'} w_{h'}^{M, \pi}(s', a' | s_h = s) \cdot \left[r_{h'}^{M'}(s', a') - r_{h'}^M(s', a') + \mathbb{P}_{h'}^{M'}[V_{h'+1}^{M', \pi}](s', a') \right. \\ &\quad \left. - \mathbb{P}_{h'}^M[V_{h'+1}^{M, \pi}](s', a') \right]. \end{aligned}$$

Lemma G.15. *For MDPs M, M' and policy π , we have*

$$D_{\text{KL}}(M(\pi) \| M'(\pi)) = \sum_{h=1}^H \sum_{s, a} w_h^{M, \pi}(s, a) D_{\text{KL}}(M_{hs}(a) \| M'_{hs}(a)).$$

Proof of Lemma G.15. This is a standard calculation; see e.g. (Simchowitz and Jamieson, 2019; Tirinzoni et al., 2021). □

Appendix H. Proofs and Additional Results from Appendix B

H.1. Technical Lemmas

Throughout this section, when the class \mathcal{M} is clear from context, we define

$$\Lambda(M; \varepsilon, \bar{n}) := \{\lambda \in \Delta_{\Pi} : \exists n \leq \bar{n} \text{ s.t. } \Delta^M(\lambda) \leq (1 + \varepsilon)g^M/n, I^M(\lambda; \mathcal{M}) \geq (1 - \varepsilon)/n\}. \quad (77)$$

Lemma H.1 (Derandomization). *Let $\bar{n} > 0$ be given. For any $p \in \Delta(\Delta(\Pi))$, defining $\bar{\lambda}_p = \mathbb{E}_{\lambda \sim p}[\lambda] \in \Delta(\Pi)$, we have that for all $M \in \mathcal{M}$,*

$$\mathbb{I}\{\bar{\lambda}_p \notin \Lambda(M; \varepsilon, \bar{n})\} \leq \delta^{-1} \cdot \mathbb{P}_{\lambda \sim p}[\lambda \notin \Lambda(M; \varepsilon/2, \bar{n})],$$

where $\delta := \frac{\varepsilon}{2} \cdot \min\left\{1, \frac{\mathbf{g}^M}{\bar{n}}\right\}$.

Proof of Lemma H.1. Let $M \in \mathcal{M}$ be fixed and abbreviate $I^M(\lambda) = I^M(\lambda; \mathcal{M})$. Fix $p \in \Delta(\Delta(\Pi))$. For any $\lambda \in \Lambda(M; \varepsilon/2, \bar{n})$, let $n_\lambda > 0$ denote the least $n > 0$ such that

$$\Delta^M(\lambda) \leq (1 + \varepsilon/2)\mathbf{g}^M/n, \quad \text{and} \quad I^M(\lambda; \mathcal{M}) \geq (1 - \varepsilon/2)/n.$$

Define

$$n = \left(\mathbb{E}_{\lambda \sim p} \left[\frac{1}{n_\lambda} \mid \lambda \in \Lambda(M; \varepsilon/2, \bar{n}) \right] \right)^{-1},$$

and note that by Jensen's inequality,

$$n \leq \mathbb{E}_{\lambda \sim p}[n_\lambda \mid \lambda \in \Lambda(M; \varepsilon/2, \bar{n})] \leq \bar{n}.$$

We first observe that since $\Delta^M \in [0, 1]$,

$$\begin{aligned} \Delta^M(\bar{\lambda}_p) &\leq \mathbb{E}_{\lambda \sim p}[\Delta^M(\lambda) \mid \lambda \in \Lambda(M; \varepsilon/2, \bar{n})] + \mathbb{P}_{\lambda \sim p}[\lambda \notin \Lambda(M; \varepsilon/2, \bar{n})] \\ &\leq (1 + \varepsilon/2)\mathbf{g}^M \cdot \mathbb{E}_{\lambda \sim p} \left[\frac{1}{n_\lambda} \mid \lambda \in \Lambda(M; \varepsilon/2, \bar{n}) \right] + \mathbb{P}_{\lambda \sim p}[\lambda \notin \Lambda(M; \varepsilon/2, \bar{n})] \\ &= (1 + \varepsilon/2)\frac{\mathbf{g}^M}{n} + \mathbb{P}_{\lambda \sim p}[\lambda \notin \Lambda(M; \varepsilon/2, \bar{n})]. \end{aligned}$$

Next, note that the map $\lambda \mapsto I^M(\lambda)$ is concave and non-negative (it is an infimum over non-negative linear functions), so we have

$$\begin{aligned} I^M(\bar{\lambda}_p) &\geq \mathbb{E}_{\lambda \sim p}[I^M(\lambda)] \\ &\geq \mathbb{E}_{\lambda \sim p}[I^M(\lambda) \mid \lambda \in \Lambda(M; \varepsilon/2, \bar{n})] \cdot \mathbb{P}_{\lambda \sim p}[\lambda \in \Lambda(M; \varepsilon/2, \bar{n})] \\ &\geq (1 - \varepsilon/2) \mathbb{E}_{\lambda \sim p} \left[\frac{1}{n_\lambda} \mid \lambda \in \Lambda(M; \varepsilon/2, \bar{n}) \right] \cdot \mathbb{P}_{\lambda \sim p}[\lambda \in \Lambda(M; \varepsilon/2, \bar{n})] \\ &= (1 - \varepsilon/2) \frac{1}{n} \cdot \mathbb{P}_{\lambda \sim p}[\lambda \in \Lambda(M; \varepsilon/2, \bar{n})] \\ &= (1 - \varepsilon/2) \frac{1}{n} \cdot (1 - \mathbb{P}_{\lambda \sim p}[\lambda \notin \Lambda(M; \varepsilon/2, \bar{n})]). \end{aligned}$$

It follows that as long as

$$\mathbb{P}_{\lambda \sim p}[\lambda \notin \Lambda(M; \varepsilon/2, \bar{n})] \leq \delta := \frac{\varepsilon}{2} \cdot \min\left\{1, \frac{\mathbf{g}^M}{\bar{n}}\right\} \leq \frac{\varepsilon}{2} \cdot \min\left\{1, \frac{\mathbf{g}^M}{n}\right\},$$

we have

$$\Delta^M(\bar{\lambda}_p) \leq (1 + \varepsilon)\mathbf{g}^M/n, \quad \text{and} \quad I^M(\bar{\lambda}_p; \mathcal{M}) \geq (1 - \varepsilon)/n,$$

so that $\bar{\lambda}_p \in \Lambda(M; \varepsilon, n) \subseteq \Lambda(M; \varepsilon, \bar{n})$. We conclude that

$$\mathbb{I}\{\bar{\lambda}_p \notin \Lambda(M; \varepsilon, \bar{n})\} \leq \mathbb{I}\{\mathbb{P}_{\lambda \sim p}[\lambda \notin \Lambda(M; \varepsilon/2, \bar{n})] > \delta\} \leq \delta^{-1} \cdot \mathbb{P}_{\lambda \sim p}[\lambda \notin \Lambda(M; \varepsilon/2, \bar{n})].$$

□

Lemma H.2. Let $M \in \mathcal{M}$ be given, and suppose [Assumption A.4](#) holds. Fix $T \in \mathbb{N}$ and consider an algorithm \mathbb{A} such that for all $M' \in \mathcal{M}^{\text{alt}}(M) \cup \{M\}$,

$$\mathbb{E}^{M', \mathbb{A}}[\mathbf{Reg}(T)] \leq R^{M'} \cdot \log(T).$$

for some bound $R^M \geq 2$. Define $\eta^M \in \mathbb{R}_+^{\Pi}$ via

$$\eta^M(\pi) = \mathbb{E}^{M, \mathbb{A}} \left[\frac{T(\pi)}{\log(T)} \right].$$

Then if

$$\log(T) \geq \frac{6}{\varepsilon} \log \left(\sup_{M' \in \mathcal{M}^{\text{alt}}(M) \cup \{M\}} \frac{R^{M'}}{\Delta_{\min}^{M'}} \cdot \log(T) \right),$$

we must have

$$I^M(\eta^M; \mathcal{M}) \geq (1 - \varepsilon).$$

Proof of Lemma H.2. Throughout this proof we will use that π_M is uniquely defined for all $M \in \mathcal{M}$ by [Assumption A.4](#). Note that

$$\mathbb{E}^{M, \mathbb{A}}[\mathbf{Reg}(T)] \leq R^M \log T \implies \sum_{\pi \neq \pi_M} \mathbb{E}^{M, \mathbb{A}}[T(\pi)] \leq \frac{R^M \log T}{\Delta_{\min}^M}.$$

Fix some $M' \in \mathcal{M}^{\text{alt}}(M)$. Then $\pi_M \neq \pi_{M'}$ (recall that under [Assumption A.4](#), each $M \in \mathcal{M}$ has a unique optimal), so

$$\mathbb{E}^{M, \mathbb{A}}[T(\pi_{M'})] \leq \frac{R^M \log T}{\Delta_{\min}^M}, \quad \mathbb{E}^{M', \mathbb{A}}[T(\pi_{M'})] \geq T - \frac{R^{M'} \log T}{\Delta_{\min}^{M'}}.$$

By Lemma H.1 of [Simchowitz and Jamieson \(2019\)](#), we have that for any \mathcal{H}^T -measurable variable $Z \in [0, 1]$, that

$$\sum_{\pi} \mathbb{E}^{M, \mathbb{A}}[T(\pi)] D_{\text{KL}}(M(\pi), M'(\pi)) \geq d(\mathbb{E}^{M, \mathbb{A}}[Z], \mathbb{E}^{M', \mathbb{A}}[Z])$$

for $d(x, y) = x \log \frac{x}{y} + (1 - x) \log \frac{1-x}{1-y}$. Choosing $Z = T(\pi_{M'})/T$, and using that

$$d(x, y) \geq (1 - x) \log \frac{1}{1 - y} - \log 2,$$

we have

$$\begin{aligned} \sum_{\pi} \mathbb{E}^{M, \mathbb{A}}[T(\pi)] D_{\text{KL}}(M(\pi), M'(\pi)) &\geq \left(1 - \frac{R^M \log T}{\Delta_{\min}^M T} \right) \log \frac{T}{T - (T - \frac{R^{M'} \log T}{\Delta_{\min}^{M'}})} - \log 2 \\ &= \left(1 - \frac{R^M \log T}{\Delta_{\min}^M T} \right) \left(\log T - \log \frac{R^{M'} \log T}{\Delta_{\min}^{M'}} \right) - \log 2. \end{aligned}$$

Now, if

$$\log(T) \geq \frac{2}{\varepsilon} \log \left(\frac{\sup_{M' \in \mathcal{M}^{\text{alt}}(M) \cup \{M\}} R^{M'} \log T}{\Delta_{\min}^{M'}} \vee 2 \right),$$

we have $\log(2) \leq \varepsilon \log(T)$,

$$\log T - \log \left(\frac{R^{M'} \log T}{\Delta_{\min}^{M'}} \right) \geq (1 - \varepsilon) \log(T),$$

and $1 - \frac{R^M \log T}{\Delta_{\min}^M T} \geq (1 - \varepsilon)$, so we can lower bound

$$\sum_{\pi} \mathbb{E}^{M, \mathbb{A}} [T(\pi)] D_{\text{KL}}(M(\pi), M'(\pi)) \geq ((1 - \varepsilon)^2 - \varepsilon) \log(T) \geq (1 - 3\varepsilon) \log(T).$$

As this is true for every $M' \in \mathcal{M}^{\text{alt}}(M)$, the result follows. \square

Lemma B.1. *Let $\varepsilon \in (0, 2)$, and suppose that [Assumption A.4](#) holds. Fix $T \in \mathbb{N}$ and consider an algorithm \mathbb{A} such that for all $M \in \mathcal{M}$,*

$$\mathbb{E}^{M, \mathbb{A}} [\mathbf{Reg}(T)] \leq (1 + \varepsilon) \mathbf{g}^M \cdot \log(T).$$

For each $M \in \mathcal{M}$, define $\eta^M \in \mathbb{R}_+^{\Pi}$ via $\eta^M(\pi) = \mathbb{E}^{M, \mathbb{A}} \left[\frac{T(\pi)}{\log(T)} \right]$, where $T(\pi)$ denotes the number of pulls of decision π , and define $\lambda^M = \eta^M / \|\eta^M\|_1$. Then if

$$\log(T) \geq \frac{6}{\varepsilon} \log \left(\sup_{M \in \mathcal{M}} \frac{2\mathbf{g}^M}{\Delta_{\min}^M} \cdot \log(T) \right),$$

we have that for all $M \in \mathcal{M}$,

$$\lambda^M \in \Lambda(M; \varepsilon). \tag{35}$$

Proof of Lemma B.1. Immediate consequence of [Lemma H.2](#). \square

H.2. Proof of Theorem B.1

Theorem H.1 (Full Statement of Theorem B.1). *Let $\varepsilon > 0$ and $\mathcal{M}_0 \subseteq \mathcal{M}$ be given. Let $\{\bar{\mathbf{n}}^M\}_{M \in \mathcal{M}_0}$ be a collection of non-negative scalars indexed by \mathcal{M}_0 , and set $\delta := \frac{\varepsilon}{2} \cdot \min \left\{ 1, \inf_{M \in \mathcal{M}_0} \frac{\mathbf{g}^M}{\bar{\mathbf{n}}^M} \right\}$. Unless*

$$T > \frac{\delta}{8} \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0, \bar{M}),$$

any algorithm must have, for some $M \in \mathcal{M}_0$:

$$\mathbb{P}^{M, \mathbb{A}} \left[\hat{\lambda} \notin \Lambda(M; \varepsilon, \bar{\mathbf{n}}^M) \right] \geq \frac{\delta}{6}.$$

Proof of Theorem H.1. Fix $\varepsilon > 0$, and let an algorithm \mathbb{A} be given. For any $\bar{M} \in \mathcal{M}^+$, define $q^{\bar{M}} = \mathbb{P}^{\bar{M}, \mathbb{A}}(\widehat{\lambda} = \cdot) \in \Delta(\Delta(\Pi))$, and let $\omega^{\bar{M}} := \mathbb{E}^{\bar{M}, \mathbb{A}}\left[\frac{1}{T} \sum_{t=1}^T p^t\right] \in \Delta(\Pi)$.

Fix $\alpha > 0$ and $\bar{M} \in \mathcal{M}^+$ be fixed. Define

$$M = \arg \max_{M \in \mathcal{M}_0} \left\{ \mathbb{P}_{\lambda \sim q^{\bar{M}}}[\lambda \notin \Lambda(M; \varepsilon, \bar{n}^M)] \mid \mathbb{E}_{\pi \sim \omega^{\bar{M}}}[D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\};$$

we assume that such an $M \in \mathcal{M}_0$ does exist, as otherwise the claim we will prove is trivial. It is immediate from this definition that we have

$$\begin{aligned} & \mathbb{P}^{\bar{M}, \mathbb{A}}[\widehat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M)] \\ &= \mathbb{P}_{\lambda \sim q^{\bar{M}}}[\lambda \notin \Lambda(M; \varepsilon, \bar{n}^M)] \\ &= \sup_{M \in \mathcal{M}_0} \left\{ \mathbb{P}_{\lambda \sim q^{\bar{M}}}[\lambda \notin \Lambda(M; \varepsilon, \bar{n}^M)] \mid \mathbb{E}_{\pi \sim \omega^{\bar{M}}}[D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\} \\ &\geq \inf_{q \in \Delta(\Delta(\Pi)), \omega \in \Delta(\Pi)} \sup_{M \in \mathcal{M}_0} \left\{ \mathbb{P}_{\lambda \sim q}[\lambda \notin \Lambda(M; \varepsilon, \bar{n}^M)] \mid \mathbb{E}_{\pi \sim \omega}[D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\} \\ &=: \text{opt}, \end{aligned}$$

with the convention that this value is zero if the set $\{M \in \mathcal{M}_0 \mid \mathbb{E}_{\pi \sim \omega}[D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2\}$ is empty. In addition, we have

$$\mathbb{E}_{\pi \sim \omega^{\bar{M}}}[D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2. \quad (78)$$

Now, define $\delta := \frac{\varepsilon}{2} \cdot \min\left\{1, \inf_{M \in \mathcal{M}_0} \frac{\mathfrak{g}^M}{\bar{n}^M}\right\}$, and let $\bar{\lambda}_q := \mathbb{E}_{\lambda \sim q}[\lambda]$. By Lemma H.1, we have

$$\begin{aligned} \text{opt} &= \inf_{q \in \Delta(\Delta(\Pi)), \omega \in \Delta(\Pi)} \sup_{M \in \mathcal{M}_0} \left\{ \mathbb{P}_{\lambda \sim q}[\lambda \notin \Lambda(M; \varepsilon, \bar{n}^M)] \mid \mathbb{E}_{\pi \sim \omega}[D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\} \\ &\geq \delta \cdot \inf_{q \in \Delta(\Delta(\Pi)), \omega \in \Delta(\Pi)} \sup_{M \in \mathcal{M}_0} \left\{ \mathbb{I}\{\bar{\lambda}_q \notin \Lambda(M; 2\varepsilon, \bar{n}^M)\} \mid \mathbb{E}_{\pi \sim \omega}[D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\} \\ &\geq \delta \cdot \inf_{\lambda \in \Delta(\Pi), \omega \in \Delta(\Pi)} \sup_{M \in \mathcal{M}_0} \left\{ \mathbb{I}\{\lambda \notin \Lambda(M; 2\varepsilon, \bar{n}^M)\} \mid \mathbb{E}_{\pi \sim \omega}[D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\} \\ &\geq \delta \cdot \inf_{\lambda \in \Delta(\Pi), \omega \in \Delta(\Pi)} \sup_{M \in \mathcal{M}_0} \left\{ \mathbb{I}\{\lambda \notin \Lambda(M; 2\varepsilon)\} \mid \mathbb{E}_{\pi \sim \omega}[D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\} \\ &= \delta \cdot \mathbb{I}\left\{\alpha^2 \geq (\text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0, \bar{M}))^{-1}\right\}. \end{aligned} \quad (79)$$

We conclude that

$$\mathbb{P}^{\bar{M}, \mathbb{A}}[\widehat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M)] \geq \delta \cdot \mathbb{I}\left\{\alpha^2 \geq (\text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0, \bar{M}))^{-1}\right\}. \quad (80)$$

To proceed, using Lemma A.11 of Foster et al. (2021), we have

$$\begin{aligned} \mathbb{P}^{\bar{M}, \mathbb{A}}[\widehat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M)] &\geq \frac{1}{3} \mathbb{P}^{\bar{M}, \mathbb{A}}[\widehat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M)] - \frac{4}{3} D_{\text{KL}}(\mathbb{P}^{\bar{M}, \mathbb{A}} \parallel \mathbb{P}^{M, \mathbb{A}}) \\ &\geq \frac{\delta}{3} \mathbb{I}\left\{\alpha^2 \geq (\text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0, \bar{M}))^{-1}\right\} - \frac{4}{3} D_{\text{KL}}(\mathbb{P}^{\bar{M}, \mathbb{A}} \parallel \mathbb{P}^{M, \mathbb{A}}). \end{aligned}$$

Using (78) gives

$$D_{\text{KL}}(\mathbb{P}^{\bar{M}, \mathbb{A}} \parallel \mathbb{P}^{M, \mathbb{A}}) = \mathbb{E}^{M, \mathbb{A}} \left[\sum_{t=1}^T \mathbb{E}_{\pi \sim p^t} D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi)) \right] = T \cdot \mathbb{E}_{\pi \sim \omega_{\bar{M}}} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 T,$$

so we have

$$\mathbb{P}^{M, \mathbb{A}} \left[\hat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M) \right] \geq \frac{\delta}{3} \mathbb{I} \left\{ \alpha^2 \geq (\text{aec}_{2\varepsilon}^M(\mathcal{M}_0, \bar{M}))^{-1} \right\} - \frac{4}{3} \alpha^2 T.$$

We set $\alpha^2 = \frac{\delta}{8T}$, so that

$$\begin{aligned} \mathbb{P}^{M, \mathbb{A}} \left[\hat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M) \right] &\geq \frac{\delta}{6} \cdot \mathbb{I} \left\{ \alpha^2 \geq (\text{aec}_{2\varepsilon}^M(\mathcal{M}_0, \bar{M}))^{-1} \right\} \\ &= \frac{\delta}{6} \cdot \mathbb{I} \left\{ T \leq \frac{\delta}{8} \cdot \text{aec}_{2\varepsilon}^M(\mathcal{M}_0, \bar{M}) \right\}. \end{aligned}$$

By taking the supremum over all possible choices for $\bar{M} \in \mathcal{M}^+$, we conclude that unless

$$T > \frac{\delta}{8} \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_{2\varepsilon}^M(\mathcal{M}_0, \bar{M}),$$

the algorithm must have $\mathbb{P}^{M, \mathbb{A}} \left[\hat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M) \right] \geq \frac{\delta}{6}$. \square

H.3. Proof of Theorem B.2

Recall that for $M \in \mathcal{M}^+$ we define

$$\pi_M = \arg \max_{\pi \in \Pi} f^M(\pi)$$

as the set $\pi_M \subseteq \Pi$ of all optimal decisions π with $f^M(\pi) = \max_{\pi' \in \Pi} f^M(\pi')$. Unless otherwise stated, the results in this subsection do not make use of [Assumption A.4](#). For $\bar{M} \in \mathcal{M}^+$ and $\mathcal{M}_0 \subseteq \mathcal{M}$, we define

$$\mathcal{M}_0^{\text{opt}}(\bar{M}) = \{M \in \mathcal{M}_0 \mid \pi_M \subseteq \pi_{\bar{M}}, D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi)) = 0 \forall \pi \in \pi_{\bar{M}}\}.$$

For a subset $\Pi' \subseteq \Pi$, let

$$N_{\neg \Pi'} = |\{t \in [T] \mid \pi^t \notin \Pi'\}|.$$

Note that for all $M \in \mathcal{M}_0^{\text{opt}}(\bar{M})$, since $\pi_M \subseteq \pi_{\bar{M}}$, we have

$$N_{\neg \pi_{\bar{M}}} \leq N_{\neg \pi_M}.$$

Theorem B.2 (Main lower bound—strong variant). *Let $\varepsilon > 0$, $n_{\max} > 0$, and $\mathcal{M}_0 \subseteq \mathcal{M}$ be given, and define $\delta = \frac{\varepsilon}{2} \cdot \min\{1, \inf_{M \in \mathcal{M}_0} \mathfrak{g}^M / n_{\max}\}$. Unless*

$$\sup_{M \in \mathcal{M}_0} \frac{\mathfrak{g}^M}{\Delta_{\min}^M} \cdot \log(T) \geq \Omega(\delta^2) \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_{2\varepsilon}^M(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}),$$

there is no algorithm that simultaneously ensures that

$$1. \mathbb{E}^{M, \mathbb{A}}[\mathbf{Reg}(T)] \leq 2 \cdot \mathbf{g}^M \log(T), \quad \forall M \in \mathcal{M}_0.$$

$$2. \mathbb{P}^{M, \mathbb{A}}[\widehat{\lambda} \notin \Lambda(M; \varepsilon, \mathbf{n}_{\max})] \leq \frac{\delta}{12}, \quad \forall M \in \mathcal{M}_0.$$

Proof of Theorem B.2. For each $M \in \mathcal{M}$, if $\mathbb{E}^{M, \mathbb{A}}[\mathbf{Reg}(T)] \leq 2\mathbf{g}^M \log(T)$, then $\mathbb{E}^{M, \mathbb{A}}[N_{\neg\pi_M}] \leq 2 \frac{\mathbf{g}^M}{\Delta_{\min}^M} \log(T)$. The result now follows by appealing to [Theorem H.2](#) with $\bar{\mathbf{n}}^M = \mathbf{n}_{\max}$ and $R = 2 \sup_{M \in \mathcal{M}_0} \frac{\mathbf{g}^M}{\Delta_{\min}^M} \log(T)$. \square

Theorem H.2. Let $T \in \mathbb{N}$, $\varepsilon > 0$, $R \geq 1$, and $\mathcal{M}_0 \subseteq \mathcal{M}$ be given. Let $\{\bar{\mathbf{n}}^M\}_{M \in \mathcal{M}_0}$ be a collection of non-negative scalars indexed by \mathcal{M}_0 . Define $\delta = \frac{\varepsilon}{2} \cdot \min\left\{1, \inf_{M \in \mathcal{M}_0} \frac{\mathbf{g}^M}{\bar{\mathbf{n}}^M}\right\}$. Unless

$$R \geq \frac{\delta^2}{192} \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}),$$

there is no algorithm that simultaneously ensures that

$$1. \mathbb{E}^{M, \mathbb{A}}[N_{\neg\pi_M}] \leq R, \quad \forall M \in \mathcal{M}_0.$$

$$2. \mathbb{P}^{M, \mathbb{A}}[\widehat{\lambda} \notin \Lambda(M; \varepsilon, \bar{\mathbf{n}}^M)] \leq \frac{\delta}{12}, \quad \forall M \in \mathcal{M}_0.$$

Proof of Theorem H.2. Let $\varepsilon > 0$ be fixed. To prove the result, it suffices to lower bound the constrained minimax value

$$\mathfrak{M} := \sup_{\bar{M} \in \mathcal{M}^+} \inf_{\mathbb{A}} \left\{ \sup_{M \in \mathcal{M}_0^{\text{opt}}(\bar{M})} \mathbb{P}^{M, \mathbb{A}}[\widehat{\lambda} \notin \Lambda(M; \varepsilon, \bar{\mathbf{n}}^M)] \mid \mathbb{E}^{M', \mathbb{A}}[N_{\neg\pi_{M'}}] \leq R \quad \forall M' \in \mathcal{M}_0^{\text{opt}}(\bar{M}) \right\}. \quad (81)$$

We begin by appealing to the following technical lemma.

Lemma H.3. Let $\bar{M} \in \mathcal{M}^+$ and $T \in \mathbb{N}$ be given. Consider any algorithm \mathbb{A} with the property that for all $M \in \mathcal{M}_0^{\text{opt}}(\bar{M})$,

$$\mathbb{E}^{M, \mathbb{A}}[N_{\neg\pi_M}] \leq R$$

for some $R \geq 1$. For any $\beta \in (0, 1)$, there exists a modified algorithm \mathbb{A}' with the following properties:

- $\mathbb{P}^{M, \mathbb{A}'}[N_{\neg\pi_{\bar{M}}} > \lceil \frac{R}{\beta} \rceil] = 0$ for all models $M \in \mathcal{M}^+$.
- For all $M \in \mathcal{M}_0^{\text{opt}}(\bar{M})$,

$$\mathbb{P}^{M, \mathbb{A}}[\widehat{\lambda} \notin \Lambda(M; \varepsilon, \bar{\mathbf{n}}^M)] \geq \mathbb{P}^{M, \mathbb{A}'}[\widehat{\lambda} \notin \Lambda(M; \varepsilon, \bar{\mathbf{n}}^M)] - \beta.$$

By [Lemma H.3](#), for any $\beta \in (0, 1)$, we have

$$\mathfrak{M} \geq \sup_{\bar{M} \in \mathcal{M}^+} \inf_{\mathbb{A}} \left\{ \sup_{M \in \mathcal{M}_0^{\text{opt}}(\bar{M})} \mathbb{P}^{M, \mathbb{A}}[\widehat{\lambda} \notin \Lambda(M; \varepsilon, \bar{\mathbf{n}}^M)] \mid \mathbb{P}^{M', \mathbb{A}}[N_{\neg\pi_{\bar{M}}} > \lceil \frac{R}{\beta} \rceil] = 0 \quad \forall M' \in \mathcal{M}^+ \right\} - \beta.$$

Now, consider an arbitrary choice for \bar{M} above. We lower bound the minimax value using another technical lemma.

Lemma H.4. *Let $T \in \mathbb{N}$ and $\varepsilon > 0$ be given. Let $\{\bar{\pi}^M\}_{M \in \mathcal{M}}$ be a collection of non-negative scalars indexed by \mathcal{M} . Consider any algorithm \mathbb{A} with the property that*

$$\mathbb{P}^{\bar{M}, \mathbb{A}}[N_{\neg \pi_{\bar{M}}} > R] = 0$$

for some $R \geq 1$. For any $\bar{M} \in \mathcal{M}^+$, if we set $\delta := \frac{\varepsilon}{2} \cdot \min\left\{1, \inf_{M \in \mathcal{M}_0^{\text{opt}}(\bar{M})} \frac{g^M}{\bar{\pi}^M}\right\}$, then unless

$$R > \frac{\delta}{8} \cdot \text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}),$$

the algorithm must have

$$\sup_{M \in \mathcal{M}_0^{\text{opt}}(\bar{M})} \mathbb{P}^{M, \mathbb{A}}[\hat{\lambda} \notin \Lambda(M; \varepsilon, \bar{\pi}^M)] \geq \frac{\delta}{6}.$$

Bounding $\left[\frac{R}{\beta}\right] \leq \frac{2R}{\beta}$, it follows from Lemma H.4 that unless

$$\frac{2R}{\beta} > \frac{\delta}{8} \cdot \text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}),$$

where $\delta := \frac{\varepsilon}{2} \cdot \min\left\{1, \inf_{M \in \mathcal{M}_0^{\text{opt}}(\bar{M})} \frac{g^M}{\bar{\pi}^M}\right\}$, we have

$$\mathfrak{M} \geq \inf_{\mathbb{A}} \left\{ \sup_{M \in \mathcal{M}_0^{\text{opt}}(\bar{M})} \mathbb{P}^{M, \mathbb{A}}[\hat{\lambda} \notin \Lambda(M; \varepsilon, \bar{\pi}^M)] \mid \mathbb{P}^{M', \mathbb{A}}\left[N_{\neg \pi_{\bar{M}}} > \left\lceil \frac{R}{\beta} \right\rceil\right] = 0 \ \forall M' \in \mathcal{M}^+ \right\} - \beta \geq \frac{\delta}{6} - \beta.$$

To conclude, we set $\beta = \frac{\delta}{12}$ and maximize over $\bar{M} \in \mathcal{M}^+$. □

Proof of Lemma H.3. Fix $\beta \in (0, 1)$ and let $C := \lceil \frac{R}{\beta} \rceil$. Fix $\mathbb{A} = (p, q)$ and consider the algorithm $\mathbb{A}' = (p', q')$ defined implicitly as follows. For $t = 1, \dots, T$:

- Sample $\pi^t \sim p^t(\cdot \mid \mathcal{H}^{t-1})$.
- If $|\{i \leq t \mid \pi^i \notin \pi_{\bar{M}}\}| = C$, break and play an arbitrary decision $\pi \in \pi_{\bar{M}}$ until round T .

Return $\hat{\lambda} \sim q(\cdot \mid \mathcal{H}^T)$.

It is immediate from this construction that $\mathbb{A}' = (p', q')$ has $N_{\neg \pi_{\bar{M}}} \leq C$ almost surely under all possible models $M \in \mathcal{M}^+$. We now focus on bounding the performance. Let T_0 be the greatest value of t for which $|\{i \leq t \mid \pi^i \notin \pi_{\bar{M}}\}| \leq C$. First, observe that for all $M \in \mathcal{M}_0^{\text{opt}}(\bar{M})$, since the algorithms behave identically in law whenever $T_0 = T$,

$$\begin{aligned} \mathbb{P}^{M, \mathbb{A}'}[\hat{\lambda} \in \Lambda(M; \varepsilon, \bar{\pi}^M)] &\geq \mathbb{P}^{M, \mathbb{A}'}[\hat{\lambda} \in \Lambda(M; \varepsilon, \bar{\pi}^M) \wedge T_0 = T] \\ &= \mathbb{P}^{M, \mathbb{A}}[\hat{\lambda} \in \Lambda(M; \varepsilon, \bar{\pi}^M) \wedge T_0 = T] \\ &= \mathbb{P}^{M, \mathbb{A}}[\hat{\lambda} \in \Lambda(M; \varepsilon, \bar{\pi}^M) \wedge N_{\neg \pi_{\bar{M}}} \leq C]. \end{aligned}$$

By the union bound, we have

$$\mathbb{P}^{M, \mathbb{A}} \left[\widehat{\lambda} \in \Lambda(M; \varepsilon, \bar{n}^M) \wedge N_{\neg \pi_{\bar{M}}} \leq C \right] \geq \mathbb{P}^{M, \mathbb{A}} \left[\widehat{\lambda} \in \Lambda(M; \varepsilon, \bar{n}^M) \right] - \mathbb{P}^{M, \mathbb{A}} [N_{\neg \pi_{\bar{M}}} > C].$$

Finally, we observe that by Markov's inequality we have

$$\mathbb{P}^{M, \mathbb{A}} [N_{\neg \pi_{\bar{M}}} > C] \leq \frac{\mathbb{E}^{M, \mathbb{A}} [N_{\neg \pi_{\bar{M}}}]}{C} \leq \frac{\mathbb{E}^{M, \mathbb{A}} [N_{\neg \pi_M}]}{C} \leq \beta,$$

where we have used that $N_{\neg \pi_{\bar{M}}} \leq N_{\neg \pi_M}$, since $\pi_M \subseteq \pi_{\bar{M}}$. Rearranging, we obtain

$$\mathbb{P}^{M, \mathbb{A}} \left[\widehat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M) \right] \geq \mathbb{P}^{M, \mathbb{A}'} \left[\widehat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M) \right] - \beta.$$

□

Proof of Lemma H.4. Fix $\varepsilon > 0$, and let an algorithm \mathbb{A} be given. For any $\bar{M} \in \mathcal{M}^+$, define $q^{\bar{M}} = \mathbb{P}^{\bar{M}, \mathbb{A}}(\widehat{\lambda} = \cdot) \in \Delta(\Delta(\Pi))$, and let $\omega^{\bar{M}} := \mathbb{E}^{\bar{M}, \mathbb{A}} \left[\frac{1}{N_{\neg \pi_{\bar{M}}}} \sum_{t: \pi^t \notin \pi_{\bar{M}}} p^t \right] \in \Delta(\Pi)$, with the convention that the value inside the expectation is zero whenever $N_{\neg \pi_{\bar{M}}} = 0$.⁹

Fix $\alpha > 0$ and $\bar{M} \in \mathcal{M}^+$ be fixed. Define

$$M = \arg \max_{M \in \mathcal{M}_0^{\text{opt}}(\bar{M})} \left\{ \mathbb{P}_{\lambda \sim q^{\bar{M}}} [\lambda \notin \Lambda(M; \varepsilon, \bar{n}^M)] \mid \mathbb{E}_{\pi \sim \omega^{\bar{M}}} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\};$$

we assume that such an $M \in \mathcal{M}_0^{\text{opt}}(\bar{M})$ does exist, as otherwise the claim we will prove is trivial. It is immediate from this definition that we have

$$\begin{aligned} & \mathbb{P}^{\bar{M}, \mathbb{A}} \left[\widehat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M) \right] \\ &= \mathbb{P}_{\lambda \sim q^{\bar{M}}} [\lambda \notin \Lambda(M; \varepsilon, \bar{n}^M)] \\ &= \sup_{M \in \mathcal{M}_0^{\text{opt}}(\bar{M})} \left\{ \mathbb{P}_{\lambda \sim q^{\bar{M}}} [\lambda \notin \Lambda(M; \varepsilon, \bar{n}^M)] \mid \mathbb{E}_{\pi \sim \omega^{\bar{M}}} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\} \\ &\geq \inf_{q \in \Delta(\Delta(\Pi)), \omega \in \Delta(\Pi)} \sup_{M \in \mathcal{M}_0^{\text{opt}}(\bar{M})} \left\{ \mathbb{P}_{\lambda \sim q} [\lambda \notin \Lambda(M; \varepsilon, \bar{n}^M)] \mid \mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\} =: \text{opt}, \end{aligned}$$

with the convention that this value is zero if the set $\left\{ M \in \mathcal{M}_0^{\text{opt}}(\bar{M}) \mid \mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\}$ is empty. In addition, we have

$$\mathbb{E}_{\pi \sim \omega^{\bar{M}}} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2. \quad (82)$$

Now, define $\delta := \frac{\varepsilon}{2} \cdot \min \left\{ 1, \inf_{M \in \mathcal{M}_0^{\text{opt}}(\bar{M})} \frac{g^M}{\bar{n}^M} \right\}$, and let $\bar{\lambda}_q := \mathbb{E}_{\lambda \sim q}[\lambda]$. By Lemma H.1, we have

$$\text{opt} = \inf_{q \in \Delta(\Delta(\Pi)), \omega \in \Delta(\Pi)} \sup_{M \in \mathcal{M}_0^{\text{opt}}(\bar{M})} \left\{ \mathbb{P}_{\lambda \sim q} [\lambda \notin \Lambda(M; \varepsilon, \bar{n}^M)] \mid \mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\}$$

9. If $N_{\neg \pi_{\bar{M}}} = 0$ almost surely under \bar{M} , we can take $R = 0$, in which case the statement of the lemma is vacuous.

$$\begin{aligned}
 &\geq \delta \cdot \inf_{q \in \Delta(\Delta(\Pi)), \omega \in \Delta(\Pi)} \sup_{M \in \mathcal{M}_0^{\text{opt}}(\bar{M})} \left\{ \mathbb{I}\{\bar{\lambda}_q \notin \Lambda(M; 2\varepsilon, \bar{n}^M)\} \mid \mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\} \\
 &\geq \delta \cdot \inf_{\lambda \in \Delta(\Pi), \omega \in \Delta(\Pi)} \sup_{M \in \mathcal{M}_0^{\text{opt}}(\bar{M})} \left\{ \mathbb{I}\{\lambda \notin \Lambda(M; 2\varepsilon, \bar{n}^M)\} \mid \mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\} \\
 &\geq \delta \cdot \inf_{\lambda \in \Delta(\Pi), \omega \in \Delta(\Pi)} \sup_{M \in \mathcal{M}_0^{\text{opt}}(\bar{M})} \left\{ \mathbb{I}\{\lambda \notin \Lambda(M; 2\varepsilon)\} \mid \mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 \right\} \\
 &= \delta \cdot \mathbb{I}\left\{ \alpha^2 \geq \left(\text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}) \right)^{-1} \right\}. \tag{83}
 \end{aligned}$$

Hence, we have

$$\mathbb{P}^{\bar{M}, \mathbb{A}} \left[\hat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M) \right] \geq \delta \cdot \mathbb{I}\left\{ \alpha^2 \geq \left(\text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}) \right)^{-1} \right\}. \tag{84}$$

To proceed, using Lemma A.11 of Foster et al. (2021), we have

$$\begin{aligned}
 \mathbb{P}^{\bar{M}, \mathbb{A}} \left[\hat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M) \right] &\geq \frac{1}{3} \mathbb{P}^{\bar{M}, \mathbb{A}} \left[\hat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M) \right] - \frac{4}{3} D_{\text{KL}}(\mathbb{P}^{\bar{M}, \mathbb{A}} \parallel \mathbb{P}^{M, \mathbb{A}}) \\
 &\geq \frac{\delta}{3} \mathbb{I}\left\{ \alpha^2 \geq \left(\text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}) \right)^{-1} \right\} - \frac{4}{3} D_{\text{KL}}(\mathbb{P}^{\bar{M}, \mathbb{A}} \parallel \mathbb{P}^{M, \mathbb{A}}).
 \end{aligned}$$

Now, recall that from the definition, we have that for all $M \in \mathcal{M}_0^{\text{opt}}(\bar{M})$,

$$\begin{aligned}
 D_{\text{KL}}(\mathbb{P}^{\bar{M}, \mathbb{A}} \parallel \mathbb{P}^{M, \mathbb{A}}) &= \mathbb{E}^{M, \mathbb{A}} \left[\sum_{t: \pi^t \notin \pi_{\bar{M}}} \mathbb{E}_{\pi \sim p^t} D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi)) \right] \\
 &= \mathbb{E}^{M, \mathbb{A}} \left[\frac{N - \pi_{\bar{M}}}{N - \pi_{\bar{M}}} \sum_{t: \pi^t \notin \pi_{\bar{M}}} \mathbb{E}_{\pi \sim p^t} D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi)) \right] \\
 &\leq R \cdot \mathbb{E}^{M, \mathbb{A}} \left[\frac{1}{N - \pi_{\bar{M}}} \sum_{t: \pi^t \notin \pi_{\bar{M}}} \mathbb{E}_{\pi \sim p^t} D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi)) \right] \\
 &= R \cdot \mathbb{E}_{\pi \sim \omega_{\bar{M}}} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))] \leq \alpha^2 R,
 \end{aligned}$$

where the first inequality uses that $\mathbb{P}^{\bar{M}, \mathbb{A}}[N - \pi_{\bar{M}} > R] = 0$, and the second inequality uses (82).

With this, we have

$$\mathbb{P}^{\bar{M}, \mathbb{A}} \left[\hat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M) \right] \geq \frac{\delta}{3} \mathbb{I}\left\{ \alpha^2 \geq \left(\text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}) \right)^{-1} \right\} - \frac{4}{3} \alpha^2 R.$$

We set $\alpha^2 = \frac{\delta}{8R}$, so that

$$\begin{aligned}
 \mathbb{P}^{\bar{M}, \mathbb{A}} \left[\hat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M) \right] &\geq \frac{\delta}{6} \cdot \mathbb{I}\left\{ \alpha^2 \geq \left(\text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}) \right)^{-1} \right\} \\
 &= \frac{\delta}{6} \cdot \mathbb{I}\left\{ R \leq \frac{\delta}{8} \cdot \text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}) \right\}.
 \end{aligned}$$

We conclude that unless

$$R > \frac{\delta}{8} \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_{2\varepsilon}^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}),$$

the algorithm must have $\mathbb{P}^{M, \mathbb{A}}[\hat{\lambda} \notin \Lambda(M; \varepsilon, \bar{n}^M)] \geq \frac{\delta}{8}$.

□

H.4. Proofs for Lower Bound Examples

Proof of Example B.2. Let $\Delta \in (0, 1)$ and $A \geq 2$ be given and set

$$\mathcal{M} = \left\{ M(\pi) = \mathcal{N}(f^M(\pi), 1/2) \mid f^M \in [0, 1]^A \right\}$$

and

$$\mathcal{M}_0 = \{ M \in \mathcal{M} : \Delta_{\min}^M \geq \Delta/2 \}$$

Define $\bar{M} \in \mathcal{M}$ via $f^{\bar{M}}(\pi) = \Delta \mathbb{I}\{\pi = A\}$. Fix $\varepsilon \in (0, 1/2)$ and define a subclass

$$\mathcal{M}' = \{ \bar{M} \} \cup \{ M_i \}_{i \in [A-1]}$$

via

$$f^{M_i}(\pi) = \Delta \mathbb{I}\{\pi = A\} + \varepsilon \cdot \Delta \mathbb{I}\{\pi = i\}.$$

Since $\varepsilon \leq 1/2$, we have $\mathcal{M}' \subseteq \mathcal{M}^{\text{opt}}(\bar{M})$ and $\mathcal{M}' \subseteq \mathcal{M}_0$. In addition, we have

$$D_{\text{KL}}(\bar{M}(\pi) \parallel M_i(\pi)) = (f^{\bar{M}}(\pi) - f^{M_i}(\pi))^2.$$

Let $\mathcal{M}'' \subseteq \mathcal{M}$ denote the set of instances such that, for $M' \in \mathcal{M}''$, $D_{\text{KL}}(M_i(A) \parallel M'(A)) = 0$, and $M' \in \mathcal{M}^{\text{alt}}(M_i)$, for all $i \in [A-1]$. Then,

$$\begin{aligned} I^{M_i}(\lambda; \mathcal{M}) &= \inf_{M' \in \mathcal{M}^{\text{alt}}(M_i)} \mathbb{E}_{\pi \sim \lambda} [D_{\text{KL}}(M_i(\pi) \parallel M'(\pi))] \\ &\leq \inf_{M' \in \mathcal{M}''} \mathbb{E}_{\pi \sim \lambda} [D_{\text{KL}}(M_i(\pi) \parallel M'(\pi))] \\ &= \min_{j \in [A-1]} \left\{ \lambda_j \cdot (1 - \varepsilon)^2 \Delta^2 \mathbb{I}\{j = i\} + \lambda_j \cdot \Delta^2 \mathbb{I}\{j \neq i\} \right\}. \end{aligned} \tag{85}$$

We also have that

$$\mathbf{g}^{M_i} = \mathbf{g} := \frac{(A-2)}{\Delta} + \frac{1}{(1-\varepsilon)\Delta},$$

where we have used again that $\varepsilon \leq 1/2$.

Fix any pair $\lambda, \omega \in \Delta_{\Pi}$ and consider the value

$$\sup_{M \in \mathcal{M}' \setminus \mathcal{M}_{\varepsilon}^{\text{gl}}(\lambda)} \left\{ \frac{1}{\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \parallel M(\pi))]} \right\}.$$

Pick $\delta \leq \varepsilon/32$, and let $\mathcal{J} \subseteq [A-1]$ be the set of models i for which $\lambda \in \Lambda(M_i; \delta)$. For each such model, by definition, there exists $n_i > 0$ such that

$$\Delta^{M_i}(\lambda) \leq (1 + \delta) \frac{\mathbf{g}}{n_i}, \quad \text{and} \quad I^{M_i}(\lambda; \mathcal{M}) \geq (1 - \delta) \frac{1}{n_i}.$$

Define $\bar{n} = \max_{i \in \mathcal{J}} n_i$ and $\underline{n} = \min_{i \in \mathcal{J}} n_i$. Using Eq. (85), it is immediate that for all $j \in [A-1]$, $\lambda_j \geq (1 - \delta) \frac{\Delta^{-2}}{\underline{n}}$, and that for all $i \in \mathcal{J}$,

$$\lambda_i \geq (1 - \delta) \frac{(1 - \varepsilon)^{-2} \Delta^{-2}}{n_i}.$$

In particular, this implies that for all $i \in \mathcal{J}$,

$$\begin{aligned} (1 - \delta) \frac{(A-2)\Delta^{-1}}{\underline{n}} + (1 - \delta) \frac{(1 - \varepsilon)^{-1} \Delta^{-1}}{n_i} &\leq \Delta^{M_i}(\lambda) \\ &\leq (1 + \delta) \frac{\mathbf{g}}{n_i} \\ &= (1 + \delta) \frac{(A-2)\Delta^{-1}}{n_i} + (1 + \delta) \frac{(1 - \varepsilon)^{-1} \Delta^{-1}}{n_i}, \end{aligned}$$

or by rearranging,

$$\begin{aligned} (1 - \delta) \frac{(A-2)\Delta^{-1}}{\underline{n}} &\leq (1 + \delta) \frac{(A-2)\Delta^{-1}}{n_i} + 2\delta \frac{(1 - \varepsilon)^{-1} \Delta^{-1}}{n_i}, \\ &\leq (1 + \delta) \frac{(A-2)\Delta^{-1}}{n_i} + 4\delta \frac{\Delta^{-1}}{n_i}, \\ &\leq (1 + 2\delta) \frac{(A-2)\Delta^{-1}}{n_i} \end{aligned}$$

as long as $A \geq 6$. Since this holds uniformly for all $i \in \mathcal{J}$, rearranging once more gives

$$\bar{n} \leq \frac{(1 + 2\delta)}{1 - \delta} \underline{n} \leq (1 + 2\delta)^2 \underline{n} \leq (1 + 8\delta) \underline{n},$$

where we have used that $\delta \leq 1/2$.

Now, observe that for all $i \in \mathcal{J}$, we have

$$\begin{aligned} \Delta^{M_i}(\lambda) &\geq (1 - \delta)(A - |\mathcal{J}| - 1) \frac{1}{\Delta \underline{n}} + (1 - \delta) |\mathcal{J}| \frac{1}{(1 - \varepsilon)^2 \Delta \bar{n}} + (1 - \delta) \frac{1}{(1 - \varepsilon) \Delta n_i}, \\ &\geq (1 - \delta)(A - |\mathcal{J}| - 1) \frac{1}{\Delta \underline{n}} + \frac{(1 - \delta)}{1 + 8\delta} |\mathcal{J}| \frac{1}{(1 - \varepsilon)^2 \Delta \underline{n}} + (1 - \delta) \frac{1}{(1 - \varepsilon) \Delta n_i}, \\ &\geq \frac{(1 - \delta)}{\Delta \underline{n}} \left((A - |\mathcal{J}| - 1) + |\mathcal{J}| \frac{1 - 8\delta}{(1 - \varepsilon)^2} \right) + (1 - \delta) \frac{1}{(1 - \varepsilon) \Delta n_i}, \\ &\geq \frac{(1 - \delta)}{\Delta \underline{n}} \left((A - |\mathcal{J}| - 1) + |\mathcal{J}| \frac{1}{(1 - \varepsilon)} \right) + (1 - \delta) \frac{1}{(1 - \varepsilon) \Delta n_i}, \end{aligned}$$

where we have used that $\frac{1}{1+x} \geq 1 - x$, and that $\delta \leq \varepsilon/8$. Suppose that $|\mathcal{J}| \geq \frac{A}{2}$. Then we have

$$(A - |\mathcal{J}| - 1) + |\mathcal{J}| \frac{1}{(1 - \varepsilon)} \geq \frac{A}{2} \left(1 + \frac{1}{1 - \varepsilon} \right) - 1 \geq (1 + \varepsilon/2)A - 1$$

$$\geq (1 + \varepsilon/2)(A - 1),$$

so that

$$\Delta^{M_i}(\lambda) \geq \frac{(1 - \delta)(1 + \varepsilon/2)}{\Delta \underline{n}} (A - 1) + (1 - \delta) \frac{1}{(1 - \varepsilon)\Delta n_i}.$$

Noting that $\underline{n} \leq n_i$ and $\delta \leq \varepsilon/8$, we further have

$$\Delta^{M_i}(\lambda) \geq (1 + \varepsilon/4) \frac{A - 1}{\Delta} \frac{1}{n_i} + (1 - \delta) \frac{1}{(1 - \varepsilon)\Delta n_i}.$$

Observe that the right-hand side above is greater than $(1 + \delta) \frac{\underline{g}}{n_i}$ if and only if

$$(1 + \varepsilon/4)(A - 1) > (1 + \delta)(A - 1) + \frac{2\delta}{1 - \varepsilon},$$

which is satisfied if $\delta \leq \varepsilon/32$. In this case, we have

$$\Delta^{M_i}(\lambda) > (1 + \delta) \frac{\underline{g}}{n_i},$$

which contradicts the assumption that $i \in \mathcal{J}$. It follows that we must have $|\mathcal{J}| < A/2$.

Now, to conclude, select $i = \arg \min_{i \in [A-1] \setminus \mathcal{J}} \omega_i$, and consider the value

$$\sup_{M \in \mathcal{M}' \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)} \left\{ \frac{1}{\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \| M(\pi))]} \right\} \geq \frac{1}{\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \| M_i(\pi))]} = \frac{1}{\omega_i \cdot \varepsilon^2 \Delta^2},$$

where the first inequality follows because $\lambda \notin \Lambda(M_i; \delta)$ by definition, and the equality follows from the construction of \bar{M} and M_i . Since $\sum_{i \in [A-1] \setminus \mathcal{J}} \omega_i \leq 1$ and $|[A-1] \setminus \mathcal{J}| \geq \frac{A}{2}$, we must have $\omega_i \leq \frac{2}{A}$, so that

$$\frac{1}{\omega_i \cdot \varepsilon^2 \Delta^2} \geq \frac{A}{2\varepsilon^2 \Delta^2}$$

as desired. To complete the proof, note that this holds uniformly for all choices for λ and ω , and that

$$\begin{aligned} \text{aec}_\varepsilon^{\mathcal{M}}(\mathcal{M}_0, \bar{M}) &= \inf_{\lambda, \omega \in \Delta_\Pi} \sup_{M \in \mathcal{M}_0 \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)} \left\{ \frac{1}{\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \| M(\pi))]} \right\} \\ &\geq \inf_{\lambda, \omega \in \Delta_\Pi} \sup_{M \in \mathcal{M}' \setminus \mathcal{M}_\varepsilon^{\text{gl}}(\lambda)} \left\{ \frac{1}{\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \| M(\pi))]} \right\}. \end{aligned}$$

To obtain the parameter setting in the theorem statement, we rescale $\Delta \leftarrow 2\Delta$ and $\varepsilon \leftarrow 32\varepsilon$. \square

Proof of Example B.3. We reduce the lower bound to that of multi-armed bandits via a standard tree construction (Osband and Van Roy, 2016; Domingues et al., 2021); as the argument is standard, we only sketch the approach. Assume without loss of generality that H is a multiple of 2. Set $H = \log_2(S/2)$. Consider a sub-class $\mathcal{M}' \subseteq \mathcal{M}$ defined as follows. All models $M \in \mathcal{M}'$ have

identical, deterministic dynamics given by a binary tree. Each layer h has 2^{h-1} states, so that layer H has $S/4$ states, and the total number of states is $S - 1$. The agent begins from a root state s_1 deterministically. For $h \leq H - 1$, there are two available actions, left and right. Choosing left leads to the left successor for the current layer, and right leads to the right successor for the current layer. There are no rewards for layer $h \leq H - 1$. For layer H , there are A available actions, and rewards are arbitrary, subject to the constraint that the mean lies in $[0, 1]$ and the noise follows $\mathcal{N}(0, 1/2)$.

It is clear that the class \mathcal{M}' is equivalent to the class of multi-armed bandit instances with $SA/4$ actions. As a consequence, the lower bound follows from [Example B.2](#). \square

Proof of [Example B.1](#). Let $\Delta \in (0, 1/6)$, $\beta \in (0, 1)$, and $A, N \geq 2$ be given. Consider the reference model $\bar{M} \in \mathcal{M}^+$ defined as follows:

- For each bandit arm $k \in [A]$, we have $f^{\bar{M}}(k) = \frac{1}{2} + \Delta \mathbb{I}\{k = A\}$ and $r \sim \mathcal{N}(f^{\bar{M}}(k), 1)$. There are no observations, i.e. $o = \perp$ almost surely.
- For each revealing arm π_k° , we receive zero reward almost surely (so $f^{\bar{M}}(\pi_k^\circ) = 0$) and $o \sim \text{Unif}([A])$.

We define a subclass

$$\mathcal{M}' = \{M_j\}_{j \in [N]} \subset \mathcal{M}_0^{\text{opt}}(\bar{M})$$

as follows

- For each bandit arm $k \in [A]$, we have $f^{M_j}(k) = \frac{1}{2} + \Delta \mathbb{I}\{k = A\}$ and $r \sim \mathcal{N}(f^{M_j}(k), 1)$. There are no observations, i.e. $o = \perp$ almost surely.
- For each revealing arm π_k° , we receive zero reward almost surely (so $f^{M_j}(\pi_k^\circ) = 0$). We have

$$o \sim \begin{cases} \text{Unif}([A]), & k \neq j, \\ \beta \mathbb{I}_i + (1 - \beta) \text{Unif}([A]), & k = j. \end{cases}$$

Note that $\mathcal{M}' \subseteq \mathcal{M}_0$. For all $j \in [N]$, a direct calculation gives

$$D_{\text{KL}}(\bar{M}(\pi_j^\circ) \| M_j(\pi_j^\circ)) = \frac{A-1}{A} \log\left(\frac{1}{1-\beta}\right) + \frac{1}{A} \log\left(\frac{1}{1+\beta(A-1)}\right) =: \alpha$$

and

$$D_{\text{KL}}(\bar{M}(\pi) \| M_j(\pi)) = \alpha \cdot \mathbb{I}\{\pi = \pi_j^\circ\}. \quad (86)$$

In addition, it is straightforward to see that $\alpha \leq 2\beta$ whenever $\beta \leq 1/2$. Next we calculate that for any $j \in [N]$ and $M \in \mathcal{M}$ with $\pi_M^\circ = \pi_j^\circ$, and $\pi_M \neq A$,

$$D_{\text{KL}}(M_j(\pi_j^\circ) \| M(\pi_j^\circ)) = \beta \log\left(1 + \frac{\beta A}{1-\beta}\right) =: \gamma,$$

which has $\gamma \leq O(\beta)$ whenever $\beta \leq 1/A$ and $\gamma \geq \beta \log(1 + \beta A)$. Lastly, we have that for all $i \in [A]$, all $M, M' \in \mathcal{M}$ have

$$D_{\text{KL}}(M(i) \| M'(i)) = \frac{1}{2}(f^M(i) - f^{M'}(i))^2.$$

Let \mathcal{M}'' denote the set of instances such that for $M' \in \mathcal{M}''$, $\pi_{M'} \neq A$, and $f^{M'}(A) = \frac{1}{2} + \Delta$, so that $D_{\text{KL}}(M_j(A) \| M'(A)) = 0$ and $\mathcal{M}'' \subseteq \mathcal{M}^{\text{alt}}(M_j)$ for all $j \in [N]$. Using the above calculations and the definition of $I^{M_j}(\lambda; \mathcal{M})$, we can then compute, for all $j \in [N]$,

$$I^{M_j}(\lambda; \mathcal{M}) \leq \inf_{M' \in \mathcal{M}''} \mathbb{E}_\lambda[D_{\text{KL}}(M_j(\pi) \| M'(\pi))] = \frac{\Delta^2}{2} \cdot \min_{k \in [A-1]} \lambda_k + \gamma \cdot \lambda_{\pi_j^\circ} \quad (87)$$

and

$$g^{M_j} = g := \min \left\{ 2 \frac{A-1}{\Delta}, \left(\frac{1}{2} + \Delta \right) \frac{1}{\gamma} \right\} = \left(\frac{1}{2} + \Delta \right) \frac{1}{\gamma}$$

whenever $\gamma \geq \Delta/2(A-1)$.

Fix any pair $\lambda, \omega \in \Delta_\Pi$ and consider the value

$$\sup_{M \in \mathcal{M}' \setminus \mathcal{M}_{1/2}^{\text{gl}}(\lambda)} \left\{ \frac{1}{\mathbb{E}_{\pi \sim \omega}[D_{\text{KL}}(\bar{M}(\pi) \| M(\pi))]} \right\}.$$

Let $\mathcal{J} \subseteq [N]$ be the set of models j for which $\lambda \in \Lambda(M_j; 1/2)$. For each such model, by definition, there exists $n_j > 0$ such that

$$\Delta^{M_j}(\lambda) \leq (1 + 1/2) \frac{g}{n_j} \leq \frac{1}{\gamma n_j}, \quad \text{and} \quad I^{M_j}(\lambda; \mathcal{M}) \geq \frac{1}{2n_j}.$$

Define $\bar{n} = \max_{j \in \mathcal{J}} n_j$ and $\underline{n} = \min_{j \in \mathcal{J}} n_j$. Let us begin with some basic observations.

- Since $\Delta^{M_j} = \Delta^{\bar{M}}$ for all j , we have $\Delta^{M_j}(\lambda) \leq \frac{1}{\gamma n_{j'}}$ for all $j, j' \in \mathcal{J}$, and hence

$$\Delta^{M_j}(\lambda) \leq \frac{1}{\gamma \bar{n}}. \quad (88)$$

- Any $j \in \mathcal{J}$ must have

$$\lambda_{\pi_j^\circ} \gamma \geq \frac{1}{4n_j} \geq \frac{1}{4\bar{n}}. \quad (89)$$

To see this, observe that if it were not the case, we would need $\min_{i \in [A-1]} \lambda_i \frac{\Delta^2}{2} \geq \frac{1}{4n_j}$ to satisfy the constraint that $I^{M_j}(\lambda; \mathcal{M}) \geq \frac{1}{2n_j}$ (by (87)). But if this were to occur, we would have

$$\frac{A-1}{2\Delta n_j} \leq \Delta^{M_j}(\lambda) \leq \frac{1}{\gamma n_j},$$

which would contradict the assumption that $\gamma \geq 2\Delta/(A-1)$.

Combing the inequalities (88) and (89), it follows that any $j \in \mathcal{J}$ must have

$$\frac{|\mathcal{J}|}{8\gamma n_j} \leq \frac{1}{2} \sum_{k \in \mathcal{J}} \lambda_{\pi_k^\circ} \leq \Delta^{M_j}(\lambda) \leq \frac{1}{\gamma n_j},$$

which implies that $|\mathcal{J}| \leq 8$. Hence, as long as $N \geq 16$, we have $|[N] \setminus \mathcal{J}| \geq N/2$.

To conclude, select $k = \arg \min_{j \in [N] \setminus \mathcal{J}} \omega_{\pi_j^\circ}$, and consider the value

$$\sup_{M \in \mathcal{M}' \setminus \mathcal{M}_{1/2}^{\text{gl}}(\lambda)} \left\{ \frac{1}{\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \| M(\pi))]} \right\} \geq \frac{1}{\mathbb{E}_{\pi \sim \omega} [D_{\text{KL}}(\bar{M}(\pi) \| M_k(\pi))]} = \frac{1}{\omega_{\pi_k^\circ} \cdot \alpha},$$

where the first inequality follows because $\lambda \notin \Lambda(M_k; 1/2)$ by definition, and the equality follows from (86). Since $\sum_{j \in [N] \setminus \mathcal{J}} \omega_{\pi_j^\circ} \leq 1$ and $|[N] \setminus \mathcal{J}| \geq \frac{N}{2}$, we must have $\omega_{\pi_k^\circ} \leq \frac{2}{N}$, so that

$$\frac{1}{\omega_{\pi_k^\circ} \cdot \alpha} \geq \frac{N}{2\alpha}.$$

as desired. Since this holds uniformly for all choices for λ and ω , the proof is completed. \square

H.5. Lower Bound on Regret for Algorithms with Well-Behaved Tails

In this section, we present an additional result, [Theorem H.3](#), which shows that for algorithms for which the tail behavior is “well-behaved” in a certain sense, the Allocation-Estimation Coefficient directly leads to lower bounds on the least possible value of T for which any algorithm can achieve (approximate) instance-optimality.

Theorem H.3. *Let the time horizon $T \in \mathbb{N}$, $\varepsilon \in (0, 1/2)$, and $\mathcal{M}_0 \subseteq \mathcal{M}$ be given. Suppose that there exists an algorithm \mathbb{A} with the property that for all $M \in \mathcal{M}_0$,*

1. $\mathbb{E}^{M, \mathbb{A}}[\mathbf{Reg}(T)] \leq (1 + \varepsilon)g^M \log(T)$.
2. For all $\pi \in \Pi$, if $\mathbb{E}^{M, \mathbb{A}}[T(\pi)] \neq 0$, then $\mathbb{E}^{M, \mathbb{A}}[T(\pi)] \geq 1$.
3. $\sqrt{\mathbb{E}^{M, \mathbb{A}}[(\mathbf{Reg}(T))^2]} \leq 2g^M \log(T)$.

In addition, assume that 1) $g^M \geq 1$ for all $M \in \mathcal{M}_0$, 2) [Assumption A.4](#) holds, 3) [Assumption A.2](#) holds with parameter $V_{\mathcal{M}} \geq 1$, and 4) that

$$\log(T) \geq \frac{12}{\varepsilon} \log \left(\sup_{M \in \mathcal{M}} \frac{2g^M}{\Delta_{\min}^M} \cdot \log(T) \right).$$

Then if we define $\delta = \varepsilon \cdot \min\{1, \inf_{M \in \mathcal{M}_0} \frac{g^M}{3g^M/\Delta_{\min}^M + n_\varepsilon^M}\}$, it must be the case that

$$\log^3(T) \geq \frac{\delta^2}{C} \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_{4\varepsilon}^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}),$$

for $C \leq O\left(\left(\sup_{M \in \mathcal{M}} \frac{g^M}{\Delta_{\min}^M}\right)^4 \cdot \frac{V_{\mathcal{M}}^2 \log(\delta^{-1})}{\varepsilon^2}\right)$.

We prove [Theorem H.3](#) by combining [Theorem B.2](#) with another technical result, [Proposition H.1](#) (stated and proven in the sequel), which shows that for algorithms that satisfy the assumptions of [Theorem H.3](#), it is possible to use the empirical decision frequencies to compute an allocation that is approximately optimal with high probability.

Proof of [Theorem H.3](#). Define $\bar{n}^M = 3 \frac{\mathbf{g}^M}{\Delta_{\min}^M} + n_\varepsilon^M$, and set

$$\delta := \varepsilon \cdot \min \left\{ 1, \inf_{M \in \mathcal{M}_0} \frac{\mathbf{g}^M}{\bar{n}^M} \right\}.$$

Assume that

$$\log(T) \geq \frac{12}{\varepsilon} \log \left(\sup_{M \in \mathcal{M}} \frac{2\mathbf{g}^M}{\Delta_{\min}^M} \cdot \log(T) \right).$$

Let \mathbb{A} be the algorithm in the statement of the proposition, and let \mathbb{A}' be the modified algorithm created through [Proposition H.1](#) with parameter δ . By assumption, we have that $\sqrt{\mathbb{E}^{M, \mathbb{A}}[N_{\neg\pi_M}]} \leq R := 2 \sup_{M \in \mathcal{M}} \frac{\mathbf{g}^M}{\Delta_{\min}^M} \log(T)$. We define $n = c \cdot \frac{\log(24\delta^{-1})}{\varepsilon^2} \cdot \frac{R^3 V_M^2}{\log(T)}$ for a sufficiently large numerical constant c and $T' = T \cdot n$. [Proposition H.1](#) implies that for time T' , the algorithm \mathbb{A}' satisfies

$$\sqrt{\mathbb{E}^{M, \mathbb{A}'}[N_{\neg\pi_M}]} \leq R' := R \cdot n$$

and

$$\mathbb{P}^{M, \mathbb{A}'} \left[\hat{\lambda} \in \Lambda(M; 2\varepsilon, \bar{n}^M) \right] \geq 1 - \frac{\delta}{24}.$$

On the other hand, since the precondition of [Theorem H.2](#) is now satisfied with parameter R' , we have that unless

$$R' \geq \frac{\delta^2}{192} \cdot \sup_{\bar{M} \in \mathcal{M}^+} \text{aec}_{4\varepsilon}^{\mathcal{M}}(\mathcal{M}_0^{\text{opt}}(\bar{M}), \bar{M}), \quad (90)$$

the algorithm must have

$$\mathbb{P}^{M, \mathbb{A}'} \left[\hat{\lambda} \in \Lambda(M; 2\varepsilon, \bar{n}^M) \right] \leq 1 - \frac{\delta}{12}.$$

As $\frac{\delta}{12} > \frac{\delta}{24}$, this is a contradiction unless (90) holds. □

Proposition H.1. *Let the time horizon $T \in \mathbb{N}$ and $\mathcal{M}_0 \subseteq \mathcal{M}$ be given. Let \mathbb{A} be an algorithm with the property that for all $M \in \mathcal{M}_0$,*

1. $\mathbb{E}^{M, \mathbb{A}}[\mathbf{Reg}(T)] \leq (1 + \varepsilon)\mathbf{g}^M \log(T)$ for some $\varepsilon \in (0, 1)$.
2. For all $\pi \in \Pi$, if $\mathbb{E}^{M, \mathbb{A}}[T(\pi)] \neq 0$, then $\mathbb{E}^{M, \mathbb{A}}[T(\pi)] \geq 1$.
3. $\sqrt{\mathbb{E}^{M, \mathbb{A}}[N_{\neg\pi_M}^2]} \leq R$ for some $R \geq 2$.

In addition, assume that 1) $\mathbf{g}^M \geq 1$ for all $M \in \mathcal{M}_0$, 2) [Assumption A.4](#) holds, 3) [Assumption A.2](#) holds with parameter $V_{\mathcal{M}} \geq 1$, and 4) that

$$\log(T) \geq \frac{12}{\varepsilon} \log \left(\sup_{M \in \mathcal{M}} \frac{2\mathbf{g}^M}{\Delta_{\min}^M} \cdot \log(T) \right).$$

Then for any $\delta \in (0, e^{-1})$, if we define $n = c \cdot \frac{\log(\delta^{-1})}{\varepsilon^2} \cdot \frac{R^3 V_{\mathcal{M}}^2}{\log(T)}$ for a sufficiently large numerical constant c , there exists an algorithm \mathbb{A}' that, using $T' := T \cdot n$ rounds, returns a normalized allocation $\hat{\lambda} \in \Delta_{\Pi}$ such that

$$\mathbb{P}^{M, \mathbb{A}'} \left[\hat{\lambda} \in \Lambda(M; 2\varepsilon, \bar{n}^M) \right] \geq 1 - \delta,$$

for $\bar{n}^M \leq 3 \frac{\mathbf{g}^M}{\Delta_{\min}^M} + n_{\varepsilon}^M$ and that $\sqrt{\mathbb{E}^{M, \mathbb{A}'} [N_{\neg \pi_M}^2]} \leq R \cdot n$ and

$$\mathbb{E}^{M, \mathbb{A}'} [\mathbf{Reg}(T')] \leq (1 + \varepsilon) \mathbf{g}^M \log(T) \cdot n$$

for all $M \in \mathcal{M}_0$.

Proof of [Proposition H.1](#). We first state a technical lemma regarding robust mean estimation.

Lemma H.5. Let $X \in \mathbb{R}^d$ be a random variable with $\mu := \mathbb{E}[X]$. Assume that $\|\mu\|_0 \leq s$, where s is a known upper bound. For any $\delta \leq e^{-1}$, there exists an estimator $\hat{\mu}_n$ that, given n independent samples from X , ensures that with probability at least $1 - \delta$,

$$\|\hat{\mu}_n - \mu\|_1 \leq 24 \sqrt{\frac{2s \cdot \mathbb{E}\|X - \mu\|_2^2 \cdot \log(\delta^{-1})}{n}}.$$

In addition, $\|\hat{\mu}_n\|_0 \leq s$ with probability 1.

Throughout, we will use that since [Assumption A.4](#) holds, π_M is unique for all $M \in \mathcal{M}$. Fix $M \in \mathcal{M}_0$. Let $\hat{\eta} \in \mathbb{R}_{+}^{\Pi}$ denote the vector of empirical decision frequencies when \mathbb{A} is run with horizon T , i.e. $\hat{\eta}(\pi) = T(\pi)$. Let

$$\eta^M = \mathbb{E}^{M, \mathbb{A}}[\hat{\eta}].$$

For parameters $n \in \mathbb{N}$ and $\delta \leq e^{-1}$ we define \mathbb{A}' as follows:

- Run \mathbb{A} a total of n times independently (so that $T' = T \cdot n$), and let $\hat{\eta}^1, \dots, \hat{\eta}^n$ be the empirical decision frequencies.
- Apply the algorithm from [Lemma H.5](#) to $\hat{\eta}^1, \dots, \hat{\eta}^n$ with parameters δ and $s = 2R$, and let $\check{\eta} \in \mathbb{R}_{+}^{\Pi}$ be the resulting vector (note that we can take $\check{\eta}$ to have non-negative entries without loss of generality).
- Set $\hat{\pi} = \arg \max_{\pi \in \Pi} \check{\eta}$, and set $\tilde{\eta}(\pi) = \check{\eta}(\pi) \mathbb{I}\{\pi \neq \hat{\pi}\}$ and $\tilde{\eta}(\hat{\pi}) = n_{\varepsilon}^M \cdot \log T$ (note that n_{ε}^M is a class-dependent quantity, and so is known to the learner).
- Set $\hat{\lambda} = \tilde{\eta} / \|\tilde{\eta}\|_1$.

It is immediate that this construction satisfies $\sqrt{\mathbb{E}^{M, \mathbb{A}'} [N_{\neg\pi_M}^2]} \leq R \cdot n$ and

$$\mathbb{E}^{M, \mathbb{A}'} [\mathbf{Reg}(T')] \leq (1 + \varepsilon) \mathbf{g}^M \log(T) \cdot n,$$

so it remains to show that $\hat{\lambda}$ is a near-optimal allocation with high probability when n is chosen appropriately.

We start by applying [Lemma H.5](#). To do so, we carry out some prerequisite calculations. First, by assumption, we have $\eta^M(\pi) \geq 1$ if $\eta^M \neq 0$. Using this, along with [Assumption A.4](#), we have

$$\|\eta^M\|_0 \leq 1 + \sum_{\pi \neq \pi_M} \eta^M(\pi) \leq 1 + \mathbb{E}^{M, \mathbb{A}} [N_{\neg\pi_M}] \leq 1 + R \leq 2R.$$

Second,

$$\begin{aligned} \mathbb{E}^{M, \mathbb{A}} \|\hat{\eta} - \eta^M\|_2^2 &\leq \mathbb{E}^{M, \mathbb{A}} \left[(\hat{\eta}(\pi_M) - \eta^M(\pi_M))^2 \right] + \sum_{\pi \neq \pi_M} \mathbb{E}^{M, \mathbb{A}} [\hat{\eta}(\pi)^2] \\ &\leq \mathbb{E}^{M, \mathbb{A}} \left[(\hat{\eta}(\pi_M) - \eta^M(\pi_M))^2 \right] + \mathbb{E}^{M, \mathbb{A}} [N_{\neg\pi_M}^2]. \end{aligned}$$

Furthermore,

$$\begin{aligned} \mathbb{E}^{M, \mathbb{A}} \left[(\hat{\eta}(\pi_M) - \eta^M(\pi_M))^2 \right] &= \mathbb{E}^{M, \mathbb{A}} \left[\left(\sum_{\pi \neq \pi_M} \hat{\eta}(\pi) - \sum_{\pi \neq \pi_M} \eta^M(\pi) \right)^2 \right] \\ &\leq \mathbb{E}^{M, \mathbb{A}} \left[\left(\sum_{\pi \neq \pi_M} \hat{\eta}(\pi) \right)^2 \right] \\ &= \mathbb{E}^{M, \mathbb{A}} [N_{\neg\pi_M}^2]. \end{aligned}$$

so that

$$\mathbb{E}^{M, \mathbb{A}} \|\hat{\eta} - \eta^M\|_2^2 \leq 2 \mathbb{E}^{M, \mathbb{A}} [N_{\neg\pi_M}^2] \leq 2R^2.$$

As a result, [Lemma H.5](#) implies that with probability $1 - \delta$,

$$\|\check{\eta} - \eta^M\|_1 \leq \sqrt{C_1 \frac{\log(\delta^{-1})}{n}} =: \varepsilon_{\text{stat}} \quad (91)$$

where $C_1 = O(R^3)$.

Next, we appeal to [Lemma H.2](#), which implies that as long as

$$\log(T) \geq \frac{12}{\varepsilon} \log \left(\sup_{M \in \mathcal{M}} \frac{2\mathbf{g}^M}{\Delta_{\min}^M} \cdot \log(T) \right), \quad (92)$$

we have

$$\Delta^M(\eta^M) \leq (1 + \varepsilon/2) \mathbf{g}^M \log(T), \quad \text{and} \quad I^M(\eta^M; \mathcal{M}) \geq (1 - \varepsilon/2) \log(T).$$

Applying (91) and using that $\Delta^M \in [0, 1]$ and $D_{\text{KL}}(M(\pi) \| M'(\pi)) \leq 2V_{\mathcal{M}}$ by Lemma F.13, we have

$$\Delta^M(\check{\eta}) \leq (1 + \varepsilon/2)\mathbf{g}^M \log(T) + \varepsilon_{\text{stat}}, \quad \text{and} \quad I^M(\check{\eta}; \mathcal{M}) \geq (1 - \varepsilon/2) \log(T) - \varepsilon_{\text{stat}} \cdot 2V_{\mathcal{M}}.$$

Thus, as long as

$$\varepsilon_{\text{stat}} \leq \frac{1}{2}\varepsilon \cdot \min\{\mathbf{g}^M \log(T), (2V_{\mathcal{M}})^{-1} \log(T)\},$$

we have

$$\Delta^M(\check{\eta}) \leq (1 + \varepsilon)\mathbf{g}^M \log(T), \quad \text{and} \quad I^M(\check{\eta}; \mathcal{M}) \geq (1 - \varepsilon) \log(T). \quad (93)$$

Next, we claim that $\hat{\pi} = \pi_M$. To see this, note that since $\mathbb{E}^{M, \mathbb{A}}[\mathbf{R}\text{eg}(T)] \leq 2\mathbf{g}^M \log(T)$, we have

$$\eta^M(\pi_M) \geq T - 2 \frac{\mathbf{g}^M \log(T)}{\Delta_{\min}^M} \geq \frac{3}{4}T$$

as long as $T > 8 \frac{\mathbf{g}^M}{\Delta_{\min}^M} \log(T)$, which is implied by the condition (92). Hence, as long as $\varepsilon_{\text{stat}} \leq \frac{T}{4}$, we have

$$\check{\eta}(\pi_M) > \frac{T}{2},$$

which implies that $\hat{\pi} = \pi_M$. By definition of n_ε^M , and since $\check{\eta}$ satisfies Eq. (93) above, we then have that setting $\check{\eta}(\hat{\pi}) = n_\varepsilon^M \log(T)$ does not affect the regret, and only decreases the information gain by a factor of $\varepsilon \log(T)$. It follows that

$$\Delta^M(\tilde{\eta}) \leq (1 + \varepsilon)\mathbf{g}^M \log(T), \quad \text{and} \quad I^M(\tilde{\eta}; \mathcal{M}) \geq (1 - 2\varepsilon) \log(T).$$

We conclude that $\hat{\lambda} \in \Lambda(M; 2\varepsilon, \bar{n})$ for

$$\bar{n} = \|\tilde{\eta}\|_1 / \log(T).$$

To wrap up, we compute that

$$\|\tilde{\eta}\|_1 \leq \sum_{\pi \neq \pi_M} \eta^M(\pi) + \varepsilon_{\text{stat}} + n_\varepsilon^M \log(T) \leq 2 \frac{\mathbf{g}^M}{\Delta_{\min}^M} \log(T) + \varepsilon_{\text{stat}} + n_\varepsilon^M \log(T) \leq 3 \frac{\mathbf{g}^M}{\Delta_{\min}^M} \log(T) + n_\varepsilon^M \log(T)$$

whenever $\varepsilon_{\text{stat}} \leq \mathbf{g}^M \log(T)$. □

Proof of Lemma H.5. From Proposition 1 of Lugosi and Mendelson (2019), we have that for any $\delta \leq e^{-1}$, there exists an estimator $\tilde{\mu}_n$ that, given n independent samples of X , ensures that with probability at least $1 - \delta$,

$$\|\tilde{\mu}_n - \mu\|_2 \leq 12 \sqrt{\frac{\mathbb{E}\|X - \mu\|_2^2 \cdot \log(\delta^{-1})}{n}}.$$

Given $\tilde{\mu}_n$, we define $\hat{\mu}_n = \arg \min_{u \in \mathbb{R}^d: \|u\|_0 \leq s} \|u - \tilde{\mu}_n\|_2$. Since $\|\mu\|_0 \leq s$, we have $\|\tilde{\mu}_n - \hat{\mu}_n\|_2 = \min_{u: \|u\|_0 \leq s} \|\tilde{\mu}_n - u\|_2 \leq \|\tilde{\mu}_n - \mu\|_2$. It follows that

$$\|\hat{\mu}_n - \mu\|_2 \leq \|\tilde{\mu}_n - \hat{\mu}_n\|_2 + \|\tilde{\mu}_n - \mu\|_2 \leq 2\|\tilde{\mu}_n - \mu\|_2.$$

Finally, we note that since $\hat{\mu}_n$ and μ are both s -sparse, $\|\hat{\mu}_n - \mu\|_1 \leq \sqrt{2s} \|\hat{\mu}_n - \mu\|_2$. □