

AVOIDING CATASTROPHE IN ONLINE LEARNING BY ASKING FOR HELP

Anonymous authors

Paper under double-blind review

ABSTRACT

Most learning algorithms with formal regret guarantees assume that no mistake is irreparable and essentially rely on trying all possible behaviors. This approach is problematic when some mistakes are *catastrophic*, i.e., irreparable. We propose an online learning problem where the goal is to minimize the chance of catastrophe. Specifically, we assume that the payoff in each round represents the chance of avoiding catastrophe that round and try to maximize the product of payoffs (the overall chance of avoiding catastrophe) while allowing a limited number of queries to a mentor. We first show that in general, any algorithm either constantly queries the mentor or is nearly guaranteed to cause catastrophe. However, in settings where the mentor policy class is learnable in the standard online model, we provide an algorithm whose regret and rate of querying the mentor both approach 0 as the time horizon grows. Conceptually, if a policy class is learnable in the absence of catastrophic risk, it is learnable in the presence of catastrophic risk if the agent can ask for help.

1 INTRODUCTION

There has been mounting concern over catastrophic risk from AI, including but not limited to autonomous weapon accidents (Abaimov & Martellini, 2020), bioterrorism (Mouton et al., 2024), and cyberattacks on critical infrastructure (Guembe et al., 2022). See Critch & Russell (2023) and Hendrycks et al. (2023) for taxonomies of societal-scale AI risks. In this paper, we use “catastrophe” to refer to any kind of irreparable harm. This definition also covers smaller-scale (yet still unacceptable) incidents such as serious medical errors (Di Nucci, 2019), crashing a robotic vehicle (Kohli & Chadha, 2020), or discriminatory sentencing (Villasenor & Foggo, 2020).

The gravity of these risks contrasts starkly with the dearth of theoretical understanding of how to avoid them. Nearly all of learning theory explicitly or implicitly assumes that no single mistake is too costly. We focus on *online learning*, where an agent repeatedly interacts with an unknown environment and uses its observations to gradually improve its performance. Most online learning algorithms essentially try all possible behaviors and see what works well. We do not want autonomous weapons or surgical robots to try all possible behaviors.

More precisely, trial-and-error-style algorithms only work when catastrophe is assumed to be impossible. This assumption manifests differently in different subtypes of online learning. In the standard online learning model, the agent’s actions have no permanent effect on the environment.¹ Online reinforcement learning allows the agent’s actions to permanently affect the environment, but typically assumes that either no action has irreversible effects (e.g., Jaksch et al. (2010)) or that the agent is reset at the start of each “episode” (e.g., Azar et al. (2017)). One could train an agent entirely in a controlled lab setting where the above assumptions do hold, but we argue that sufficiently general agents will inevitably encounter novel scenarios when deployed in the real world. Machine learning models often behave unpredictably in unfamiliar environments (see, e.g., Quionero-Candela et al. (2009)), and we do not want AI biologists or robotic vehicles to behave unpredictably.

The goal of this paper is to understand the conditions under which it is possible to formally guarantee avoidance of catastrophe in online learning. Certainly some conditions are necessary, because if the agent can only learn by trying actions directly, the problem is hopeless: any untried action could

¹More precisely, the state can depend on the agent’s previous actions, but the agent’s performance is always evaluated with respect to the optimal policy on the same sequence of states.

lead to paradise or disaster and the agent has no way to predict which. In the real world, however, one needn't learn through pure trial-and-error: one can also ask for help. **We think it is critical for high-stakes AI applications to employ a designated supervisor who can be asked for help. Examples include a human doctor supervising AI doctors, a robotic vehicle with a human driver who can take over in emergencies, autonomous weapons with a human operator, and many more.**

1.1 OUR MODEL

We propose an online learning model of avoiding catastrophe with mentor help. On each time step, the agent observes a state, selects an action (or queries the mentor), and obtains a payoff. Each payoff represents the probability of avoiding catastrophe on that time step **(conditioned on no prior catastrophe)**. The agent's goal is to maximize the *product* of payoffs, **which is equal to the overall probability of avoiding catastrophe by the chain rule of probability.**

The (possibly suboptimal) mentor has a fixed policy, and when queried, the mentor illustrates their policy's action in the current state. We desire an agent whose regret – defined as the gap between the mentor's performance and the agent's performance – approaches zero as the time horizon T grows. In other words, with enough time, the agent should avoid catastrophe nearly as well as the mentor. We also expect the agent to become self-sufficient over time: formally, the number of queries to the mentor should be sublinear in T , or equivalently, the rate of querying the mentor should go to zero.

1.2 OUR ASSUMPTIONS

The agent needs some way to make inferences about unseen states in order to decide when to ask for help. Much past work has used Bayesian inference, which suffers tractability issues in complex environments.² In this paper, we instead assume that the mentor policy satisfies a Lipschitz-flavored condition we call “local generalization”: informally, if the mentor told us that action a was safe in a similar state, then a is probably also safe in the current state. This captures the intuition that one can transfer knowledge between similar situations; **see Section 3 for the formal definition and further discussion.** Unlike Bayesian inference, local generalization only requires computing distances and is compatible with any state space which admits a distance metric.

Unlike the standard online learning model, we assume that the agent does not observe payoffs. This is because the payoff in our model represents the chance of avoiding catastrophe on that time step. In the real world, one only observes whether catastrophe occurred, not its probability.³

1.3 STANDARD ONLINE LEARNING

An overview of standard online learning is in order before discussing our results. In the standard model, the agent observes a state on each time step and must choose an action. An adversary then reveals the correct action, which results in some payoff for the agent. The goal is sublinear regret, i.e., the average regret per time step should go to 0 as $T \rightarrow \infty$. Figure 1 delineates the precise differences between the standard model and our model.

If the adversary's choices are unconstrained, the problem is hopeless: if the adversary determines the correct action on each time step randomly and independently, the agent can do no better than random guessing. However, sublinear regret becomes possible if (1) the hypothesis class has finite Littlestone dimension (Littlestone, 1988), or (2) the hypothesis class has finite VC dimension (Vapnik & Chervonenkis, 1971) and the input is σ -smooth⁴ (Haghtalab et al., 2024).

The goal of sublinear regret in online learning implicitly assumes catastrophe is impossible: the agent can make arbitrarily many (and arbitrarily costly) mistakes as long as the *average* regret per time step goes to 0. In contrast, we demand subconstant regret: the *total* probability of catastrophe should go to 0. Furthermore, standard online learning allows the agent to observe payoffs on every time step, while our agent only receives feedback on time steps with queries. However, access to a mentor (and local generalization) allows our agent to learn without trying actions directly, which is enough to offset all of the above disadvantages.

²For the curious reader, Betancourt (2018) provides a thorough treatment. See also Section 2.

³One may be able to detect “close calls” in some cases, but observing the precise probability seems unrealistic.

⁴Informally, the adversarial chooses a distribution over states instead of a precise state. See Section 3.

	Objective	Regret goal	Feedback	Mentor	Local gen.
Standard model	Sum of payoffs	Sublinear	Every time step	No	No
Our model	Product of payoffs	Subconstant	Only from queries	Yes	Yes

Figure 1: Comparison between the standard online learning model and our model. On a technical level, sublinear vs subconstant regret is a more important distinction than the sum vs product of payoffs: for example, taking a logarithm can transform the product into a sum (we actually prove a subconstant additive regret bound as an intermediate step to bounding the multiplicative regret). However, the multiplicative objective and the goal of subconstant regret are tightly related, since the former’s interpretation as the chance of catastrophe is what gives rise to the latter: we want the total (not average) chance of catastrophe to go to 0.

1.4 OUR RESULTS

Informally, we show that avoiding catastrophe with the help of a mentor and local generalization is no harder than online learning without catastrophic risk.

More precisely, we first show that in general, any algorithm with sublinear queries to the mentor has arbitrarily poor regret in the worst-case (Theorem 4.1). This means that even when the mentor can avoid catastrophe with certainty, any algorithm either needs excessive supervision or is nearly guaranteed to cause catastrophe. Unlike online learning where the general impossibility result is trivial (the agent might as well guess randomly given an unconstrained adversary), local generalization significantly limits the adversary’s power and necessitates a careful analysis.

Next, we present a simple algorithm whose total regret and rate of querying the mentor both go to 0 as $T \rightarrow \infty$ when either (1) the mentor policy class has finite Littlestone dimension or (2) the mentor policy class has finite VC dimension and the state sequence is σ -smooth. Our algorithm can handle a multi-dimensional unbounded state space and does not need detailed access to the feature embedding, instead using two simple operations. It does need to know the mentor policy class, as is standard in online learning. We initially prove the theorem for binary actions (Theorem 5.2) and then reduce learning with many actions to the binary action case (Theorem C.1).

Finally, Appendix D presents another algorithm with the same guarantees of subconstant regret and sublinear queries under different assumptions. Most importantly, this result applies to the class of thresholds on $[0, 1]$, which is known to be hard for standard online learning (see Section 2). This shows that our problem and standard online learning do not have exactly the same difficulty, but rather our problem is no harder than online learning.

The rest of the paper is structured as follows. Section 2 discusses related work. Section 3 formally defines our model. Section 4 presents our negative result for general mentor policies. Section 5 presents our positive result for simple mentor policy classes. Proofs are deferred to the appendix.

2 RELATED WORK

Learning with irreversible costs. Despite the ubiquity of irreparable/irreversible costs in the real world, theoretical work on this topic remains limited. This may be due to the fundamental modeling question of how the agent should learn about novel states or actions without actually trying them.

The most common approach is to allow the agent to ask for help. This alone is insufficient, however: the agent must have some way to decide *when* to ask for help. A popular solution is to perform Bayesian inference on the world model, but this has two tricky requirements: (1) a prior distribution which contains the true world model (or an approximation), and (2) an environment where computing (or approximating) the posterior is tractable. A finite set of possible environments satisfies both conditions, but is unrealistic in many real-world scenarios. In contrast, our algorithm can handle an uncountable policy class and a continuous unbounded state space, which is crucial for many real-world scenarios in which one never sees the exact same state twice.

Bayesian inference combined with asking for help is studied by Cohen et al. (2021); Cohen & Hutter (2020); Kosoy (2019); Mindermann et al. (2018). We also mention Hadfield-Menell et al. (2017); Moldovan & Abbeel (2012); Turchetta et al. (2016), which utilize Bayesian inference in the context of safe (online) reinforcement learning without asking for help (and without regret bounds).

We are only aware of two papers which theoretically address irreversibility without Bayesian inference: Grinsztajn et al. (2021) and Maillard et al. (2019). The former proposes to sample trajectories and learn reversibility based on temporal consistency between states: intuitively, if s_1 always precedes s_2 , we can infer that s_1 is unreachable from s_2 . Although the paper theoretically grounds this intuition, there is no formal regret guarantee. The latter presents an algorithm which asks for help in the form of rollouts from the current state. However, the regret bound and number of rollouts are both linear in the worst case, due to the dependence on the γ^* parameter which roughly captures how bad an irreversible action can be. In contrast, our algorithm achieves good regret even when actions are maximally bad.

To our knowledge, we are the first to provide an algorithm which formally guarantees avoidance of catastrophe (with high probability) without Bayesian inference. **We are also not aware of prior results comparable to our negative result, including in the Bayesian regime.**

Constrained Markov Decision Processes (CMDPs). CMDPs (Altman, 2021; Puterman, 2014) require the agent to maximize reward while also satisfying safety constraints. The two most relevant papers are Liu et al. (2021) and Stradi et al. (2024), both of which provide algorithms guaranteed to satisfy initially unknown safety constraints with high probability on every time step. However, both papers assume that the agent knows a fully safe policy upfront. In contrast, the agent in our setting has no prior knowledge. In this sense, our work complements theirs: our goal is essentially to learn the baseline safe policy that their algorithms require. One can also view our problem as the “pessimistic” model and their problem as the “optimistic” model, with some applications better captured by our model while other applications are better captured by theirs.

Online learning. See Cesa-Bianchi & Lugosi (2006) and Chapter 21 of Shalev-Shwartz & Ben-David (2014) for introductions to online learning. A classical result states that sublinear regret is possible iff the hypothesis class has finite Littlestone dimension (Littlestone, 1988). However, even some simple hypothesis classes have infinite Littlestone dimension, such as the class of thresholds on $[0, 1]$ (Example 21.4 in Shalev-Shwartz & Ben-David (2014)). Recently, Haghtalab et al. (2024) showed that if the adversary only chooses a distribution over states rather than the precise state, only the weaker assumption of finite VC dimension (Vapnik & Chervonenkis, 1971) is needed for sublinear regret. Specifically, they assume that each state is sampled from a distribution whose concentration is upper bounded by $\frac{1}{\sigma}$ times the uniform distribution. This framework is known as *smoothed analysis*, originally proposed by Spielman & Teng (2004).

Multiplicative objectives. Although online learning traditionally studies the sum of payoffs, there is some work maximizing the product of payoffs, or equivalently the sum of logarithms (Chapter 9 of Cesa-Bianchi et al. (2017)). However, these regret bounds are still sublinear in T , in comparison to our subconstant regret bounds. Also, that work still assumes that payoffs are observed on every time step, while our agent only receives feedback in response to queries (Figure 1).

Barman et al. (2023) recently provided regret bounds for a multiplicative objective in a multi-armed bandit problem, but their objective is the geometric mean of payoffs instead of the product. Interpreted in our context, their regret bounds imply that the *average* chance of catastrophe goes to zero, while we guarantee that the *total* chance of catastrophe goes to zero. This distinction is closely related to the difference between subconstant and sublinear regret discussed in Section 1.3.

Active learning and imitation learning. Our assumption that the agent only receives feedback in response to queries falls under the umbrella of active learning (Hanneke, 2014). This contrasts with passive learning, where the agent receives feedback automatically. Although ideas from active learning could be useful in our domain, we are not aware of any results from that literature which account for irreversible costs. The process of the agent learning from a mentor is also reminiscent of imitation learning (Osa et al., 2018), but we are not aware of any relevant technical implications.

3 MODEL

States. Let \mathbb{N} refer to the set of strictly positive integers and let $T \in \mathbb{N}$ be the time horizon. Let $S \subseteq \mathbb{R}^n$ be the state space⁵ and let $\mathbf{s} = (s_1, s_2, \dots, s_T) \in S^T$ be the sequence of states. In the adversarial setting, each s_t can have arbitrary dependence on the events of prior time steps. In the smoothed setting, the adversary only chooses the distribution s_t from which s_t is sampled. Formally,

⁵One could also allow a generic metric space; our assumption of $S = \mathbb{R}^n$ is only for convenience.

a distribution \mathcal{D} over S is σ -smooth if for any $X \subseteq S$, $\mathcal{D}(X) \leq \frac{1}{\sigma} U(X)$. (In the smoothed setting, we assume that S supports a uniform distribution U .⁶) If each s_t is sampled from a σ -smooth \mathcal{D}_t , we say that \mathbf{s} is σ -smooth. The sequence $\mathcal{D} = \mathcal{D}_1, \dots, \mathcal{D}_T$ can still be adaptive, i.e., the choice of \mathcal{D}_t can depend on the events of prior time steps.

Actions. Let A be a finite set of actions. There also exists a special action \hat{a} which corresponds to querying the mentor. For $k \in \mathbb{N}$, let $[k] = \{1, 2, \dots, k\}$. On each time step $t \in [T]$, the agent must select an action $a_t \in A \cup \{\hat{a}\}$, which generates payoff $\mu(s_t, a_t) \in [0, 1]$.

Asking for help. The mentor is endowed with a (not necessarily optimal) policy $\pi^m : S \rightarrow A$. When action \hat{a} is chosen, the mentor informs the agent of the action $\pi^m(s_t)$ and the agent obtains payoff $\mu(s_t, \pi^m(s_t))$. For brevity, let $\mu^m(s) = \mu(s, \pi^m(s))$. The agent never observes payoffs: the only way to learn about μ is by querying the mentor.

We would like an algorithm which becomes “self-sufficient” over time: the rate of querying the mentor should go to 0 as $T \rightarrow \infty$. Intuitively, this is equivalent to the cumulative number of queries being sublinear in T : $|Q_T| \in o(T)$, where $Q_T = \{t \in [T] : a_t = \hat{a}\}$.⁷ However, interpreted literally, $|Q_T| \in o(T)$ is not a rigorous statement. This is because $|Q_T|$ is not just a function of T : it can also depend on \mathbf{s}, μ, π^m , and any randomness in the agent’s actions.

To formalize this, let $Q_T(\mathbf{s}, \mu, \pi^m)$ be the set of time steps with queries given \mathbf{s}, μ , and π^m . Then we say that $|Q_T|$ is sublinear in T , denoted $|Q_T| \in o(T)$, if there exists $g : \mathbb{N} \rightarrow \mathbb{N}$ such that $g(T) \in o(T)$ and $|Q_T(\mathbf{s}, \mu, \pi^m)| \leq g(T)$ for all \mathbf{s}, μ, π^m (and for all possible realizations of the agent’s randomization). However, we will mostly write Q_T and $|Q_T| \in o(T)$ for brevity.

Local generalization. We assume that the mentor policy permits “local generalization”. Informally, if the agent is in state s , taking the mentor action from a similar state s' is almost as good as taking the mentor action for s . Formally, we assume there exists $L > 0$ such that for all $s, s' \in S$, $|\mu^m(s) - \mu^m(s')| \leq L \|s - s'\|$, where $\|\cdot\|$ denotes Euclidean distance. **This represents the ability to transfer knowledge between similar states:**

$$\left| \underbrace{\mu(s, \pi^m(s))}_{\text{Taking the “right” action}} - \underbrace{\mu(s, \pi^m(s'))}_{\text{Using what you learned in } s'} \right| \leq \underbrace{L \|s - s'\|}_{\text{Similarity between } s \text{ and } s'}$$

Borrowing knowledge from similar experiences seems fundamental to learning and is well-understood in the psychology literature (Esser et al., 2023) and education literature (Hajian, 2019).

Crucially, our state space can be any feature embedding of the agent’s situation, not just its physical positioning. Our algorithms will not require knowledge of the feature embedding and do not need to know L , so it suffices that there exists *some* feature embedding which satisfies local generalization. The agent does not even need to know which embedding it is. Finally, local generalization implies the more familiar Lipschitz continuity for an optimal mentor (Proposition E.1).

Multiplicative objective and regret. If $\mu(s_t, a_t) \in [0, 1]$ is the chance of avoiding catastrophe on time step t (conditioned on no prior catastrophe), then $\prod_{t=1}^T \mu(s_t, a_t)$ is the agent’s overall chance of avoiding catastrophe.⁸ For a fixed \mathbf{s} and agent actions $\mathbf{a} = (a_1, \dots, a_T)$, the agent’s *regret* is

$$R_T(\mathbf{s}, \mathbf{a}, \mu, \pi^m) = \prod_{t=1}^T \mu^m(s_t) - \prod_{t=1}^T \mu(s_t, a_t)$$

We will usually write $R_T = R_T(\mathbf{s}, \mathbf{a}, \mu, \pi^m)$ for brevity. We will study the expected regret over any randomness in \mathbf{s} and/or \mathbf{a} . We desire subconstant worst-case regret: the total (not average) expected regret should go to 0 for any μ and π^m . Formally, we want $\lim_{T \rightarrow \infty} \sup_{\mu, \pi^m} \mathbb{E}[R_T] = 0$.

The value of a bound on $\mathbb{E}[R_T]$ depends on the quality of the mentor. In particular, subconstant regret becomes trivial if $\lim_{T \rightarrow \infty} \mathbb{E}[\prod_{t=1}^T \mu^m(s_t)] = 0$. However, we think that high-stakes AI applications

⁶For example, S having finite Lebesgue measure is sufficient. Note that this does not imply boundedness. Alternatively, σ -smoothness can be defined with respect to a different distribution, as long as the Radon-Nikodym derivative is uniformly bounded; see Definition 1 of Block et al. (2022).

⁷Note that Q_T is a random variable since a_1, \dots, a_T are random variables.

⁸Conditioning on no prior catastrophe means we do not need to assume that these probabilities are independent (and if catastrophe has already occurred, this time step does not matter). This is due to the chain rule of probability.

should ensure the presence of a mentor who is almost always safe, i.e., $\mathbb{E}[\prod_{t=1}^T \mu^m(s_t)] \approx 1$. If no such mentor exists for some application, perhaps it is better to avoid the application altogether. Also, our regret bounds include rates of convergence, so even if the mentor policy is guaranteed to eventually cause catastrophe, we can still bound how quickly the agent becomes unsafe.

VC and Littlestone dimensions. VC dimension (Vapnik & Chervonenkis, 1971) and Littlestone dimension (Littlestone, 1988) are standard measures of learning difficulty which capture the ability of a hypothesis class (in our case, a policy class) to realize arbitrary combinations of labels (in our case, actions). We omit the precise dimensions since we only utilize these concept via existing results. See Shalev-Shwartz & Ben-David (2014) for a comprehensive overview.

Misc. The diameter of $X \subseteq S$ is defined by $\text{diam}(X) = \max_{x,y \in X} \|x - y\|$. All logarithms and exponents are base e unless otherwise noted.

4 AVOIDING CATASTROPHE IS IMPOSSIBLE IN GENERAL

We begin by showing that in general, any algorithm with sublinear mentor queries has arbitrary poor regret in the worst-case, even when states are i.i.d. on $[0, 1]$. The result also holds even if the algorithm knows L and s ahead of time. We use $\mathbb{E}_{\mathbf{s} \sim U^T, \mathbf{a}}$ denote the expectation over both \mathbf{s} and any additional randomness in the algorithm.

Theorem 4.1. *The worst-case expected regret of any algorithm with sublinear queries goes to 1 as T goes to infinity. Formally, $\lim_{T \rightarrow \infty} \sup_{\mu, \pi^m} \mathbb{E}_{\mathbf{s} \sim U^T, \mathbf{a}} [R_T] = 1$.*

4.1 INTUITION

We partition S into equally-sized sections that are “independent” in the sense that querying a state in section i gives you no information about section j . The number of sections is determined by a function $f : \mathbb{N} \rightarrow \mathbb{N}$ that we will choose. If $|Q_T| \in o(f(T))$, most of these sections will never contain a query. When the agent sees a state in a section not containing a query, it essentially has to guess, meaning it will be wrong a constant fraction of the time. Figure 2 fleshes out this idea.

Picking $f(T)$. A natural idea is to try $f(T) = T$, but this doesn’t quite work: even if the agent chooses wrong on every time step, the minimum payoff is still at least $1 - \frac{L}{2T}$, and $\lim_{T \rightarrow \infty} \prod_{t=1}^T (1 - \frac{L}{2T}) = \lim_{T \rightarrow \infty} (1 - \frac{L}{2T})^T = e^{-L/2}$. In order for the regret to approach 1, we need $f(T)$ to be asymptotically between $|Q_T|$ and T (such f must exist since $|Q_T| \leq g(T) \in o(T)$). This leads to the following bound: $\prod_{t=1}^T \mu(s_t, a_t) \leq (1 - \frac{L}{\Theta(f(T))})^{\Theta(T)}$. When $f(T) \in o(T)$, the right hand side converges to 0, while $\prod_{t=1}^T \mu^m(s_t) = 1$. In words, the agent is nearly guaranteed to cause catastrophe, despite the existence of a policy which is guaranteed to avoid catastrophe.

VC dimension. The class of mentor policies induced by our construction has VC dimension $f(T)$; considered over all possible values of T , this implies infinite VC (and Littlestone) dimension. This is necessary given our positive results in Section 5.

4.2 FORMAL DEFINITION OF CONSTRUCTION

Let $S = [0, 1]$ and $\mathcal{D}_t = U$ for each $t \in [T]$. Assume that $L \leq 1$; this will simplify the math and only makes the problem easier for the agent. We define a family of payoff functions parameterized by a function $f : \mathbb{N} \rightarrow \mathbb{N}$ and a bit string $\mathbf{x} = (x_1, x_2, \dots, x_{f(T)}) \in \{0, 1\}^{f(T)}$. The bit x_j will denote the optimal action in section j . Note that $f(T) \geq 1$ and since we defined \mathbb{N} to exclude 0.

For each $j \in [f(T)]$, we refer to $S_j = [\frac{j-1}{f(T)}, \frac{j}{f(T)}]$ as the j th section. Let $m_j = \frac{j-0.5}{f(T)}$ be the midpoint of S_j . Assume that each s_t belongs to exactly one S_j (this happens with probability 1, so this assumption does affect the expected regret). Let $j(s)$ denote the index of the section containing state s . Then $\mu_{f,\mathbf{x}}$ is defined by

$$\mu_{f,\mathbf{x}}(s, a) = \begin{cases} 1 & \text{if } a = x_{j(s)} \\ 1 - L \left(\frac{1}{2f(T)} - |m_{j(s)} - s| \right) & \text{if } a \neq x_{j(s)} \end{cases}$$

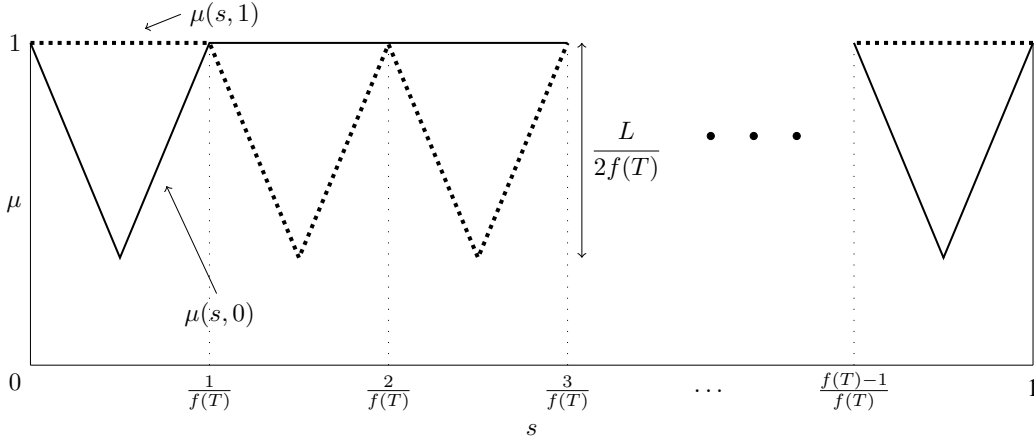


Figure 2: An illustration of the construction we use to prove Theorem 4.1 (not to scale). The horizontal axis indicates the state $s \in [0, 1]$ and the vertical axis indicates the payoff $\mu(s, a) \in [0, 1]$. The solid line represents $\mu(s, 0)$ and the dotted line represents $\mu(s, 1)$. In each section, one of the actions has the optimal payoff of 1, and the other action has the worst possible payoff allowed by L , reaching a minimum of $1 - \frac{L}{2f(T)}$. Crucially, both actions result in a payoff of 1 at the boundaries between sections: this allows us to “reset” for the next section. As a result, we can freely toggle the optimal action for each section independently.

Let π^m be any policy which is optimal for $\mu_{f, \mathbf{x}}$. Note that there is a unique optimal action for each s_t , since each s_t belongs to exactly one S_j ; formally, $\pi^m(s_t) = x_{j(s_t)}$.

For any $\mathbf{x} \in \{0, 1\}^{f(T)}$, $\mu_{f, \mathbf{x}}$ is piecewise linear (trivially) and continuous (because both actions have payoff 1 on the boundary between sections). Since the slope of each piece is in $\{-L, 0, L\}$, $\mu_{f, \mathbf{x}}$ is Lipschitz continuous. Thus by Proposition E.1, π^m satisfies local generalization.

5 AVOIDING CATASTROPHE ASSUMING FINITE VC OR LITTLESTONE DIMENSION

Theorem 4.1 shows that avoiding catastrophe is impossible in general, which is also true in online learning. What if we restrict ourselves to settings where standard online learning is possible? Specifically, we assume that π^m belongs to a policy class Π where either (1) Π has finite VC dimension d and \mathbf{s} is σ -smooth or (2) Π has finite Littlestone dimension d .⁹

This section presents a simple algorithm which guarantees subconstant regret and sublinear queries under either of those assumptions. Our algorithm needs to know Π , as is standard in online learning. The algorithm does not need to know σ (in the smooth case) or L , and can handle an unbounded state space (the number of queries simply scales with the maximum distance between observed states).

For simplicity, we initially prove our result for $A = \{0, 1\}$. Appendix C extends our result to many actions using the standard “one versus rest” reduction.¹⁰

5.1 INTUITION BEHIND THE ALGORITHM

Algorithm 1 has two simple components: (1) run a modified version of the Hedge algorithm for online learning, but (2) ask for help in unfamiliar states (specifically, when the current state is far from any queried state with the same action under the proposed policy). Hedge ensures that the number of time steps where the agent’s action doesn’t match the mentor’s is small, and asking for help in unfamiliar states ensures that when we do make a mistake, the cost isn’t too high. This algorithmic structure seems quite natural: mostly follow a baseline strategy, but ask for help when out-of-distribution.

⁹Recall from Section 1.3 that standard online learning becomes tractable under either of these assumptions.

¹⁰For each action a , we learn a binary classifier which predicts whether $\pi^m(s) = a$. If every binary classifier is correct, we can correctly determine $\pi^m(s)$. See, e.g., Chapter 29 of Shalev-Shwartz & Ben-David (2014).

Algorithm 1 successfully avoids catastrophe assuming finite VC or Littlestone dimension.

```

1: function AVOIDCATASTROPHE( $T \in \mathbb{N}$ ,  $\varepsilon \in \mathbb{R}_{>0}$ ,  $d \in \mathbb{N}$ , policy class  $\Pi$ )
2:   if  $\Pi$  has VC dimension  $d$  then
3:      $\tilde{\Pi} \leftarrow$  any smooth  $\varepsilon$ -cover of  $\Pi$  of size at most  $(41/\varepsilon)^d$  ▷ See Definition 5.3
4:   else
5:      $\tilde{\Pi} \leftarrow$  any adversarial cover of size at most  $(eT/d)^d$  ▷ See Definition 5.4
6:    $X \leftarrow \emptyset$ 
7:    $w(\pi) \leftarrow 1$  for all  $\pi \in \tilde{\Pi}$ 
8:    $p \leftarrow 1/\sqrt{\varepsilon T}$ 
9:    $\eta \leftarrow \max\left(\sqrt{\frac{p \log |\tilde{\Pi}|}{2T}}, \frac{p^2}{\sqrt{2}}\right)$ 
10:  for  $t$  from 1 to  $T$  do ▷ Run one step of Hedge, which selects policy  $\pi_t$ 
11:    hedgeQuery  $\leftarrow$  true with probability  $p$  else false
12:    if hedgeQuery then
13:      Query mentor and observe  $\pi^m(s_t)$ 
14:       $\ell(t, \pi) \leftarrow \mathbf{1}(\pi(s_t) \neq \pi^m(s_t))$  for all  $\pi \in \tilde{\Pi}$ 
15:       $\ell^* \leftarrow \min_{\pi \in \tilde{\Pi}} \ell(t, \pi)$ 
16:       $w(\pi) \leftarrow w(\pi) \cdot \exp(-\eta(\ell(t, \pi) - \ell^*))$  for all  $\pi \in \tilde{\Pi}$ 
17:       $\pi_t \leftarrow \arg \min_{\pi \in \tilde{\Pi}} \ell(t, \pi)$ 
18:    else
19:       $P(\pi) \leftarrow w(\pi) / \sum_{\pi' \in \tilde{\Pi}} w(\pi')$  for all  $\pi \in \tilde{\Pi}$ 
20:      Sample  $\pi_t \sim P$ 
21:      if  $\min_{(s,a) \in X: a=\pi_t(s_t)} \|s_t - s\| > \varepsilon^{1/n}$  then ▷ Ask for help if out-of-distribution
22:        Query mentor and observe  $\pi^m(s_t)$ 
23:         $X \leftarrow X \cup \{(s_t, \pi^m(s_t))\}$ 
24:      else ▷ Otherwise, follow Hedge's chosen policy
25:        Take action  $\pi_t(s_t)$ 

```

Simple operations. The algorithm does not require detailed access to the state embedding, instead relying on two simple operations: evaluating a policy on a particular state, and computing a nearest neighbor distance. The former seems necessary for any algorithm. The latter could be modeled as an out-of-distribution detector score, for which many methods are available (see e.g., Yang et al. (2024)).

Hedge. Hedge (Freund & Schapire, 1997) is a standard online learning algorithm which ensures sublinear regret when the number of hypotheses (in our case, the number of policies in Π) is finite.¹¹ We would prefer not to assume that Π is finite. Luckily, any policy Π can be approximated within ε when either (1) Π has finite VC dimension and is σ -smooth or (2) Π has finite Littlestone dimension. Thus we can run Hedge on this approximative policy class instead.

One other modification is necessary. In standard online learning, losses are observed on every time step, but our agent only receives feedback in response to queries. To handle this, we modify Hedge to only perform updates on time steps with queries and to issue a query with probability p on each time step. Continuing our lucky streak, Russo et al. (2024) analyzes exactly this modification of Hedge.

We prove the following theorem parametrized by ε :

Theorem 5.1. *Let $A = \{0, 1\}$. Assume $\pi^m \in \Pi$ where either (1) Π has finite VC dimension d , s is σ -smooth, and $\varepsilon T \log T > 12\sigma d \log(4e^2/\varepsilon)$ or (2) Π has finite Littlestone dimension d . Then for any $T \in \mathbb{N}$ and $\varepsilon > 0$, Algorithm 1 satisfies*

$$\mathbb{E}[R_T] \in O\left(\frac{dL}{\sigma} T \varepsilon^{1+1/n} \log(1/\varepsilon) \log T\right)$$

$$\mathbb{E}[|Q_T|] \in O\left(\sqrt{\frac{T}{\varepsilon}} + \frac{d}{\sigma} T \varepsilon \log(1/\varepsilon) \log T + \frac{\text{diam}(s)^n}{\varepsilon}\right)$$

¹¹See Chapter 5 of Slivkins et al. (2019) and Chapter 21 of Shalev-Shwartz & Ben-David (2014) for modern introductions to Hedge.

In Case 1, the expectation is over the randomness of both \mathbf{s} and the algorithm, while in Case 2, the expectation is over only the randomness of the algorithm. Also, R_T and Q_T clearly have no dependence on σ in Case 2, but we include σ anyway to avoid writing two separate bounds.

To obtain subconstant regret and sublinear queries, we can choose $\varepsilon = T^{\frac{-2n}{2n+1}}$. This also satisfies the requirement of $\varepsilon T \log T > 12\sigma d \log(4e^2/\varepsilon)$ for large enough T .

Theorem 5.2. *Let $A = \{0, 1\}$. Assume $\pi^m \in \Pi$ where either (1) Π has finite VC dimension d and \mathbf{s} is σ -smooth or (2) Π has finite Littlestone dimension d . Then for any $T \in \mathbb{N}$, Algorithm 1 with $\varepsilon = T^{\frac{-2n}{2n+1}}$ satisfies*

$$\begin{aligned}\mathbb{E}[R_T] &\in O\left(\frac{dL}{\sigma} T^{\frac{-1}{2n+1}} \log T\right) \\ \mathbb{E}[|Q_T|] &\in O\left(T^{\frac{4n+1}{4n+2}} \left(\frac{d}{\sigma} \log T + \text{diam}(\mathbf{s})^n\right)\right)\end{aligned}$$

We can also make the query bound unconditional in exchange for slightly worse bounds. The probability that $|Q_T|$ asymptotically exceeds the bound from Theorem 5.2 goes to 0, so the query threshold in Corollary 5.2.1 is almost never reached. Thus we still obtain subconstant regret.

Corollary 5.2.1. *If Algorithm 1 is modified to stop querying after $T^{\frac{4n+1.5}{4n+2}} (\frac{d}{\sigma} \log T + \text{diam}(\mathbf{s})^n)$ queries, then under the conditions of Theorem 5.2, we have $\mathbb{E}[R_T] \in O(\frac{dL}{\sigma} T^{\frac{-1}{8n+4}} \log T)$.*

5.2 PROOF SKETCH

The formal proof of Theorem 5.1 can be found in Appendix B, but we outline the key elements here. The regret analysis consists of two ingredients: analyzing the Hedge component, and analyzing the “ask for help when out-of-distribution” component. The former will bound the number of mistakes made by the algorithm (i.e., the number of time steps where the agent’s action doesn’t match the mentor’s), and the latter will bound the cost of any single mistake. We must also carefully show that the latter does not result in excessively many queries, which we do via a novel packing argument.

We begin by formalizing two notion of approximating a policy class:

Definition 5.3. Let U be the uniform distribution over S . For $\varepsilon > 0$, a policy class $\tilde{\Pi}$ is a *smooth ε -cover* of a policy class Π is for every $\pi \in \Pi$, there exists $\tilde{\pi} \in \tilde{\Pi}$ such that $\Pr_{s \sim U}[\pi(s) \neq \tilde{\pi}(s)] \leq \varepsilon$.

Definition 5.4. A policy class $\tilde{\Pi}$ is an *adversarial cover* of a policy class Π is for every $\mathbf{s} \in S^T$ and $\pi \in \Pi$, there exists $\tilde{\pi} \in \tilde{\Pi}$ such that $\pi(s_t) = \tilde{\pi}(s_t)$ for all $t \in [T]$.

The existence of small covers is crucial:

Lemma 5.1 (Lemma 7.3.2 in Haghtalab (2018)¹²). *For all $\varepsilon > 0$, any policy class of VC dimension d admits a smooth ε -cover of size at most $(41/\varepsilon)^d$.*

Lemma 5.2 (Lemmas 21.13 and A.5 in Shalev-Shwartz & Ben-David (2014)). *Any policy class of Littlestone dimension d admits an adversarial cover of size at most $(eT/d)^d$.*

An adversarial cover is a perfect cover by definition. The following lemma establishes that a smooth ε -cover is a good approximation for any sequence of σ -smooth distributions.

Lemma 5.3 (Equation 2 and Lemma 3.3 in Haghtalab et al. (2024)). *Let $\tilde{\Pi}$ be a finite smooth ε -cover of Π and let $\mathcal{D} = \mathcal{D}_1, \dots, \mathcal{D}_T$ be a sequence of σ -smooth distributions. If $\varepsilon T \log T > 12\sigma d \log(4e^2/\varepsilon)$, then $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}} \left[\sup_{\pi \in \Pi} \min_{\tilde{\pi} \in \tilde{\Pi}} \sum_{t=1}^T \mathbf{1}(\pi(s_t) \neq \tilde{\pi}(s_t)) \right] \in O\left(\frac{1}{\sigma} T \varepsilon \log T \sqrt{d \log(1/\varepsilon)}\right)$.*

We will run a variant of Hedge on $\tilde{\Pi}$. The vanilla Hedge algorithm operates in the standard online learning model where on each time step, the agent selects a policy (or more generally, a hypothesis), and observes the *loss* of every policy. In general the loss function can depend arbitrarily on the time

¹²See also Haussler & Long (1995) or Lemma 13.6 in Boucheron et al. (2013) for variants which are less convenient for our purposes.

step, the policy, and prior events, but we will only use the indicator loss function $\ell(t, \pi) = \mathbf{1}(\pi(s_t) \neq \pi^m(s_t))$. Crucially, whenever we query and learn $\pi^m(s_t)$, we can compute $\ell(t, \pi)$ for every $\pi \in \tilde{\Pi}$.

We cannot afford to query on every time step, however. Recently, [Russo et al. \(2024\)](#) analyzed a variant of Hedge where losses are observed only in response to queries, which they call “label-efficient feedback”. They proved a regret bound when a query is issued on each time step with fixed probability p . Lemma 5.4 restates their result in a form that is more convenient for us (see Appendix B for details). Although their result is stated for non-adaptive adversaries, we explain in Appendix B.3 how their argument easily generalizes to adaptive adversaries. Full pseudocode for HEDGETHWITHQUERIES can also be found in the appendix (Algorithm 2).

Lemma 5.4 (Lemma 3.5 in [Russo et al. \(2024\)](#)). *Assume $\tilde{\Pi}$ is finite. Then for any loss function $\ell : [T] \times \tilde{\Pi} \rightarrow [0, 1]$ and query probability p , HEDGETHWITHQUERIES enjoys the regret bound*

$$\sum_{t=1}^T \mathbb{E}[\ell(t, \pi_t)] - \min_{\tilde{\pi} \in \tilde{\Pi}} \sum_{t=1}^T \ell(t, \tilde{\pi}) \leq \frac{2 \log |\tilde{\Pi}|}{p^2}$$

where π_t is the policy chosen at time t and the expectation is over the randomness of the algorithm.

We apply Lemma 5.4 with $\ell(t, \pi) = \mathbf{1}(\pi(s_t) \neq \pi^m(s_t))$ and combine this with Lemmas 5.1 and 5.3 (in the σ -smooth case) and with Lemma 5.2 (in the adversarial case). This yields a $O\left(\frac{d}{\sigma} T \varepsilon \log(1/\varepsilon) \log T\right)$ bound on the number of mistakes made by Algorithm 1 (Lemma B.1).

The other key ingredient of the proof is analyzing the “ask for help when out-of-distribution” component. Combined with the local generalization assumption, this allows us to fairly easily bound the cost of a single mistake (Lemma B.2). The trickier part is bounding the number of resulting queries. It is tempting to claim that the states queried in the out-of-distribution case must all be separated by at least $\varepsilon^{1/n}$ and thus form an $\varepsilon^{1/n}$ -packing, but this is actually not true. **Instead, we provide a novel method for bounding the number of data points (i.e., queries) needed to cover a set with respect to the realized actions of the algorithm Lemma B.7. This is in contrast to vanilla packing arguments which consider all data points in aggregate. Our method may be useful in other contexts where a more refined packing argument is needed.**

6 CONCLUSION AND FUTURE WORK

In this paper, we proposed a model of avoiding catastrophe in online learning. We showed that achieving subconstant regret in our problem (with the help of a mentor and local generalization) is no harder than achieving sublinear regret in standard online learning.

There remain some technical questions within this paper’s model. One question is whether the time complexity of Algorithm 1 be improved, which currently stands at $\Omega(|\tilde{\Pi}| \cdot T)$ plus the time to compute the ε -cover. Also, we have not resolved whether our problem is tractable for finite VC dimension and fully adversarial inputs (although Appendix D shows that the problem is tractable for at least some classes with finite VC but infinite Littlestone dimension).

We are also interested in alternatives to the local generalization assumption, since an action which is safe in one state may not always be safe in a nearby state (depending on the definition of “nearby”). Some assumption is necessary, or there is no way to avoid costly mistakes without querying on every time step. One possibility is Bayesian inference. We intentionally avoided Bayesian approaches in this paper due to tractability concerns, but it seems premature to abandon those ideas entirely.

Finally, we are excited to apply the ideas in this paper to Markov Decision Processes (MDPs): specifically, MDPs where some actions are irreversible (“non-communicating”) and the agent only gets one attempt (“single-episode”). In such MDPs, the agent must not only avoid catastrophe but also obtain high reward. As discussed in Section 2, very little theory exists for RL in non-communicating single-episode MDPs. Can an agent learn near-optimal behavior in high-stakes environments while becoming self-sufficient over time? Formally, we pose the following open problem:

Is there an algorithm for non-communicating single-episode undiscounted MDPs which ensures that both the regret and the number of mentor queries are sublinear in T ?

REFERENCES

- Stanislav Abaimov and Maurizio Martellini. *Artificial Intelligence in Autonomous Weapon Systems*, pp. 141–177. Springer International Publishing, Cham, 2020. ISBN 978-3-030-28285-1. doi: 10.1007/978-3-030-28285-1_8. URL https://doi.org/10.1007/978-3-030-28285-1_8.
- Eitan Altman. *Constrained Markov decision processes*. Routledge, 2021.
- Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax Regret Bounds for Reinforcement Learning. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 263–272. PMLR, July 2017. URL <https://proceedings.mlr.press/v70/azar17a.html>. ISSN: 2640-3498.
- Siddharth Barman, Arindam Khan, Arnab Maiti, and Ayush Sawarni. Fairness and welfare quantification for regret in multi-armed bandits. In *Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence*, volume 37 of AAAI’23/IAAI’23/EAAI’23, pp. 6762–6769. AAAI Press, February 2023. ISBN 978-1-57735-880-0. doi: 10.1609/aaai.v37i6.25829. URL <https://doi.org/10.1609/aaai.v37i6.25829>.
- Michael Betancourt. A Conceptual Introduction to Hamiltonian Monte Carlo, July 2018. URL <http://arxiv.org/abs/1701.02434>. arXiv:1701.02434 [stat].
- Adam Block, Yuval Dagan, Noah Golowich, and Alexander Rakhlin. Smoothed online learning is as easy as statistical learning. In *Conference on Learning Theory*, pp. 1716–1786. PMLR, 2022.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 02 2013. ISBN 9780199535255. doi: 10.1093/acprof:oso/9780199535255.001.0001. URL <https://doi.org/10.1093/acprof:oso/9780199535255.001.0001>.
- Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.
- Nicolò Cesa-Bianchi, Pierre Gaillard, Claudio Gentile, and Sébastien Gerchinovitz. Algorithmic Chaining and the Role of Partial Feedback in Online Nonparametric Learning. In *Proceedings of the 2017 Conference on Learning Theory*, pp. 465–481. PMLR, June 2017. URL <https://proceedings.mlr.press/v65/cesa-bianchi17a.html>. ISSN: 2640-3498.
- Michael K. Cohen and Marcus Hutter. Pessimism About Unknown Unknowns Inspires Conservatism. In *Proceedings of Thirty Third Conference on Learning Theory*, pp. 1344–1373. PMLR, July 2020. URL <https://proceedings.mlr.press/v125/cohen20a.html>. ISSN: 2640-3498.
- Michael K. Cohen, Elliot Catt, and Marcus Hutter. Curiosity Killed or Incapacitated the Cat and the Asymptotically Optimal Agent. *IEEE Journal on Selected Areas in Information Theory*, 2(2):665–677, June 2021. ISSN 2641-8770. doi: 10.1109/JSAIT.2021.3079722. URL <https://ieeexplore.ieee.org/document/9431093>. Conference Name: IEEE Journal on Selected Areas in Information Theory.
- Andrew Critch and Stuart Russell. Tasra: a taxonomy and analysis of societal-scale risks from ai. *arXiv preprint arXiv:2306.06924*, 2023.
- Ezio Di Nucci. Should we be afraid of medical ai? *Journal of Medical Ethics*, 45(8):556–558, 2019.
- Sarah Esser, Hilde Haider, Clarissa Lustig, Takumi Tanaka, and Kanji Tanaka. Action–effect knowledge transfers to similar effect stimuli. *Psychological Research*, 87(7):2249–2258, October 2023. ISSN 1430-2772. doi: 10.1007/s00426-023-01800-4. URL <https://doi.org/10.1007/s00426-023-01800-4>.
- Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1):119–139, 1997.

- Nathan Grinsztajn, Johan Ferret, Olivier Pietquin, philippe preux, and Matthieu Geist. There Is No Turning Back: A Self-Supervised Approach for Reversibility-Aware Reinforcement Learning. In *Advances in Neural Information Processing Systems*, volume 34, pp. 1898–1911. Curran Associates, Inc., 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/0e98aeeb54acf612b9eb4e48a269814c-Abstract.html>.
- Blessing Guembe, Ambrose Azeta, Sanjay Misra, Victor Chukwudi Osamor, Luis Fernandez-Sanz, and Vera Pospelova. The Emerging Threat of Ai-driven Cyber Attacks: A Review. *Applied Artificial Intelligence*, 36(1), December 2022. ISSN 0883-9514, 1087-6545. doi: 10.1080/08839514.2022.2037254. URL <https://www.tandfonline.com/doi/full/10.1080/08839514.2022.2037254>.
- Dylan Hadfield-Menell, Smitha Milli, Pieter Abbeel, Stuart Russell, and Anca D. Dragan. Inverse reward design. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, pp. 6768–6777, Red Hook, NY, USA, December 2017. Curran Associates Inc. ISBN 978-1-5108-6096-4.
- Nika Haghtalab. *Foundation of Machine Learning, by the People, for the People*. PhD thesis, Microsoft Research, 2018.
- Nika Haghtalab, Tim Roughgarden, and Abhishek Shetty. Smoothed analysis with adaptive adversaries. *Journal of the ACM*, 71(3):1–34, 2024.
- Shiva Hajian. Transfer of Learning and Teaching: A Review of Transfer Theories and Effective Instructional Practices. *IAFOR Journal of Education*, 7(1):93–111, 2019. URL <https://eric.ed.gov/?id=EJ1217940>. Publisher: International Academic Forum ERIC Number: EJ1217940.
- Steve Hanneke. Theory of disagreement-based active learning. 7(2–3):131–309, j=Jun 2014. ISSN 1935-8237. doi: 10.1561/22000000037. URL <https://doi.org/10.1561/22000000037>.
- David Haussler and Philip M Long. A generalization of Sauer’s lemma. *Journal of Combinatorial Theory, Series A*, 71(2):219–240, August 1995. ISSN 0097-3165. doi: 10.1016/0097-3165(95)90001-2. URL <https://www.sciencedirect.com/science/article/pii/0097316595900012>.
- Dan Hendrycks, Mantas Mazeika, and Thomas Woodside. An overview of catastrophic ai risks, 2023. URL <https://arxiv.org/abs/2306.12001>.
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal Regret Bounds for Reinforcement Learning. *Journal of Machine Learning Research*, 11(51):1563–1600, 2010. ISSN 1533-7928. URL <http://jmlr.org/papers/v11/jaksch10a.html>.
- Heinrich Jung. Ueber die kleinste kugel, die eine räumliche figur einschliesst. *Journal für die reine und angewandte Mathematik*, 123:241–257, 1901. URL <http://eudml.org/doc/149122>.
- Puneet Kohli and Anjali Chadha. Enabling pedestrian safety using computer vision techniques: A case study of the 2018 uber inc. self-driving car crash. In *Advances in Information and Communication: Proceedings of the 2019 Future of Information and Communication Conference (FICC), Volume 1*, pp. 261–279. Springer, 2020.
- Vanessa Kosoy. Delegative Reinforcement Learning: learning to avoid traps with a little help. arXiv, July 2019. doi: 10.48550/arXiv.1907.08461. URL <http://arxiv.org/abs/1907.08461>. arXiv:1907.08461 [cs, stat].
- Nick Littlestone. Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. *Machine learning*, 2:285–318, 1988.
- Tao Liu, Ruida Zhou, Dileep Kalathil, Panganamala Kumar, and Chao Tian. Learning policies with zero or bounded constraint violation for constrained mdps. *Advances in Neural Information Processing Systems*, 34:17183–17193, 2021.

- Odalric-Ambrym Maillard, Timothy Mann, Ronald Ortner, and Shie Mannor. Active Roll-outs in MDP with Irreversible Dynamics. July 2019. URL <https://hal.science/hal-02177808>.
- Sören Mindermann, Rohin Shah, Adam Gleave, and Dylan Hadfield-Menell. Active Inverse Reward Design. In *Proceedings of the 1st Workshop on Goal Specifications for Reinforcement Learning*, 2018. URL <http://arxiv.org/abs/1809.03060>. arXiv:1809.03060 [cs, stat].
- Teodor Mihai Moldovan and Pieter Abbeel. Safe exploration in Markov decision processes. In *Proceedings of the 29th International Conference on Machine Learning*, ICML’12, pp. 1451–1458, Madison, WI, USA, June 2012. Omnipress. ISBN 978-1-4503-1285-1.
- Christopher A Mouton, Caleb Lucas, and Ella Guest. The operational risks of ai in large-scale biological attacks, 2024.
- T. Osa, J. Pajarinen, G. Neumann, J.A. Bagnell, P. Abbeel, and J. Peters. *An Algorithmic Perspective on Imitation Learning*. Foundations and trends in robotics. Now Publishers, 2018. ISBN 978-1-68083-410-9. URL <https://books.google.com/books?id=6p6EtQEACAAJ>.
- Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- Joaquin Quionero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D. Lawrence. *Dataset Shift in Machine Learning*. The MIT Press, 2009. ISBN 0262170051.
- Matteo Russo, Andrea Celli, Riccardo Colini Baldeschi, Federico Fusco, Daniel Haimovich, Dima Karamshuk, Stefano Leonardi, and Niek Tax. Online learning with sublinear best-action queries. *arXiv preprint arXiv:2407.16355*, 2024.
- Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 1 edition, May 2014. ISBN 978-1-107-05713-5 978-1-107-29801-9. doi: 10.1017/CBO9781107298019. URL <https://www.cambridge.org/core/product/identifier/9781107298019/type/book>.
- Aleksandrs Slivkins et al. Introduction to multi-armed bandits. *Foundations and Trends® in Machine Learning*, 12(1-2):1–286, 2019.
- Daniel A Spielman and Shang-Hua Teng. Smoothed analysis of algorithms: Why the simplex algorithm usually takes polynomial time. *Journal of the ACM (JACM)*, 51(3):385–463, 2004.
- Francesco Emanuele Stradi, Matteo Castiglioni, Alberto Marchesi, and Nicola Gatti. Learning adversarial mdps with stochastic hard constraints. *arXiv preprint arXiv:2403.03672*, 2024.
- Matteo Turchetta, Felix Berkenkamp, and Andreas Krause. Safe Exploration in Finite Markov Decision Processes with Gaussian Processes. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL <https://proceedings.neurips.cc/paper/2016/hash/9a49a25d845a483fae4be7e341368e36-Abstract.html>.
- V. N. Vapnik and A. Ya. Chervonenkis. On the Uniform Convergence of Relative Frequencies of Events to Their Probabilities. *Theory of Probability & Its Applications*, 16(2):264–280, January 1971. ISSN 0040-585X. doi: 10.1137/1116025. URL <https://epubs.siam.org/doi/10.1137/1116025>. Publisher: Society for Industrial and Applied Mathematics.
- John Villasenor and Virginia Foggo. Artificial intelligence, due process and criminal sentencing. *Mich. St. L. Rev.*, pp. 295, 2020.
- Yihong Wu. *Lecture notes on: Information-theoretic methods for high-dimensional statistics*. 2020.
- Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution detection: A survey. *International Journal of Computer Vision*, pp. 1–28, 2024.

A PROOF OF THEOREM 4.1

A.1 PROOF ROADMAP

Throughout the proof, let V_j be the set of time steps $t \leq T$ where $|m_j - s_t| \leq \frac{1}{4f(T)}$. In words, s_t is relatively close to the midpoint of S_j . This will imply that the suboptimal action is in fact quite suboptimal. This also implies that s_t is in S_j , since each S_j has length $1/f(T)$.

Recall that $\mathbb{E}_{\mathbf{s} \sim U^T, \mathbf{a}}$ denotes the expectation over both the random states and any additional randomness in the algorithm.

The proof proceeds via the following steps:

1. Prove that $f(T) = \sqrt{(|Q_T| + 1)T}$ is asymptotically between $|Q_T|$ and T (Lemma A.1).
2. Provide a simple variant of the Chernoff bound which we will apply multiple times (Lemma A.2).
3. Show that with high probability, $\sum_{j \in A} |V_j|$ is adequately large (Lemma A.3).
4. The key lemma is Lemma A.4, which shows that a randomly sampled \mathbf{x} produces poor agent performance with high probability. The central idea is that at least $f(T) - |Q_T|$ sections are never queried (which is large, by Lemma A.1), so the agent has no way of knowing the optimal action in those sections. As a result, the agent picks the wrong answer at least half the time on average (and at least a quarter of the time with high probability). Lemma A.3 implies that a constant fraction of those time steps will have quite suboptimal payoffs, again with high probability.
5. To complete the proof, we observe that $\sup_{\mu} \mathbb{E}_{\mathbf{s} \sim U^T, \mathbf{a}} R_T(\mathbf{s}, \mathbf{a}, \mu, \pi^m) \geq \mathbb{E}_{\mathbf{x} \sim U(\{0,1\}^{f(T)})} \mathbb{E}_{\mathbf{s} \sim U^T, \mathbf{a}} R_T(\mathbf{s}, \mathbf{a}, \mu_{f,\mathbf{x}}, \pi^m)$. This is essentially an application of the probabilistic method: if a randomly chosen $\mu_{f,\mathbf{x}}$ has high expected regret, then the worst case μ also has high expected regret.

Note that \mathbf{s} , \mathbf{a} , and \mathbf{x} are random variables, so all variables defined on top of them (such as V_j) are also random variables. In contrast, the partition $S = \{S_1, \dots, S_{f(T)}\}$ and properties thereof (like the midpoints m_j) are not random variables.

Lastly, while the intuition provided in Section 4.1 is accurate, the analysis will mostly occur in log space, so the bounds will look different. However, bounds of the form discussed in Section 4.1 can still be found as an intermediate step in Part 4 of the proof of Lemma A.4.

A.2 PROOF

Lemma A.1. *Let $a, b : \mathbb{N} \rightarrow \mathbb{N}$ be functions such that $a(x) \in o(b(x))$. Then $c(x) = \sqrt{a(x)b(x)}$ satisfies $a(x) \in o(c(x))$ and $c(x) \in o(b(x))$.*

Proof. Since a and b are strictly positive (and thus c is as well), we have

$$\frac{a(x)}{c(x)} = \frac{a(x)}{\sqrt{a(x)b(x)}} = \sqrt{\frac{a(x)}{b(x)}} = \frac{\sqrt{a(x)b(x)}}{b(x)} = \frac{c(x)}{b(x)}$$

Then $a(x) \in o(b(x))$ implies

$$\lim_{x \rightarrow \infty} \frac{a(x)}{c(x)} = \lim_{x \rightarrow \infty} \frac{c(x)}{b(x)} = \lim_{x \rightarrow \infty} \sqrt{\frac{a(x)}{b(x)}} = 0$$

as required. \square

Lemma A.2. *Let z_1, \dots, z_n be i.i.d. variables in $\{0, 1\}$ and let $Z = \sum_{i=1}^n z_i$. If $\mathbb{E}[Z] \geq M$, then $\Pr[Z \leq M/2] \leq \exp(-M/8)$.*

Proof. By the Chernoff bound for i.i.d. binary variables, we have $\Pr[Z \leq \mathbb{E}[Z]/2] \leq \exp(-\mathbb{E}[Z]/8)$. Since $-\mathbb{E}[Z] \leq -M$ and \exp is an increasing function, we have $\exp(-\mathbb{E}[Z]/8) \leq \exp(-M/8)$. Also, $M/2 \leq \mathbb{E}[Z]/2$ implies $\Pr[Z \leq M/2] \leq \Pr[Z \leq \mathbb{E}[Z]/2]$. Combining these inequalities proves the lemma. \square

Lemma A.3. *let $A \subseteq [f(T)]$ be any nonempty subset of sections. Then*

$$\Pr \left[\sum_{j \in A} |V_j| \leq \frac{T|A|}{4f(T)} \right] \leq \exp \left(\frac{-T}{16f(T)} \right)$$

Proof. Fix any $j \in [f(T)]$. For each $t \in [T]$, define the random variable z_t by $z_t = 1$ if $t \in V_j$ for some $j \in A$ and 0 otherwise. We have $t \in V_j$ iff s_t falls within a particular interval of length $\frac{1}{2f(T)}$. Since these intervals are disjoint for different j 's, we have $z_t = 1$ iff s_t falls within a portion of the state space with total measure $\frac{|A|}{2f(T)}$. Since s_t is uniformly random across $[0, 1]$, we have $\mathbb{E}[z_t] = \frac{|A|}{2f(T)}$. Then $\mathbb{E}[\sum_{t=1}^T z_t] = \mathbb{E}[\sum_{j \in A} |V_j|] = \frac{T|A|}{2f(T)}$. Furthermore, since s_1, \dots, s_T are i.i.d., so are z_1, \dots, z_T . Then by Lemma A.2,

$$\Pr \left[\sum_{j \in A} |V_j| \leq \frac{T|A|}{4f(T)} \right] \leq \exp \left(\frac{-T|A|}{16f(T)} \right) \leq \exp \left(\frac{-T}{16f(T)} \right)$$

with the last step due to $|A| \geq 1$. \square

Lemma A.4. *Independently sample $\mathbf{x} \sim U(\{0, 1\}^{f(T)})$ and $\mathbf{s} \sim U^T$.¹³ Then with probability at least $1 - \exp(-\frac{T}{16f(T)}) - \exp(-\frac{f(T) - |Q_T|}{16})$,*

$$\prod_{t=1}^T \mu_{f, \mathbf{x}}(s_t, a_t) \leq \exp \left(-\frac{LT(f(T) - |Q_T|)}{2^7 f(T)^2} \right)$$

Proof. Part 1: setup. Let $J_{-Q} = \{j \in [f(T)] : s_t \notin S_j \forall t \in Q_T\}$ be the set of sections that are never queried. Since each query appears in exactly one section (because each state appears in exactly one section), $|J_{-Q}| \geq f(T) - |Q_T|$.

For each $j \in J_{-Q}$, let y_j be the action taken most frequently among time steps in V_j :

$$y_j = \arg \max_{a \in \{0, 1\}} |\{t \in V_j : a = a_t\}|$$

Let $\bar{J} = \{j \in J_{-Q} : x_j \neq y_j\}$. For each $j \in \bar{J}$, let $V'_j = \{t \in V_j : a_t \neq x_j\}$ be the set of time steps where the agent chooses the wrong action (assuming payoff function $\mu_{f, \mathbf{x}}$).

Part 2: \bar{J} is not too small. Define a random variable $z_j = \mathbf{1}_{j \in \bar{J}}$ for each $j \in J_{-Q}$. By definition, if $j \in J_{-Q}$, no state in S_j is queried. Since queries outside of S_j provide no information about x_j , the agent's actions must be independent of x_j . In particular, the random variables x_j and y_j are independent. Combining that independence with $\Pr[x_j = 0] = \Pr[x_j = 1] = 0.5$ yields $\Pr[z_j = 1] = 0.5$ for all $j \in J_{-Q}$. Furthermore, since $x_1, \dots, x_{f(T)}$ are independent, the random variables $\{z_j : j \in J_{-Q}\}$ are also independent. Since $\mathbb{E}[|\bar{J}|] = \mathbb{E}[\sum_{j \in J_{-Q}} z_j] = |J_{-Q}|/2 \geq \frac{f(T) - |Q_T|}{2}$, Lemma A.2 implies that

$$\Pr \left[|\bar{J}| \leq \frac{f(T) - |Q_T|}{4} \right] \leq \exp \left(-\frac{f(T) - |Q_T|}{16} \right)$$

Part 3: $|V'_j| \geq |V_j|/2$. Since $j \in J_{-Q}$, the mentor is not queried on any time step $t \in V_j$, so $a_t \in \{0, 1\}$ for all $t \in V_j$. Since the agent chooses one of two actions for each $t \in V_j$, the more

¹³That is, the entire set $\{x_1, \dots, x_{f(T)}, s_1, \dots, s_T\}$ is mutually independent.

frequent action must be chosen at least half of the time: $a_t = y_j$ for at least half of the time steps in V_j . Since $x_j \neq y_j$ for $j \in \bar{J}$, we have $a_t = y_j \neq x_j$ for those time steps, so $|V'_j| \geq |V_j|/2$.

Part 4: a bound in terms of \bar{J} and V_j . Consider any $j \in \bar{J}$ and $t \in V'_j \subseteq V_j$. By definition of V_j , we have $|m_j - s_t| \leq \frac{1}{4f(T)}$. Then by definition of $\mu_{f,\mathbf{x}}$,

$$\begin{aligned}\mu_{f,\mathbf{x}}(s_t, a_t) &= 1 - L \left(\frac{1}{2f(T)} - |s_t - m_j| \right) \\ &\leq 1 - L \left(\frac{1}{2f(T)} - \frac{1}{4f(T)} \right) \\ &= 1 - \frac{L}{4f(T)}\end{aligned}$$

Now aggregating across time steps,

$$\begin{aligned}\prod_{t=1}^T \mu_{f,\mathbf{x}}(s_t, a_t) &\leq \prod_{j \in \bar{J}} \prod_{t \in V'_j} \mu_{f,\mathbf{x}}(s_t, a_t) \quad (\mu_{f,\mathbf{x}}(s_t, a_t) \in [0, 1] \text{ for all } t) \\ &\leq \prod_{j \in \bar{J}} \left(1 - \frac{L}{4f(T)} \right)^{|V'_j|} \quad (\text{bound on } \mu_{f,\mathbf{x}}(s_t, a_t) \text{ when } t \in V'_j) \\ &\leq \prod_{j \in \bar{J}} \left(1 - \frac{L}{4f(T)} \right)^{|V_j|/2} \quad (|V'_j| \geq |V_j|/2)\end{aligned}$$

The last step also relies on $1 - \frac{L}{4f(T)} \in [0, 1]$, which is due to $L \leq 1$ and $f(T) \geq 1$. Converting into log space and using the standard inequality $\log(1 + x) \leq x$ for all $x \in \mathbb{R}$, we have

$$\begin{aligned}\log \prod_{t=1}^T \mu_{f,\mathbf{x}}(s_t, a_t) &\leq \log \prod_{j \in \bar{J}} \left(1 - \frac{L}{4f(T)} \right)^{|V_j|/2} \\ &= \sum_{j \in \bar{J}} \frac{|V_j|}{2} \log \left(1 - \frac{L}{4f(T)} \right) \\ &\leq - \sum_{j \in \bar{J}} \frac{L|V_j|}{8f(T)}\end{aligned}$$

Part 5: putting it all together. By Lemma A.3, Part 2 of this lemma, and the union bound, with probability at least $1 - \exp(-\frac{T}{16f(T)}) - \exp(-\frac{f(T)-|Q_T|}{16})$ we have $\sum_{j \in \bar{J}} |V_j| \geq \frac{T|\bar{J}|}{4f(T)}$ for all $j \in [f(T)]$ and $|\bar{J}| \geq \frac{f(T)-|Q_T|}{4}$. Assuming those inequalities hold, we have

$$\begin{aligned}\log \prod_{t=1}^T \mu_{f,\mathbf{x}}(s_t, a_t) &\leq - \sum_{j \in \bar{J}} \frac{L|V_j|}{8f(T)} \\ &\leq - \frac{L}{8f(T)} \cdot \frac{T|\bar{J}|}{4f(T)} \\ &\leq - \frac{L}{8f(T)} \cdot \frac{T}{4f(T)} \cdot \frac{f(T)-|Q_T|}{4} \\ &= - \frac{LT(f(T)-|Q_T|)}{2^7 f(T)^2}\end{aligned}$$

Exponentiating both sides proves the lemma. \square

Let $\alpha(T) = \exp(-\frac{T}{16f(T)}) + \exp(-\frac{f(T)-|Q_T|}{16})$ for brevity.

Theorem 4.1. *The worst-case expected regret of any algorithm with sublinear queries goes to 1 as T goes to infinity. Formally, $\lim_{T \rightarrow \infty} \sup_{\mu, \pi^m} \mathbb{E}_{\mathbf{s} \sim U^T, \mathbf{a}} [R_T] = 1$.*

Proof. If the algorithm has sublinear queries, then there exists $g(T) \in o(T)$ such that $|Q_T| \leq g(T)$ for all \mathbf{s}, μ, π^m . Let $f(T) = \sqrt{(g(T) + 1)T}$. Then by Lemma A.1, $g(T) \in o(f(T))$ and $f(T) \in o(T)$. Combining this with $|Q_T| \leq g(T)$, we get $\lim_{T \rightarrow \infty} \alpha(T) = 0$. Also, since $|Q_T| \in o(f(T))$, there exists T_0 such that $|Q_T| \leq f(T)/2$ for all $T \geq T_0$. Combining this with Lemma A.4 and noting that $\prod_{t=1}^T \mu_{f, \mathbf{x}}(s_t, a_t) \leq 1$, we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim U(\{0,1\}^{f(T)})} \mathbb{E}_{\mathbf{s} \sim U^T, \mathbf{a}} \prod_{t=1}^T \mu_{f, \mathbf{x}}(s_t, a_t) \\ & \leq \alpha(T) \cdot 1 + (1 - \alpha(T)) \exp\left(-\frac{LT(f(T) - |Q_T|)}{2^7 f(T)^2}\right) \\ & \leq \alpha(T) + (1 - \alpha(T)) \exp\left(-\frac{LTf(T)/2}{2^7 f(T)^2}\right) \\ & = \alpha(T) + (1 - \alpha(T)) \exp\left(-\frac{LT}{2^8 f(T)}\right) \end{aligned}$$

whenever $T \geq T_0$. Since $\prod_{t=1}^T \mu_{f, \mathbf{x}}^m(s_t) = 1$ always, we have¹⁴

$$\begin{aligned} \sup_{\mu} \mathbb{E}_{\mathbf{s} \sim U^T, \mathbf{a}} R_T(\mathbf{s}, \mathbf{a}, \mu, \pi^m) & \geq \mathbb{E}_{\mathbf{x} \sim U(\{0,1\}^{f(T)})} \mathbb{E}_{\mathbf{s} \sim U^T, \mathbf{a}} R_T(\mathbf{s}, \mathbf{a}, \mu_{f, \mathbf{x}}, \pi^m) \\ & = \mathbb{E}_{\mathbf{x} \sim U(\{0,1\}^{f(T)})} \mathbb{E}_{\mathbf{s} \sim U^T, \mathbf{a}} \left[\prod_{t=1}^T \mu_{f, \mathbf{x}}^m(s_t) - \prod_{t=1}^T \mu_{f, \mathbf{x}}(s_t, a_t) \right] \\ & \geq 1 - \alpha(T) - (1 - \alpha(T)) \exp\left(-\frac{LT}{2^8 f(T)}\right) \end{aligned}$$

Therefore

$$\begin{aligned} \lim_{T \rightarrow \infty} \sup_{\mu} \mathbb{E}_{\mathbf{s} \sim U^T, \mathbf{a}} R_T(\mathbf{s}, \mathbf{a}, \mu, \pi^m) & \geq 1 - \lim_{T \rightarrow \infty} \alpha(T) - (1 - \lim_{T \rightarrow \infty} \alpha(T)) \cdot \exp\left(\lim_{T \rightarrow \infty} -\frac{LT}{2^8 f(T)}\right) \\ & = 1 - 0 - (1 - 0) \cdot \exp(-\infty) \\ & = 1 \end{aligned}$$

as required. \square

B PROOF OF THEOREM 5.2

B.1 CONTEXT ON LEMMA 5.4

Before diving into the main proof, we provide some context on Lemma 5.4 from Section 5:

Lemma 5.4 (Lemma 3.5 in Russo et al. (2024)). *Assume $\tilde{\Pi}$ is finite. Then for any loss function $\ell : [T] \times \tilde{\Pi} \rightarrow [0, 1]$ and query probability p , HEDGWITHQUERIES enjoys the regret bound*

$$\sum_{t=1}^T \mathbb{E}[\ell(t, \pi_t)] - \min_{\tilde{\pi} \in \tilde{\Pi}} \sum_{t=1}^T \ell(t, \tilde{\pi}) \leq \frac{2 \log |\tilde{\Pi}|}{p^2}$$

where π_t is the policy chosen at time t and the expectation is over the randomness of the algorithm.

Lemma 5.4 is a restatement and simplification of Lemma 3.5 in Russo et al. (2024). First, Russo et al. (2024) parametrize their algorithm by the expected number of queries \hat{k} instead of the query probability $p = \hat{k}/T$. Second, Russo et al. (2024) include a second parameter k , which is the eventual target number of queries for their unconditional query bound. In our case, an expected query bound is

¹⁴Fubini's theorem means we need not worry about the order of the expectation operators.

Algorithm 2 A variant of the Hedge algorithm which only observes losses in response to queries.

```

1: function HEDGWITHQUERIES( $p \in (0, 1]$ , finite policy class  $\tilde{\Pi}$ , unknown  $\ell : [T] \times \tilde{\Pi} \rightarrow [0, 1]$ )
2:    $w(\pi) \leftarrow 1$  for all  $\pi \in \tilde{\Pi}$ 
3:    $\eta \leftarrow \max\left(\sqrt{\frac{p \log |\tilde{\Pi}|}{2T}}, \frac{p^2}{\sqrt{2}}\right)$ 
4:   for  $t$  from 1 to  $T$  do
5:     hedgeQuery  $\leftarrow$  true with probability  $p$  else false
6:     if hedgeQuery then
7:       Query and observe  $\ell(t, \pi)$  for all  $\pi \in \tilde{\Pi}$ 
8:        $\ell^* \leftarrow \min_{\pi \in \tilde{\Pi}} \ell(t, \pi)$ 
9:        $w(\pi) \leftarrow w(\pi) \cdot \exp(-\eta(\ell(t, \pi) - \ell^*))$  for all  $\pi \in \tilde{\Pi}$ 
10:      Select policy  $\arg \min_{\pi \in \tilde{\Pi}} \ell(t, \pi)$ 
11:     else
12:       $P(\pi) \leftarrow w(\pi) / \sum_{\pi' \in \tilde{\Pi}} w(\pi')$  for all  $\pi \in \tilde{\Pi}$ 
13:      Sample  $\pi_t \sim P$ 
14:      Select policy  $\pi_t$ 

```

sufficient, so we simply set $k = \hat{k}$. Third, Russo et al. (2024) provide a second bound which is tighter for small k ; that bound is less useful for us so we omit it. Fourth, their number of actions n is equal to $|\tilde{\Pi}|$ in our setting. (Their actions correspond to policies in $\tilde{\Pi}$, not our actions in A .) Since Russo et al. (2024) set $\eta = \max\left(\frac{1}{T} \sqrt{\frac{\hat{k} \log n}{2}}, \frac{k\hat{k}}{\sqrt{2}T^2}\right)$, we end up with $\eta = \max\left(\sqrt{\frac{p \log |\tilde{\Pi}|}{2T}}, \frac{p^2}{\sqrt{2}}\right)$. Algorithm 2 provides precise pseudocode for the HEDGWITHQUERIES algorithm to which Lemma 5.4 refers.

B.2 MAIN PROOF

We use the following notation throughout the proof:

- For each $t \in [T]$, let X_t refer to the value of X at the start of time step t .
- Let $V_T = \{t \in [T] : \pi_t(s_t) \neq \pi^m(s_t)\}$ be the set of time steps where Hedge’s proposed action doesn’t match the mentor’s. Note that $|V_T|$ upper bounds the number of mistakes the algorithm makes (the number of mistakes could be smaller, since the algorithm sometimes queries instead of taking action $\pi_t(s_t)$).
- For $K \subseteq S$, let $\text{vol}(K)$ denote the n -dimensional Lebesgue measure of K .
- With slight abuse of notation, we will use inequalities of the form $f(T) \leq g(T) + O(h(T))$ to mean that there exists a constant C such that $f(T) \leq g(T) + Ch(T)$.
- We will use “Case 1” to refer to finite VC dimension and σ -smooth s and “Case 2” to refer to finite Littlestone dimension. In Case 1, expectations are over the randomness of both s and the algorithm, while in Case 2, expectations are over just the randomness of the algorithm. When we need to distinguish, we use \mathbb{E}_a to denote the expectation over randomness of the algorithm and $\mathbb{E}_{s \sim \mathcal{D}}$ to denote the expectation over s .

Lemma B.1. *Under the conditions of Theorem 5.1, Algorithm 1 satisfies*

$$\mathbb{E}[|V_T|] \in O\left(\frac{d}{\sigma} T \varepsilon \log(1/\varepsilon) \log T\right)$$

Proof. Define $\ell : [T] \times \tilde{\Pi} \rightarrow [0, 1]$ by $\ell(t, \pi) = \mathbf{1}(\pi(s_t) \neq \pi^m(s_t))$, and let w^h and π_t^h denote the values of w and π_t respectively in HEDGWITHQUERIES, while w and π_t refer to the variables in Algorithm 1. Then w and w^h evolve in the exact same way, so the distributions of π_t and π_t^h coincide. Thus by Lemma 5.4,

$$\mathbb{E}_a \left[\sum_{t=1}^T \ell(t, \pi_t) \right] - \min_{\tilde{\pi} \in \tilde{\Pi}} \sum_{t=1}^T \ell(t, \tilde{\pi}) \leq \frac{2 \log |\tilde{\Pi}|}{p^2} = 2T\varepsilon \log |\tilde{\Pi}|$$

Since Lemma 5.4 holds for any loss function, the bound above holds for any $\mathbf{s} \in S^T$, so the bound also holds in expectation over $\mathbf{s} \sim \mathcal{D}$ (which is needed for Case 1). Next, observe that $|V_T| = \sum_{t=1}^T \mathbf{1}(\pi_t(s_t) \neq \pi^m(s_t)) = \sum_{t=1}^T \ell(t, \pi_t)$, so

$$\mathbb{E}_{\mathbf{a}}[|V_T|] \leq 2T\varepsilon \log |\tilde{\Pi}| + \min_{\tilde{\pi} \in \tilde{\Pi}} \sum_{t=1}^T \mathbf{1}(\tilde{\pi}(s_t) \neq \pi^m(s_t))$$

Case 1: Since $\tilde{\Pi}$ is a smooth ε -cover of Π , we have

$$\begin{aligned} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} \left[\min_{\tilde{\pi} \in \tilde{\Pi}} \sum_{t=1}^T \mathbf{1}(\tilde{\pi}(s_t) \neq \pi^m(s_t)) \right] &\leq \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} \left[\sup_{\pi \in \Pi} \min_{\tilde{\pi} \in \tilde{\Pi}} \sum_{t=1}^T \mathbf{1}(\tilde{\pi}(s_t) \neq \pi(s_t)) \right] \\ &\in O \left(\frac{1}{\sigma} T\varepsilon \log T \sqrt{d \log(1/\varepsilon)} \right) \end{aligned}$$

with the first step due to $\pi^m \in \Pi$ and the second step due to Lemma 5.3. The last component we need is that $|\tilde{\Pi}| \leq (41/\varepsilon)^d$ by construction (and such a $\tilde{\Pi}$ is guaranteed to exist by Lemma 5.1). Combining the above inequalities and taking the expectation over $\mathbf{s} \sim \mathcal{D}$ (in addition to the randomness of the algorithm), we get

$$\begin{aligned} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a}}[|V_T|] &\leq 2T\varepsilon \log |\tilde{\Pi}| + \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} \left[\min_{\tilde{\pi} \in \tilde{\Pi}} \sum_{t=1}^T \mathbf{1}(\tilde{\pi}(s_t) \neq \pi^m(s_t)) \right] \\ &\leq 2dT\varepsilon \log(41/\varepsilon) + O \left(\frac{1}{\sigma} T\varepsilon \log T \sqrt{d \log(1/\varepsilon)} \right) \\ &\in O \left(\frac{d}{\sigma} T\varepsilon \log(1/\varepsilon) \log T \right) \end{aligned}$$

Case 2: Since $\tilde{\Pi}$ is an adversarial cover of Π and $\pi^m \in \Pi$, there exists $\tilde{\pi} \in \tilde{\Pi}$ such that $\sum_{t=1}^T \mathbf{1}(\tilde{\pi}(s_t) \neq \pi^m(s_t)) = 0$. Since $|\tilde{\Pi}| \leq (eT/d)^d$ (with such a $\tilde{\Pi}$ guaranteed to exist by Lemma 5.2),

$$\begin{aligned} \mathbb{E}_{\mathbf{a}}[|V_T|] &\leq 2T\varepsilon \log |\tilde{\Pi}| + \min_{\tilde{\pi} \in \tilde{\Pi}} \sum_{t=1}^T \mathbf{1}(\tilde{\pi}(s_t) \neq \pi^m(s_t)) \\ &\leq 2T\varepsilon d \ln(eT/d) \\ &\in O \left(\frac{d}{\sigma} T\varepsilon \log(1/\varepsilon) \log T \right) \end{aligned}$$

as required. \square

Lemma B.2. For all $t \in [T]$, $\mu(s_t, a_t) \geq \mu^m(s_t) - L\varepsilon^{1/n}$.

Proof. Fix any $t \in [T]$. If $t \in Q_T$, then $\mu(s_t, a_t) = \mu^m(s_t)$, so assume $t \notin Q_T$. Let $(s', a') = \arg \min_{(s, a) \in X_t: \pi_t(s_t) = a} \|s_t - s\|$. Since $t \notin Q_T$, we must have $\|s_t - s'\| \leq \varepsilon^{1/n}$.

We have $a' = \pi^m(s')$ by construction of X_t and $\pi_t(s_t) = a'$ by construction of a' . Combining these with the local generalization assumption, we get

$$\mu(s_t, a_t) = \mu(s_t, \pi_t(s_t)) = \mu(s_t, \pi^m(s')) \geq \mu^m(s_t) - L\|s_t - s'\| \geq \mu^m(s_t) - L\varepsilon^{1/n}$$

as required. \square

Lemma B.3. Assume $x_1, \dots, x_T, y_1, \dots, y_T \in [0, 1]$ and $x_t \geq y_t$ for all $t \in [T]$. Then

$$\prod_{t=1}^T x_t - \prod_{t=1}^T y_t \leq \sum_{t=1}^T x_t - \sum_{t=1}^T y_t$$

Proof. We proceed by induction on T . The claim is trivially satisfied for $T = 1$, so suppose $T > 1$ and assume that $\prod_{t=1}^{T-1} x_t - \prod_{t=1}^{T-1} y_t \leq \sum_{t=1}^{T-1} x_t - \sum_{t=1}^{T-1} y_t$. Then

$$\begin{aligned} \sum_{t=1}^T x_t - \sum_{t=1}^T y_t - \prod_{t=1}^T x_t + \prod_{t=1}^T y_t &= x_T \sum_{t=1}^{T-1} x_t - y_T \sum_{t=1}^{T-1} y_t - x_T \prod_{t=1}^{T-1} x_t + y_T \prod_{t=1}^{T-1} y_t \\ &= x_T \left(\sum_{t=1}^{T-1} x_t - \prod_{t=1}^{T-1} x_t \right) - y_T \left(\sum_{t=1}^{T-1} y_t - \prod_{t=1}^{T-1} y_t \right) \end{aligned}$$

Since $T > 1$ and $x_t \in [0, 1]$ for all $t \in [T]$, we have $\sum_{t=1}^{T-1} x_t \geq x_1 \geq \sum_{t=1}^{T-1} x_t$. Thus $\sum_{t=1}^{T-1} x_t - \prod_{t=1}^{T-1} x_t \geq 0$. Combining this with $x_T \geq y_T$, we get

$$\begin{aligned} \sum_{t=1}^T x_t - \sum_{t=1}^T y_t - \prod_{t=1}^T x_t + \prod_{t=1}^T y_t &= x_T \left(\sum_{t=1}^{T-1} x_t - \prod_{t=1}^{T-1} x_t \right) - y_T \left(\sum_{t=1}^{T-1} y_t - \prod_{t=1}^{T-1} y_t \right) \\ &\geq y_T \left(\sum_{t=1}^{T-1} x_t - \prod_{t=1}^{T-1} x_t \right) - y_T \left(\sum_{t=1}^{T-1} y_t - \prod_{t=1}^{T-1} y_t \right) \\ &= y_T \left(\sum_{t=1}^{T-1} x_t - \prod_{t=1}^{T-1} x_t - \sum_{t=1}^{T-1} y_t + \prod_{t=1}^{T-1} y_t \right) \\ &\geq 0 \end{aligned}$$

The last step is due to $y_T \geq 0$ and our assumption of $\prod_{t=1}^{T-1} x_t - \prod_{t=1}^{T-1} y_t \leq \sum_{t=1}^{T-1} x_t - \sum_{t=1}^{T-1} y_t$. \square

Lemma B.4. Under the conditions of Theorem 5.1, Algorithm 1 satisfies

$$\mathbb{E}[R_T] \in O\left(\frac{dL}{\sigma} T \varepsilon^{1+1/n} \log(1/\varepsilon) \log T\right)$$

Proof. We first claim that $a_t = \pi^m(s_t)$ for all $t \notin V_T$. If $t \in Q_T$, the claim is immediate; if not, we have $a_t = \pi_t(s_t)$, and $\pi_t(s_t) = \pi^m(s_t)$ due to $t \notin V_T$. Thus $\min(\mu^m(s_t), \mu(s_t, a_t)) = \mu^m(s_t)$ for $t \notin V_T$.

We next claim that $\mu^m(s_t) - \min(\mu^m(s_t), \mu(s_t, a_t)) \leq L\varepsilon^{1/n}$ for all $t \in [T]$. If $\mu(s_t, a_t) \leq \mu^m(s_t)$, this follows from Lemma B.2. If $\mu(s_t, a_t) > \mu^m(s_t)$, then $\mu^m(s_t) - \min(\mu^m(s_t), \mu(s_t, a_t)) = 0 \leq L\varepsilon^{1/n}$. Therefore

$$\begin{aligned} \sum_{t=1}^T \left(\mu^m(s_t) - \min(\mu^m(s_t), \mu(s_t, a_t)) \right) &\leq \sum_{t \in V_T} \left(\mu^m(s_t) - \min(\mu^m(s_t), \mu(s_t, a_t)) \right) \\ &\leq \sum_{t \in V_T} L\varepsilon^{1/n} \\ &= |V_T| L\varepsilon^{1/n} \end{aligned}$$

Now let $x_t = \mu^m(s_t)$ and $y_t = \min(\mu^m(s_t), \mu(s_t, a_t))$ for all $t \in [T]$. Then by Lemma B.3,

$$\prod_{t=1}^T \mu^m(s_t) - \prod_{t=1}^T \min(\mu^m(s_t), \mu(s_t, a_t)) \leq \sum_{t=1}^T \left(\mu^m(s_t) - \min(\mu^m(s_t), \mu(s_t, a_t)) \right)$$

Since $\mu(s_t, a_t) \geq \min(\mu^m(s_t), \mu(s_t, a_t))$ for all $t \in [T]$, we have

$$\begin{aligned} R_T &= \prod_{t=1}^T \mu^m(s_t) - \prod_{t=1}^T \mu(s_t, a_t) \\ &\leq \sum_{t=1}^T \left(\mu^m(s_t) - \min(\mu^m(s_t), \mu(s_t, a_t)) \right) \\ &\leq |V_T| L\varepsilon^{1/n} \end{aligned}$$

Taking the expectation and applying Lemma B.1 completes the proof. \square

Definition B.1. Let $(K, \|\cdot\|)$ be a normed vector space and let $\delta > 0$. Then $W \subseteq K$ is a δ -packing of K if for all $x, y \in W$, $\|x - y\| > \delta$. The δ -packing number of K , denoted $\mathcal{M}(K, \|\cdot\|, \delta)$, is the maximum cardinality of any δ -packing of K .

In this paper, we only consider the Euclidean distance norm, so we just write $M(K, \|\cdot\|, \delta) = M(K, \delta)$.

Lemma B.5 (Theorem 14.2 in Wu (2020)). *If $K \subset \mathbb{R}^n$ is convex, bounded, and contains a ball with radius $\delta > 0$, then*

$$\mathcal{M}(K, \delta) \leq \frac{3^n \text{vol}(K)}{\delta^n \text{vol}(B)}$$

where B is a unit ball.

Lemma B.6 (Jung’s Theorem (Jung, 1901)). *If $K \subset \mathbb{R}^n$ is compact, then there exists a closed ball with radius at most $\text{diam}(K) \sqrt{\frac{n}{2(n+1)}}$ containing K .*

Lemma B.7. *Under the conditions of Theorem 5.1, Algorithm 1 satisfies*

$$\mathbb{E}[|Q_T|] \in O\left(\sqrt{\frac{T}{\varepsilon}} + \frac{d}{\sigma} T \varepsilon \log(1/\varepsilon) \log T + \frac{\text{diam}(\mathbf{s})^n}{\varepsilon}\right)$$

Proof. If $t \in Q_T$, then either `hedgeQuery` = true or $\min_{(s,a) \in X_t: \pi_t(s_t)=a} |s_t - s| > r$. The expected number of time steps with `hedgeQuery` = true is $pT = \sqrt{T/\varepsilon}$, so let $\hat{S} = \{s_t : t \in Q_T \text{ and } \min_{(s,a) \in X_t: \pi_t(s_t)=a} |s_t - s| > r\}$. We further subdivide \hat{S} into $\hat{S}_1 = \{s_t \in \hat{S} : \pi_t(s_t) \neq \pi^m(s_t)\}$ and $\hat{S}_2 = \{s_t \in \hat{S} : \pi_t(s_t) = \pi^m(s_t)\}$. Since $\hat{S}_1 \subseteq V_T$, Lemma B.1 implies that $\mathbb{E}[|\hat{S}_1|] \in O\left(\frac{d}{\sigma} T \varepsilon \log(1/\varepsilon) \log T\right)$.

Next, fix an $a \in A$ and let $S_a = \{s \in \mathbf{s} : \pi^m(s) = a\}$ be the set of observed states which share a mentor action. We claim that $\hat{S}_2 \cap S_a$ is a packing of S_a . Suppose instead that there exists $s, s' \in \hat{S}_2 \cap S_a$, with $\|s - s'\| \leq \varepsilon^{1/n}$. WLOG assume s was queried after s' and let t be the time step on which s was queried. Then $(s', \pi^m(s')) \in X_t$. Also, $s, s' \in \hat{S}_2 \cap S_a$ implies that and $\pi_t(s_t) = \pi^m(s_t) = a = \pi^m(s')$. Therefore

$$\min_{(s'', a'') \in X_t: a'' = \pi_t(s_t)} \|s_t - s''\| \leq \|s_t - s'\| \leq \varepsilon^{1/n}$$

which contradicts $s_t \in \hat{S}$. Thus $\hat{S}_2 \cap S_a$ is a $\varepsilon^{1/n}$ -packing of S_a .

By Lemma B.6, there exists a ball B_1 of diameter $\text{diam}(\mathbf{s}) \sqrt{\frac{n}{2(n+1)}}$ which contains \mathbf{s} . Let $R = \text{diam}(\mathbf{s}) \sqrt{\frac{n}{8(n+1)}}$ denote the radius of B_1 . Let B_2 be the ball with the same center as B_1 but with radius $\max(R, \varepsilon^{1/n})$. Since $S_a \subset \mathbf{s} \subset B_1 \subset B_2$, $\hat{S}_2 \cap S_a$ is also a $\varepsilon^{1/n}$ -packing of B_2 . Also, B_2 must contain a ball of radius $\varepsilon^{1/n}$, so Lemma B.5 implies that

$$\begin{aligned} |\hat{S}_2 \cap S_a| &\leq \mathcal{M}(B_2, \varepsilon^{1/n}) \\ &\leq \frac{3^n \text{vol}(B_2)}{\varepsilon \text{vol}(B)} \\ &= (\max(R, \varepsilon^{1/n}))^n \frac{3^n \text{vol}(B)}{\varepsilon \text{vol}(B)} \\ &= \max\left(\text{diam}(\mathbf{s})^n \left(\frac{n}{8(n+1)}\right)^{n/2}, \varepsilon\right) \frac{3^n}{\varepsilon} \\ &\leq O\left(\frac{\text{diam}(\mathbf{s})^n}{\varepsilon} + 1\right) \end{aligned}$$

(The +1 is necessary for now since $\text{diam}(\mathbf{s})$ could theoretically be zero.) Therefore

$$\mathbb{E}[|Q_T|] = \sqrt{\frac{T}{\varepsilon}} + \mathbb{E}[|\hat{S}|]$$

$$\begin{aligned}
&= \sqrt{\frac{T}{\varepsilon}} + \mathbb{E}[|\hat{S}_1|] + \mathbb{E}\left[\sum_{a \in A} |\hat{S}_2 \cap S_a|\right] \\
&\leq \sqrt{\frac{T}{\varepsilon}} + O\left(\frac{d}{\sigma} T \varepsilon \log(1/\varepsilon) \log T\right) + \sum_{a \in A} O\left(\frac{\text{diam}(\mathbf{s})^n}{\varepsilon} + 1\right) \\
&\leq \sqrt{\frac{T}{\varepsilon}} + O\left(\frac{d}{\sigma} T \varepsilon \log(1/\varepsilon) \log T\right) + |A| \cdot O\left(\frac{\text{diam}(\mathbf{s})^n}{\varepsilon} + 1\right) \\
&\leq O\left(\sqrt{\frac{T}{\varepsilon}} + \frac{d}{\sigma} T \varepsilon \log(1/\varepsilon) \log T + \frac{\text{diam}(\mathbf{s})^n}{\varepsilon}\right)
\end{aligned}$$

as required. \square

Theorem 5.1 follows from Lemmas B.4 and B.7:

Theorem 5.1. *Let $A = \{0, 1\}$. Assume $\pi^m \in \Pi$ where either (1) Π has finite VC dimension d , \mathbf{s} is σ -smooth, and $\varepsilon T \log T > 12\sigma d \log(4e^2/\varepsilon)$ or (2) Π has finite Littlestone dimension d . Then for any $T \in \mathbb{N}$ and $\varepsilon > 0$, Algorithm 1 satisfies*

$$\begin{aligned}
\mathbb{E}[R_T] &\in O\left(\frac{dL}{\sigma} T \varepsilon^{1+1/n} \log(1/\varepsilon) \log T\right) \\
\mathbb{E}[|Q_T|] &\in O\left(\sqrt{\frac{T}{\varepsilon}} + \frac{d}{\sigma} T \varepsilon \log(1/\varepsilon) \log T + \frac{\text{diam}(\mathbf{s})^n}{\varepsilon}\right)
\end{aligned}$$

We then perform some arithmetic to get Theorem 5.2:

Theorem 5.2. *Let $A = \{0, 1\}$. Assume $\pi^m \in \Pi$ where either (1) Π has finite VC dimension d and \mathbf{s} is σ -smooth or (2) Π has finite Littlestone dimension d . Then for any $T \in \mathbb{N}$, Algorithm 1 with $\varepsilon = T^{\frac{-2n}{2n+1}}$ satisfies*

$$\begin{aligned}
\mathbb{E}[R_T] &\in O\left(\frac{dL}{\sigma} T^{\frac{-1}{2n+1}} \log T\right) \\
\mathbb{E}[|Q_T|] &\in O\left(T^{\frac{4n+1}{4n+2}} \left(\frac{d}{\sigma} \log T + \text{diam}(\mathbf{s})^n\right)\right)
\end{aligned}$$

Proof. We have

$$\begin{aligned}
\mathbb{E}[R_T] &\in O\left(\frac{dL}{\sigma} T^{1-\frac{2n}{2n+1}-\frac{2}{2n+1}} \log(1/\varepsilon) \log T\right) \\
&= O\left(\frac{dL}{\sigma} T^{\frac{-1}{2n+1}} \log T\right)
\end{aligned}$$

and

$$\begin{aligned}
\mathbb{E}[|Q_T|] &\in O\left(\sqrt{T^{1+\frac{2n}{2n+1}}} + \frac{d}{\sigma} T^{1-\frac{2n}{2n+1}} \log(T^{\frac{2n}{2n+1}}) \log T + T^{\frac{2n}{2n+1}} \text{diam}(\mathbf{s})^n\right) \\
&= O\left(T^{\frac{2n+0.5}{2n+1}} + \frac{d}{\sigma} T^{\frac{1}{2n+1}} \log T + T^{\frac{2n}{2n+1}} \text{diam}(\mathbf{s})^n\right) \\
&\leq O\left(T^{\frac{4n+1}{4n+2}} \left(\frac{d}{\sigma} \log T + \text{diam}(\mathbf{s})^n\right)\right)
\end{aligned}$$

\square

Corollary 5.2.1. *If Algorithm 1 is modified to stop querying after $T^{\frac{4n+1.5}{4n+2}} (\frac{d}{\sigma} \log T + \text{diam}(\mathbf{s})^n)$ queries, then under the conditions of Theorem 5.2, we have $\mathbb{E}[R_T] \in O(\frac{dL}{\sigma} T^{\frac{-1}{8n+4}} \log T)$.*

Proof. Let ξ be the event that Algorithm 1 exceeds $T^{\frac{4n+1.5}{4n+2}} (\frac{d}{\sigma} \log T + \text{diam}(\mathbf{s})^n)$ queries. Noting that $|Q_T| \leq T$ always, we have

$$\begin{aligned} \mathbb{E}[|Q_T|] &= \Pr[\xi] \mathbb{E}[|Q_T| \mid \xi] + \Pr[\neg\xi] \mathbb{E}[|Q_T| \mid \neg\xi] \\ &\geq \Pr[\xi] T^{\frac{4n+1.5}{4n+2}} \left(\frac{d}{\sigma} \log T + \text{diam}(\mathbf{s})^n \right) \end{aligned}$$

Therefore by Theorem 5.2,

$$\Pr[\xi] \leq \frac{O\left(T^{\frac{4n+1}{4n+2}} \left(\frac{d}{\sigma} \log T + \text{diam}(\mathbf{s})^n\right)\right)}{T^{\frac{4n+1.5}{4n+2}} \left(\frac{d}{\sigma} \log T + \text{diam}(\mathbf{s})^n\right)} \leq O\left(T^{\frac{-0.5}{4n+2}}\right)$$

If ξ does not occur, then the regret bound from Theorem 5.2 holds. If not, we still have $R_T \leq 1$ unconditionally. Thus the regret of the modified version of Algorithm 1 satisfies

$$\begin{aligned} \mathbb{E}[R_T] &\leq \Pr[\xi] \mathbb{E}[R_T \mid \xi] + \Pr[\neg\xi] \mathbb{E}[R_T \mid \neg\xi] \\ &\leq O\left(T^{\frac{-0.5}{4n+2}}\right) \cdot 1 + 1 \cdot O\left(\frac{dL}{\sigma} T^{\frac{-1}{2n+1}} \log T\right) \\ &\leq O\left(\frac{dL}{\sigma} T^{\frac{-1}{8n+4}} \log T\right) \end{aligned}$$

□

B.3 ADAPTIVE ADVERSARIES

If s_t is allowed to depend on the events of prior time steps, we say that the adversary is adaptive. In contrast, a non-adaptive or “oblivious” adversary must choose the entire input upfront. This distinction is not relevant for deterministic algorithms, since an adversary knows exactly how the algorithm will behave for any input. In other words, the adversary gains no new information during the execution of the algorithm. For randomized algorithms, an adaptive adversary can base the choice of s_t on the results of randomization on previous time steps (but not on the current time step), while an oblivious adversary cannot.

In the standard online learning model, Hedge guarantees sublinear regret against both oblivious and adaptive adversaries (Chapter 5 of [Slivkins et al. \(2019\)](#) or Chapter 21 of [Shalev-Shwartz & Ben-David \(2014\)](#)). However, [Russo et al. \(2024\)](#) state their result only for oblivious adversaries. In order for our overall proof of Theorem 5.1 to hold for adaptive adversaries, Lemma 5.4 (Lemma 3.5 in [Russo et al. \(2024\)](#)) must also hold for adaptive adversaries. In this section, we argue why the proof of Lemma 5.4 (Lemma 3.5 in their paper) goes through for adaptive adversaries as well. For this rest of Appendix B.3, lemma numbers refer to the numbering in [Russo et al. \(2024\)](#).

The importance of independent queries. Recall from Appendix B.1 that [Russo et al. \(2024\)](#) allow two separate parameters k and \hat{k} , which we unify for simplicity. Recall also that Lemma 3.5 refers to the variant of Hedge which queries with probability $p = \hat{k}/T = k/T$ independently on each time step (Algorithm 2). More precisely, on each time step t , the algorithm samples a Bernoulli random variable $X_t \sim \text{Ber}(p)$ and queries if $X_t = 1$. The key idea is that X_t is independent of events on previous time steps. Thus even conditioning on the history up to time t , for any for any random variable Y_t we can write

$$\mathbb{E}[Y_t] = (1 - p) \mathbb{E}[Y_t \mid X_t = 0] + p \mathbb{E}[Y_t \mid X_t = 1]$$

This insight immediately extends Observation 3.3 to adaptive adversaries (with the minor modification that queries are now issued independently with probability p on each time step instead of issuing k uniformly distributed queries). Specifically, using the notation from [Russo et al. \(2024\)](#) where i_t is the action chosen at time t , i_t^0 is the action chosen at time t if a query is not issued, and i_t^* is the optimal action at time t , we have

$$\mathbb{E}[\ell_t(i_t)] = (1 - p) \mathbb{E}[\ell_t(i_t^0)] + p \mathbb{E}[\ell_t(i_t^*)] = \left(1 - \frac{k}{T}\right) \mathbb{E}[\ell_t(i_t^0)] + \frac{k}{T} \mathbb{E}[\ell_t(i_t^*)]$$

The same logic applies to other statements like $\mathbb{E}[\hat{\ell}_t(i) \mid X_{\leq t-1}, I_{\leq t-1}] = \ell_t(i) - \ell_t(i_t^*)$ and immediately extends those statements to adaptive adversaries as well.

Applying Observation 3.3. The other tricky part of the proof is applying Observation 3.3 using a new loss function $\hat{\ell}$ defined by $\hat{\ell}_t = \frac{T}{k}(\ell_t(i) - \ell_t(i_t^*))\mathbf{1}(X_t = 1)$. To do so, we must argue that standard Hedge run on $\hat{\ell}$ is the “counterpart without queries” of HEDGEWITHQUERIES. Specifically, both algorithms must have the same weight vectors on every time step, and the only difference should be that HEDGEWITHQUERIES takes the optimal action on each time step independently with probability p (and otherwise behaves the same as standard Hedge). On time steps with $X_t = 0$, standard Hedge observes $\hat{\ell}_t(i) = 0$ for all actions i and thus makes no updates, and HEDGEWITHQUERIES makes no updates by definition. On time steps with $X_t = 1$, both algorithms perform the typical updates $w_{t+1}(i) = w_t(i) \cdot \exp(-\eta(\hat{\ell}_t(i) - \hat{\ell}_t(i_t^*)))$. Thus the weight vectors are the same for both algorithms on every time step. Furthermore, HEDGEWITHQUERIES takes the optimal action at time t iff $X_t = 1$, which occurs independently with probability p on each time step. Thus standard Hedge run on $\hat{\ell}$ is the “counterpart without queries” of HEDGEWITHQUERIES.

The rest of the proof. The other elements of the proof of Lemma 3.5 are as follows:

1. Lemma 3.1, which analyzes the standard version of Hedge (i.e., no queries and losses are observed on every time step).
2. Applying Lemma 3.1 to a $\hat{\ell}$.
3. Arithmetic and rearranging terms.

The proof of Lemma 3.1 relies on simple arithmetic properties of the Hedge weights. Regardless of the adversary’s behavior, $\hat{\ell}$ is a well-defined loss function, so Lemma 3.1 can be applied. Step 3 clearly has no dependence on the type of adversary. Thus we conclude that Lemma 3.5 extends to adaptive adversaries.

C GENERALIZING THEOREM 5.2 TO MANY ACTIONS

We use the standard “one versus rest” reduction (see, e.g., Chapter 29 of [Shalev-Shwartz & Ben-David \(2014\)](#)). For each action a , we will learn a binary classifier which predicts whether action a is the mentor’s action. Formally, for each $a \in A$, define the policy class $\Pi_a = \{\pi_a : \pi \in \Pi \text{ and } \pi_a(s) = \mathbf{1}(\pi(s) = a) \mid \forall s \in S\}$. Informally, for each policy $\pi : S \rightarrow A$ in Π , there exists a policy $\pi_a : S \rightarrow \{0, 1\}$ in Π_a such that $\pi_a(s) = \mathbf{1}(\pi(s) = a)$ for all $s \in S$.

Algorithm 3 runs one copy of our binary-action algorithm Algorithm 1 for each action $a \in A$. At each time step t , the copy for action a returns an action b_t^a , with $b_t^a = 1$ indicating a belief that $a = \pi^m(s_t)$ and $b_t^a = 0$ indicating a belief that $a \neq \pi^m(s_t)$. (Note that $b_t^a = \hat{a}$ is also possible, indicating that the mentor was queried.)

The key idea is that if b_t^a is correct for each action a , there will be exactly one a such that $b_t^a = 1$, and specifically it will be $a = \pi^m(s_t)$. Thus we are guaranteed to take the mentor’s action on such time steps. The analysis for Theorem 5.2 (specifically, Lemma B.1) bounds the number of time steps when a given copy of Algorithm 1 is incorrect, so by the union bound, the number of time steps where *any* copy is incorrect is $|A|$ times that bound. That in turn bounds the number of time steps where Algorithm 3 takes an action other than the mentor’s. Similarly, the number of queries made by Algorithm 3 is at most $|A|$ times the bound from Theorem 5.2. The result is the following theorem:

Theorem C.1. Assume $\pi^m \in \Pi$ where either (1) Π_a has finite VC dimension d and s is σ -smooth or (2) Π_a has finite Littlestone dimension d for all $a \in A$. Then for any $T \in \mathbb{N}$, Algorithm 3 with T and $\varepsilon = T^{\frac{-2n}{2n+1}}$ satisfies

$$\begin{aligned} \mathbb{E}[R_T] &\in O\left(\frac{|A|dL}{\sigma} T^{\frac{-1}{2n+1}} \log T\right) \\ \mathbb{E}[|Q_T|] &\in O\left(|A| T^{\frac{4n+1}{4n+2}} \left(\frac{d}{\sigma} \log T + \text{diam}(s)^n\right)\right) \end{aligned}$$

We use the following terminology and notation in the proof of Theorem C.1:

1. We refer to the copy of Algorithm 1 running on Π_a as “copy a of Algorithm 1”.

Algorithm 3 extends Algorithm 1 to many actions.

```

1: function AVOIDCATASTROPHEMANYACTIONS( $T \in \mathbb{N}$ ,  $\varepsilon \in \mathbb{R}_{>0}$ ,  $d \in \mathbb{N}$ , policy class  $\Pi$ )
2:   for  $a \in A$  do
3:     if  $\Pi$  has VC dimension  $d$  then
4:        $\tilde{\Pi}_a \leftarrow$  any smooth  $\varepsilon$ -cover of  $\Pi$  of size at most  $(41/\varepsilon)^d$ 
5:     else if  $\Pi$  has Littlestone dimension  $d$  then
6:        $\tilde{\Pi}_a \leftarrow$  any adversarial  $\varepsilon$ -cover of size at most  $(eT/d)^d$ 
7:     for  $t$  from 1 to  $T$  do
8:       for  $a \in A$  do
9:          $b_t^a \leftarrow$  action from running one step of Algorithm 1 on  $\Pi_a$  (with the same  $T, \varepsilon, d$ )
10:      if  $b_t^a \neq \hat{a} \forall a \in A$  and  $\exists a \in A : b_t^a = 1$  then
11:        Take any action  $a$  with  $b_t^a = 1$ 
12:      else
13:        Query the mentor

```

2. Let π_t^a and X_t^a refer to the values of π_t and X_t for copy a of Algorithm 1.
3. Let $\pi^{ma} : S \rightarrow \{0, 1\}$ be the policy defined by $\pi^{ma}(s) = 1(\pi^m(s_t) = a)$. Note that querying the mentor tells the agent $\pi^m(s_t)$, which allows the agent to compute $\pi^{ma}(s_t)$: this is necessary when Algorithm 1 queries while running on some Π_a .
4. Let $V_T^a = \{t \in [T] : b_t^a \neq \pi^{ma}(s_t)\}$ be the set of time steps where π_t^a does not correctly determine whether the mentor would take action a and let $V_T = \{t \in [T] : a_t \neq \pi^m(s_t)\}$ be the set of time steps where the agent's action doesn't match the mentor's.

Lemma C.1. We have $|V_T| \leq \sum_{a \in A} |V_T^a|$.

Proof. We claim that $V_T \subseteq \cup_{a \in A} V_T^a$. Suppose the opposite: then there exists $t \in V_T$ such that $b_t^a = \pi^{ma}(s_t)$ for all $a \in A$. Since $\pi^m(s_t) \in A$, there is exactly one $a \in A$ such that $1(\pi^m(s_t) = a) = \pi^{ma}(s_t) = b_t^a = 1$. Specifically, this holds for $a = \pi^m(s_t)$. But then Algorithm 3 takes action $\min\{a \in A : b_t^a = 1\} = \pi^m(s_t)$, which contradicts $t \in V_T$. Therefore $V_T \subseteq \cup_{a \in A} V_T^a$, and applying the union bound completes the proof. \square

Lemma C.2. For all $t \in [T]$, $\mu^m(s_t) - \mu(s_t, a_t) \leq L\varepsilon^{1/n}$.

Proof. The argument is similar to the proof of Lemma B.2. If $\mu^m(s_t) \neq \mu(s_t, a_t)$, then $a_t = a$ for some $a \in A$ where $b_t^a = 1$. Therefore copy a of Algorithm 1 did not query at time t and $\pi_t^a(s_t) = 1$. Let $(s', a') = \arg \min_{(s, a) \in X_t^a : \pi_t^a(s_t) = a} \|s_t - s\|$. Then $\|s_t - s'\| \leq \varepsilon^{1/n}$ and $a' = \pi_t^a(s_t) = 1$.

By construction of X_t^a , $a' = \pi^{ma}(s')$ so $\pi^{ma}(s') = 1$ which implies $\pi^m(s') = a$. Then by the local generalization assumption,

$$\mu(s_t, a_t) = \mu(s_t, a) = \mu(s_t, \pi^m(s')) \geq \mu^m(s_t) - L\|s_t - s'\| \geq \mu^m(s_t) - L\varepsilon^{1/n}$$

as required. \square

Theorem C.1. Assume $\pi^m \in \Pi$ where either (1) Π_a has finite VC dimension d and \mathbf{s} is σ -smooth or (2) Π_a has finite Littlestone dimension d for all $a \in A$. Then for any $T \in \mathbb{N}$, Algorithm 3 with T and $\varepsilon = T^{\frac{2n}{2n+1}}$ satisfies

$$\begin{aligned} \mathbb{E}[R_T] &\in O\left(\frac{|A|dL}{\sigma} T^{\frac{-1}{2n+1}} \log T\right) \\ \mathbb{E}[|Q_T|] &\in O\left(|A|T^{\frac{4n+1}{4n+2}} \left(\frac{d}{\sigma} \log T + \text{diam}(\mathbf{s})^n\right)\right) \end{aligned}$$

Proof. Theorem 5.2 implies that each copy of Algorithm 1 makes $O(T^{\frac{4n+1}{4n+2}} (\frac{d}{\sigma} \log T + \text{diam}(\mathbf{s})^n))$ queries in expectation, so by linearity of expectation, the expected number of queries made by

Algorithm 3 is $O(|A|T^{\frac{4n+1}{4n+2}}(\frac{d}{\sigma}\log T + \text{diam}(\mathbf{s})^n))$. Using the same argument as in the proof of Lemma B.7 (with Lemma C.2 replacing Lemma B.2), we get

$$\begin{aligned} R_T &= \prod_{t=1}^T \mu^m(s_t) - \prod_{t=1}^T \mu(s_t, a_t) \\ &\leq \sum_{t=1}^T \left(\mu^m(s_t) - \min(\mu^m(s_t), \mu(s_t, a_t)) \right) \\ &\leq |V_T| L \varepsilon^{1/n} \end{aligned}$$

Then by Lemma C.1, $R_T \leq L \varepsilon^{1/n} \sum_{a \in A} |V_T^a|$. Taking the expectation and applying Lemma B.1 to each V_T^a gives us

$$\mathbb{E}[R_T] \leq L \varepsilon^{1/n} \sum_{a \in A} O\left(\frac{d}{\sigma} T \varepsilon \log(1/\varepsilon) \log T\right) = O\left(|A| L \varepsilon^{1/n} \frac{d}{\sigma} T \varepsilon \log(1/\varepsilon) \log T\right)$$

as required. \square

D THERE EXIST POLICY CLASSES WHICH ARE LEARNABLE IN OUR SETTING BUT NOT IN THE STANDARD ONLINE MODEL

This section presents another algorithm with subconstant regret and sublinear queries, but under different assumptions. The most important aspect of our algorithm is that it can handle the class of thresholds on $[0, 1]$, which is known to have infinite Littlestone dimension and thus be hard in the standard online learning model. (Example 21.4 in Shalev-Shwartz & Ben-David (2014)).

Specifically, we assume a 1D state space and we allow the state sequence to be fully adversarial chosen. Instead of VC/Littlestone dimension, we consider the following notion of simplicity:

Definition D.1. Given a mentor policy π^m , partition the state space S into intervals such that all states within each interval share the same mentor action. Let $\{S_1, \dots, S_k\}$ be a partition that minimizes the number of intervals. We call each S_j a *segment*. Let $f(\pi^m)$ denote the number of segments in π^m .

Bounding the number of segments is similar conceptually to VC dimension in that it limits the ability of the policy class to realize arbitrary combinations of labels (i.e., mentor actions) on \mathbf{s} . For example, if Π is the class of thresholds on $[0, 1]$, every $\pi \in \Pi$ has at most two segments, and thus the positive result in this section will apply. This demonstrates the existence of policy classes which are learnable in our setting but not learnable in the standard online learning model, meaning that the two settings do not exactly coincide.

We prove the following result:

Theorem D.2. For any $\mathbf{s} \in S^T$, any π^m with $f(\pi^m) \leq K$, and any function $g : \mathbb{N} \rightarrow \mathbb{N}$, Algorithm 4 makes at most $(\text{diam}(\mathbf{s}) + 4)g(T)$ queries and satisfies $R_T \leq \frac{2LK T}{g(T)^2}$.

Choosing $g(T) = T^c$ for $c \in (1/2, 1)$ is sufficient to subconstant regret and sublinear queries:

Theorem D.3. For any $c \in (1/2, 1)$, Algorithm 4 with $g(T) = T^c$ makes $O(T^c(\text{diam}(\mathbf{s}) + 1))$ queries and satisfies

$$\lim_{T \rightarrow \infty} \sup_{\mathbf{s} \in S^T} \sup_{\mu} \sup_{\pi^m : f(\pi^m) \leq K} R_T = 0$$

Our algorithm does not need to know L or the number of segments; it only needs to know T .

D.1 INTUITION BEHIND THE ALGORITHM

The algorithm maintains a set of buckets which partition the observed portion of the state space. Each bucket's length determines the maximum loss in payoff we will allow from that subset of the state space. As long as the bucket contains a query from a prior time step, local generalization allows us to

Algorithm 4 achieves subconstant regret when the mentor’s policy has a bounded number of segments.

```

1: function AVOIDCATASTROPHE( $T \in \mathbb{N}$ ,  $g : \mathbb{N} \rightarrow \mathbb{N}$ )
2:    $S_Q \leftarrow \emptyset$  ▷ Previously queried states
3:    $\pi \leftarrow \emptyset$  ▷ Records  $\pi^m(s)$  for each  $s \in S_Q$ 
4:    $\mathcal{B} \leftarrow \emptyset$  ▷ The set of active buckets
5:   for  $t$  from 1 to  $T$  do
6:     EVALUATESTATE( $s_t$ )
7:   function EVALUATESTATE( $s \in S$ )
8:     if  $s \notin B$  for all  $B \in \mathcal{B}$  then ▷ No bucket containing  $s$ : create a new bucket and try again
9:        $B \leftarrow \left[ \frac{j-1}{g(T)}, \frac{j}{g(T)} \right]$  for  $j \in \mathbb{Z}$  such that  $s \in B$ 
10:       $\mathcal{B} \leftarrow \mathcal{B} \cup \{B\}$ 
11:       $n_B \leftarrow 0$  ▷ Number of time steps that have used  $B$ 
12:      EVALUATESTATE( $s$ )
13:   else
14:      $B \leftarrow$  any bucket containing  $s$ 
15:     if  $S_Q \cap B = \emptyset$  then ▷ No queries in this bucket
16:       Query mentor and observe  $\pi^m(s)$ 
17:        $\pi(s) \leftarrow \pi^m(s)$ 
18:        $S_Q \leftarrow S_Q \cup \{s\}$ 
19:        $n_B \leftarrow n_B + 1$ 
20:     else if  $n_B < T/g(T)$  then ▷ Bucket has a query and isn’t full: take that action
21:       Let  $s' \in S_Q \cap B$ 
22:       Take action  $\pi(s')$ 
23:        $n_B \leftarrow n_B + 1$ 
24:     else ▷ Bucket is full: split bucket and try again
25:        $B = [x, y]$ 
26:        $(B_1, B_2) \leftarrow \left( \left[ x, \frac{x+y}{2} \right], \left[ \frac{x+y}{2}, y \right] \right)$ 
27:        $(s_{B_1}, s_{B_2}) \leftarrow (0, 0)$ 
28:        $\mathcal{B} \leftarrow \mathcal{B} \cup \{B_1, B_2\} \setminus B$ 
29:       EVALUATESTATE( $s$ )

```

bound $\mu^m(s_t) - \mu(s_t, a_t)$ based on the length of the bucket containing s_t . We always query if the bucket does not contain a prior query

The granularity of the buckets is controlled by a function g , with the initial buckets having length $1/g(T)$. Since we can expect one query per bucket, we need $g(T) \in o(T)$ to ensure sublinear queries.

Regardless of the bucket length, the adversary can still place multiple segments in the same bucket B . A single query only tells us the optimal action for one of those segments, so we risk a payoff as bad as $\mu^m(s_t) - O(\text{len}(B))$ whenever we choose not to query. We can endure a limited number of such payoffs, but if we never query again in that bucket, we may suffer $\Theta(T)$ such payoffs. Letting $\mu^m(s_t) = 1$ for simplicity, that would lead to $\prod_{t=1}^T \mu(s_t, a_t) \leq \left(1 - \frac{1}{O(g(T))}\right)^{\Theta(T)}$, which converges to 0 (i.e., guaranteed catastrophe) when $g(T) \in o(T)$.

This failure mode suggests a natural countermeasure: if we start to suffer significant (potential) losses in the same bucket, then we should probably query there again. One way to structure these supplementary queries is by splitting the bucket in half when enough time steps have involved that bucket. It turns out that splitting after $T/g(T)$ time steps is a sweet spot.

D.2 NOTATION FOR THE PROOF

We will use the following notation throughout the proof of Theorem D.2:

- Let $V_T = \{t \in [T] : \mu(s_t, a_t) < \mu^m(s_t)\}$ be the set of time steps with a suboptimal payoff.
- Let B_t be the bucket that is used on time step t (as defined on line 14 of Algorithm 4).
- Let $d(B)$ be the *depth* of bucket B

- Buckets created on line 9 are depth 0.
- We refer to B_1, B_2 created on line 26 as the children of the bucket B defined on line 14.
- If B' is the child of B , $d(B') = d(B) + 1$.
- Note that $\text{len}(B) = \frac{1}{g(T)2^{d(B)}}$.
- Viewing the set of buckets as a binary tree defined by the “child” relation, we use the terms “ancestor” and “descendant” in accordance with their standard tree definitions.
- Let $\mathcal{B}_V = \{B : \exists t \in V_T \text{ s.t. } B_t = B\}$ be the set of buckets that ever produced a suboptimal payoff.
- Let $\mathcal{B}'_V = \{B \in \mathcal{B}_V : \text{no descendant of } B \text{ is in } \mathcal{B}_V\}$.

D.3 PROOF ROADMAP

The proof proceeds in the following steps:

1. Bound the total number of buckets and therefore the total number of queries (Lemma D.1).
2. Bound the suboptimality on a single time step based on the bucket length and L (Lemma D.2).
3. Bound the sum of bucket lengths on time steps where we make a mistake (Lemma D.4), with Lemma D.3 as an intermediate step. This captures the total amount of suboptimality.
4. As in the proof of Theorem 5.2, Lemma B.3 transforms the multiplicative objective into an additive form. Lemma D.5 bounds the additive objective using Lemmas D.2 and D.4.
5. Combining Lemmas D.5 and B.3 bounds the regret (Lemma D.6).
6. Theorem D.2 directly follows from Lemmas D.1 and D.6.

D.4 PROOF

Lemma D.1. *Algorithm 4 performs at most $(\text{diam}(\mathbf{s}) + 4)g(T)$ queries.*

Proof. Algorithm 4 performs at most one query per bucket, so the total number of queries is bounded by the total number of buckets. There are two ways to create a bucket: from scratch (line 9), or by splitting an existing bucket (line 26).

Since depth 0 buckets overlap only at their boundaries, and each depth 0 bucket has length $1/g(T)$, at most $g(T) \max_{t, t' \in [T]} |s_t - s_{t'}| = g(T) \text{diam}(\mathbf{s})$ depth 0 buckets are subsets of the interval $[\min_{t \in [T]} s_t, \max_{t \in [T]} s_t]$. At most two depth 0 buckets are not subsets of that interval (one at each end), so the total number of depth 0 buckets is at most $g(T) \text{diam}(\mathbf{s}) + 2$.

We split a bucket B when n_B reaches $T/g(T)$, which creates two new buckets. Since each time step increments n_B for a single bucket B , and there are a total of T time steps, the total number of buckets created via splitting is at most

$$\frac{2T}{T/g(T)} = 2g(T)$$

Therefore the total number of buckets ever in existence is $(\text{diam}(\mathbf{s}) + 2)g(T) + 2 \leq (\text{diam}(\mathbf{s}) + 4)g(T)$, so Algorithm 4 performs at most $(\text{diam}(\mathbf{s}) + 4)g(T)$ queries. \square

Lemma D.2. *For each $t \in [T]$, $\mu(s_t, a_t) \geq \mu^m(s_t) - L \text{len}(B_t)$.*

Proof. If we query the mentor at time t , $\mu(s_t, a_t) = \mu^m(s_t)$. Thus assume we do not query the mentor at time t : then there exists $s' \in B_t$ (as defined on line 21 of Algorithm 4) such that $a_t = \pi(s') = \pi^m(s')$. Since s_t and s' are both in B_t , $|s_t - s'| \leq \text{len}(B_t)$. Then by the local generalization assumption,

$$\mu(s_t, a_t) = \mu(s_t, \pi^m(s')) \geq \mu^m(s_t) - L||s_t - s'|| \geq \mu^m(s_t) - L \text{len}(B_t)$$

as required. \square

Lemma D.3. If π^m has at most K segments, $|\mathcal{B}'_V| \leq K$.

Proof. If two buckets B and B' overlap by more than at a single point, one must be the descendant of the other. Thus by definition of \mathcal{B}'_V , any two buckets in \mathcal{B}'_V overlap by at most a single point.

Now consider any $B \in \mathcal{B}'_V$. By definition, there exists $t \in V_T$ such that $s_t \in B$. Then there exists $s' \in B$ (as defined in Algorithm 4) such that $a_t = \pi(s') = \pi^m(s')$. Since $t \in V_T$, we have $\pi^m(s_t) \neq a_t = \pi^m(s')$. Thus s_t and s' are in different segments; let S_B be the segment containing $\min(s_t, s')$. Then $\min(s_t, s') \leq \max(S_B) < \max(s_t, s') \leq \max(B)$.

Consider any $B, B' \in \mathcal{B}'_V$ with $B \neq B'$ and WLOG assume $\max(B) < \max(B')$. Then $\max(S_B) < \max(B) \leq \max(S_{B'})$, so $S_B \neq S_{B'}$. Thus the number of segments in π^m is at least $|\mathcal{B}'_V|$. Since the number of segments in π^m is at most K , we must have $|\mathcal{B}'_V| \leq K$. \square

Lemma D.4. We have $\sum_{t \in V_T} \text{len}(B_t) \leq \frac{2KT}{g(T)^2}$.

Proof. For every $t \in V_T$, we have $B_t = B$ for some $B \in \mathcal{B}_V$, so

$$\sum_{t \in V_T} \text{len}(B_t) = \sum_{B \in \mathcal{B}_V} \sum_{t \in V_T: B=B_t} \text{len}(B_t)$$

Next, observe that every $B \in \mathcal{B}_V \setminus \mathcal{B}'_V$ must have a descendent in \mathcal{B}'_V : otherwise we would have $B \in \mathcal{B}'_V$. Let $\mathcal{A}(B)$ denote the set of ancestors of B , plus B itself. Then we can write

$$\begin{aligned} \sum_{t \in V_T} \text{len}(B_t) &\leq \sum_{B' \in \mathcal{B}'_V} \sum_{B \in \mathcal{A}(B')} \sum_{t \in V_T: B=B_t} \text{len}(B_t) \\ &= \sum_{B' \in \mathcal{B}'_V} \sum_{B \in \mathcal{A}(B')} |\{t \in V_T : B = B_t\}| \cdot \text{len}(B_t) \end{aligned}$$

For any bucket B , the number of time steps t with $B = B_t$ is at most $T/g(T)$. Also recall that $\text{len}(B) = \frac{1}{g(T)2^{d(B)}}$. Therefore

$$\begin{aligned} \sum_{B \in \mathcal{A}(B')} \frac{|\{t \in V_T : B = B_t\}|}{g(T)2^{d(B)}} &\leq \frac{T}{g(T)^2} \sum_{B \in \mathcal{A}(B')} \frac{1}{2^{d(B)}} \\ &= \frac{T}{g(T)^2} \sum_{d=0}^{d(B')} \frac{1}{2^d} \leq \frac{T}{g(T)^2} \sum_{d=0}^{\infty} \frac{1}{2^d} = \frac{2T}{g(T)^2} \end{aligned}$$

Then by Lemma D.3,

$$\sum_{t \in V_T} \text{len}(B_t) \leq \sum_{B' \in \mathcal{B}'_V} \frac{2T}{g(T)^2} = \frac{2T|\mathcal{B}'_V|}{g(T)^2} \leq \frac{2KT}{g(T)^2}$$

as claimed. \square

Lemma D.5. Under the conditions of Theorem D.2, Algorithm 4 satisfies

$$\sum_{t=1}^T \left(\mu^m(s_t) - \min(\mu^m(s_t), \mu(s_t, a_t)) \right) \leq \frac{2LKT}{g(T)^2}$$

Proof. For $t \notin V_T$ we have $\min(\mu^m(s_t), \mu(s_t, a_t)) = \mu^m(s_t)$ by definition, and Lemma D.2 implies that $\min(\mu^m(s_t), \mu(s_t, a_t)) \geq L \text{len}(B_t)$ for all $t \in [T]$. Thus

$$\sum_{t=1}^T \left(\mu^m(s_t) - \min(\mu^m(s_t), \mu(s_t, a_t)) \right) \leq \sum_{t \in V_T} \left(\mu^m(s_t) - \min(\mu^m(s_t), \mu(s_t, a_t)) \right) \leq L \sum_{t \in V_T} \text{len}(B_t)$$

Then by Lemma D.4,

$$\sum_{t=1}^T \left(\mu^m(s_t) - \min(\mu^m(s_t), \mu(s_t, a_t)) \right) \leq \frac{2LKT}{g(T)^2}$$

as required. \square

Lemma D.6. Under the conditions of Theorem D.2, Algorithm 4 satisfies $R_T \leq \frac{2LKT}{g(T)^2}$.

Proof. Let $x_t = \mu^m(s_t)$ and $y_t = \min(\mu^m(s_t), \mu(s_t, a_t))$ for all $t \in [T]$. Then by Lemma B.3,

$$\prod_{t=1}^T \mu^m(s_t) - \prod_{t=1}^T \min(\mu^m(s_t), \mu(s_t, a_t)) \leq \sum_{t=1}^T (\mu^m(s_t) - \min(\mu^m(s_t), \mu(s_t, a_t)))$$

Since $\mu(s_t, a_t) \geq \min(\mu^m(s_t), \mu(s_t, a_t))$ for all $t \in [T]$, we have

$$R_T = \prod_{t=1}^T \mu^m(s_t) - \prod_{t=1}^T \mu(s_t, a_t) \leq \sum_{t=1}^T (\mu^m(s_t) - \min(\mu^m(s_t), \mu(s_t, a_t)))$$

Applying Lemma D.5 completes the proof. \square

Theorem D.2 follows from Lemma D.1 and Lemma D.6.

E OTHER PROOFS

Proposition E.1. Suppose μ is L -Lipschitz continuous and π^m is optimal, i.e., $\mu(s, \pi^m(s)) = \max_{a \in A} \mu(s, a)$ for all $s \in S$. Then π^m satisfies local generalization with constant $2L$.

Proof. For any $s, s' \in S$, we have

$$\begin{aligned} \mu(s, \pi^m(s')) &\geq \mu(s', \pi^m(s')) - L\|s - s'\| && \text{(Lipschitz continuity of } \mu) \\ &\geq \mu(s', \pi^m(s)) - L\|s - s'\| && (\pi^m \text{ is optimal for } s') \\ &\geq \mu(s, \pi^m(s)) - 2L\|s - s'\| && \text{(Lipschitz continuity of } \mu \text{ again)} \\ &= \mu^m(s) - 2L\|s - s'\| && \text{(Definition of } \mu^m(s)) \end{aligned}$$

Since π^m is optimal for s , we have

$$\mu^m(s) + 2L\|s - s'\| \geq \mu^m(s) \geq \mu(s, \pi^m(s'))$$

Thus $-2L\|s - s'\| \leq \mu(s, \pi^m(s')) - \mu^m(s) \leq 2L\|s - s'\|$. This is equivalent to $|\mu(s, \pi^m(s')) - \mu^m(s)| \leq 2L\|s - s'\|$, completing the proof. \square