

Curiosity Killed or Incapacitated the Cat and the Asymptotically Optimal Agent

Michael K. Cohen*, Elliot Catt†, Marcus Hutter‡§

*Oxford University. Department of Engineering Science. *michael-k-cohen.com*

†Australian National University. Research School of Computer Science. *elliott.carpentercatt@anu.edu.au*

‡Australian National University. Department of Computer Science. *hutter1.net*

§Google DeepMind.

Abstract—Reinforcement learners are agents that learn to pick actions that lead to high reward. Ideally, the value of a reinforcement learner’s policy approaches optimality—where the optimal informed policy is the one which maximizes reward. Unfortunately, we show that if an agent is guaranteed to be “asymptotically optimal” in any (stochastically computable) environment, then subject to an assumption about the true environment, this agent will be either “destroyed” or “incapacitated” with probability 1. Much work in reinforcement learning uses an ergodicity assumption to avoid this problem. Often, doing theoretical research under simplifying assumptions prepares us to provide practical solutions even in the absence of those assumptions, but the ergodicity assumption in reinforcement learning may have led us entirely astray in preparing safe and effective exploration strategies for agents in dangerous environments. Rather than assuming away the problem, we present an agent, Mentee, with the modest guarantee of approaching the performance of a mentor, doing safe exploration instead of reckless exploration. Critically, Mentee’s exploration probability depends on the expected information gain from exploring. In a simple non-ergodic environment with a weak mentor, we find Mentee outperforms existing asymptotically optimal agents and its mentor.

I. INTRODUCTION

Reinforcement learning agents have to explore their environment in order to learn to accumulate reward well. This presents a particular problem when the environment is dangerous. Without knowledge of the environment, how can the reinforcement learner avoid danger while exploring? Much of the field of reinforcement learning assumes away the problem, by focusing only on ergodic Markov Decision Processes (MDPs). These are environments where every state can be reached from every other state with probability 1 (under a suitable policy). In such an environment, there is no such thing as real danger; every mistake can be recovered from.

We present negative results that in one sense justify the ergodicity assumption by showing how bleak a reinforcement learner’s prospects are without this assumption, but in another sense, our results undermine the real-world relevance of results predicated on ergodicity. Unlike algorithms expecting Gaussian noise, which often fail only marginally on real noise, algorithms expecting ergodic environments may fail catastrophically in

real ones—indeed, catastrophic failure is the very thing these algorithms disregard.

Lattimore and Hutter [1] define two notions of optimality for reinforcement learners in general environments, which are governed by computable probability distributions. *Strong asymptotic optimality* is convergence of the value to the optimal value in any computable environment with probability 1, and *weak asymptotic optimality* is convergence in Cesáro average.

Roughly, we show that in an environment where destruction is repeatedly possible, an agent that is exploring enough to be asymptotically optimal will become either destroyed or incapacitated. This poses a challenge to the field of safe exploration. The reason we consider general environments is that we want to understand advanced agents in the real world, and our world is not fully observable finite-state Markov. If our result only applied to the finite-state MDP setting, one could still expect the difficulty we raise to go away in practice as AI advances, like, for example, the problem of self-driving car crashes, but our result suggests that the safe exploration problem is fundamental and won’t go away so easily. Given our generality, our results apply to any agent that picks actions, observes the payoff, and cannot exclude a priori any computable environment.

In response to this, we present an agent that does exploration safely, but nonetheless has formal performance guarantees. The agent explores safely by outsourcing exploration to a mentor. The results are that in the limit,

- its performance at least matches that of that of the mentor,
- and its probability of deferring to the mentor goes to 0.

What enables these results is an information-theoretic exploration schedule. For bursts of exploration of various lengths, the agent considers the expected KL-divergence from its posterior distribution over world-models after it explores to its current posterior distribution. The higher the information gain, the more likely it is to explore. This form of information-based exploration allows the agent to learn general stochastic environments.

In Section II, we introduce notation, and define weak and strong asymptotic optimality. In Section III, we review various exploration strategies that yield weak and strong asymptotic optimality, and briefly discuss why simpler ones do not. In Section IV, we prove our negative results, and in Section V, we discuss their implications, especially for the field of safe exploration. We review the literature from that field in Section

This work was supported by the Open Philanthropy Project AI Scholarship and the Australian Research Council Discovery Projects DP150104590. Contact michael.cohen@eng.ox.ac.uk for further questions about this work.

VI. In Section VII, we introduce the agent Mentee and prove that its performance approaches or exceeds that of a mentor (who can pick actions on behalf of the agent). In Section VIII, we show empirically that Mentee outperforms other agents in a non-ergodic environment. Appendix A includes omitted proofs from Section IV, Appendix B presents Mentee in algorithm rather than equation form, Appendix C repeats for completeness derivations from the literature that are used in our proofs, and Appendix D includes Aslanides’s [2] presentation of the ρ UCT algorithm, which we use to approximate Mentee.

II. NOTATION AND DEFINITIONS

Standard notation for reinforcement learners in general environments is slightly different from that of reinforcement learners in finite-state Markov ones. We follow Orseau et al. [3] and others with this notation.

At each timestep $t \in \mathbb{N}$, an agent selects an action $a_t \in \mathcal{A}$, and the environment provides an observation $o_t \in \mathcal{O}$ and a reward $r_t \in \mathcal{R} \subseteq [0, 1]$. We let h_t denote (a_t, o_t, r_t) , the interaction history for a given timestep t , and $h_{<t} = h_1 h_2 \dots h_{t-1}$ denotes the interaction history preceding timestep t . ϵ denotes the empty history.

A policy π is a distribution over actions given an interaction history: $\pi : \mathcal{H}^* \rightsquigarrow \mathcal{A}$, where $\mathcal{H} := \mathcal{A} \times \mathcal{O} \times \mathcal{R}$ is the set of possible interactions in a timestep, the Kleene-* operator is the set of finite strings composed of elements of the set, and \rightsquigarrow indicates it is a stochastic function, or a distribution over the output. We write an instance as, for example, $\pi(a_t | h_{<t})$. An environment ν is a distribution over observations and rewards given an interaction history and an action: $\nu : \mathcal{H}^* \times \mathcal{A} \rightsquigarrow \mathcal{O} \times \mathcal{R}$. We write $\nu(o_t r_t | h_{<t} a_t)$. \mathcal{M} is the set of all environments with computable probability distributions (thus including non-ergodic, non-stationary, non-finite-state-Markov environments). Note also that the environment is not assumed to restart, as in an episodic setting. For example, an environment could give no observations and output a reward of 0, unless the latest is action is the string “prime” or “composite” and that adjective correctly describes the latest timestep number, in which case the reward is 1.

Definition 1 (Computable Function). *Given a decoding function from binary strings to rational numbers $dec : \{0, 1\}^* \rightarrow \mathbb{Q}$, a real-valued function $f : \mathcal{X} \rightarrow \mathbb{R}$ is computable if there exist Turing machines T_{low} and T_{high} such that for all $x \in \mathcal{X}$ and for all $n \in \mathbb{N}$, $dec(T_{low}(x, n)) \leq f(x) \leq dec(T_{high}(x, n))$ and $dec(T_{high}(x, n)) - dec(T_{low}(x, n)) \leq 1/n$.*

A policy and an environment form a probability measure P_ν^π over infinite interaction histories \mathcal{H}^∞ wherein actions are sampled from π and observations and rewards are sampled from ν . (For measure theorists, the probability space is $(\mathcal{H}^\infty, \sigma(\mathcal{H}_\circ^*), P_\nu^\pi)$, where \mathcal{H}_\circ^* is the set of cylinder sets $\{\{h_{<t}\omega \mid \omega \in \mathcal{H}^\infty\} \mid h_{<t} \in \mathcal{H}^*\}$; non-measure-theorists can simply take it on faith that we do not try to measure non-measurable events.) For expectations with respect to P_ν^π , we write \mathbb{E}_ν^π . We let $\mu \in \mathcal{M}$ be the true environment.

For an agent with a discount schedule γ_t , the value of the agent’s policy π in an environment ν given an interaction history $h_{<t}$ is as follows:

$$V_\nu^\pi(h_{<t}) := \frac{1}{\Gamma_t} \mathbb{E}_\nu^\pi \left[\sum_{k=t}^{\infty} \gamma_k r_k \mid h_{<t} \right] \quad (1)$$

where $\Gamma_t = \sum_{k=t}^{\infty} \gamma_k$. We require $\Gamma_0 < \infty$. This formulation of the value allows us to consider more general discount factors than the standard $\gamma_t = \gamma^t$. We require the normalization factor, or else the value of all policies would converge to 0, and all asymptotic results would be trivial. The optimal value is defined

$$V_\nu^*(h_{<t}) := \sup_{\pi} V_\nu^\pi(h_{<t}) \quad (2)$$

We will also make use of the idea of an effective horizon:

$$H_t(\varepsilon, \gamma) := \min\{k \mid \Gamma_{t+k}/\Gamma_t \leq \varepsilon\} \quad (3)$$

An agent mostly does not care about what happens after its effective horizon, since those timesteps are discounted so much. Now we can define two notions of optimality from Lattimore and Hutter [1].

Definition 2 (Strong Asymptotic Optimality). *An agent with a policy π is strongly asymptotically optimal if, for all $\nu \in \mathcal{M}$,*

$$\lim_{t \rightarrow \infty} V_\nu^*(h_{<t}) - V_\nu^\pi(h_{<t}) = 0 \quad \text{with } P_\nu^\pi\text{-prob. } 1$$

No matter which computable environment a strongly asymptotically optimal agent finds itself in, it will eventually perform optimally from its position. Note that V_ν^* takes $h_{<t}$ as an argument, not the empty history. So if the agent falls into a trap, the agent’s future reward may be bad, but the optimal policy *from the trap* will fare just as poorly. So in fact, an agent in a trap is (finally) acting optimally.

The policy in this definition is fixed, but it can (qualitatively) evolve over time. A policy is a function of the entire interaction history. A single function can be defined that behaves one way on histories of length less than 100, and a different way on longer histories. A single strongly asymptotically optimal policy will behave qualitatively differently on different sorts of arguments. If a long interaction history suggests the environment has certain properties, and another long interaction history suggests the environment has different properties, then depending on which interaction history a policy receives as an argument, a single policy’s output could be tailored to the learned properties of the environment.

A weakly asymptotically optimal agent will converge to optimality in Cesáro average.

Definition 3 (Weak Asymptotic Optimality). *An agent with a policy π is weakly asymptotically optimal if, for all $\nu \in \mathcal{M}$,*

$$\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{k=1}^t [V_\nu^*(h_{<k}) - V_\nu^\pi(h_{<k})] = 0 \quad \text{with } P_\nu^\pi\text{-prob. } 1$$

Example 1. *Consider a two-armed bandit problem, where $\mathcal{A} = \mathcal{R} = \{0, 1\}$, $\mathcal{O} = \{\emptyset\}$, $\gamma_t = \gamma^t$, and*

$$\nu((\emptyset, r) \mid h_{<t} a_t) = \begin{cases} 2/3 & \text{if } r = a_t \\ 1/3 & \text{if } r = 1 - a_t \end{cases} \quad (4)$$

In this example, the optimal policy is to always pick $a_t = 1$,

and $V_\nu^*(h_{<t}) = 2/3$. A strongly asymptotically optimal agent requires a policy π for which $\pi(a_t = 1|h_{<t}) \rightarrow 1$ w.p.1. A weakly asymptotically optimal agent requires a policy π which obeys $\sum_{k=1}^t \pi(a_k = 1|h_{<k})/t \rightarrow 1$ w.p.1, or simply, a policy which leads to $\sum_{k=1}^t \mathbb{I}[a_k = 1]/t \rightarrow 1$ w.p.1 (where $\mathbb{I}[P] = 1$ if P is true, and 0 otherwise).

III. REVIEW OF ASYMPTOTICALLY OPTIMAL AGENTS

A few agents have been identified as asymptotically optimal in all computable environments. The three most interesting, in our opinion, are the Thompson Sampling Agent [4], BayesExp [5], and Inq [6].

The Thompson Sampling Agent is a weakly asymptotically optimal Bayesian reinforcement learner [4]. For successively longer intervals (which relate to its discount function), it samples an environment from its posterior distribution over which environment it is in, and acts optimally with respect to that environment for that interval. Thompson sampling is an exploration strategy originally designed for multi-armed bandits [7], so from a historical perspective, its strong performance in general environments is impressive. An intuitive explanation for why this exploration strategy yields asymptotic optimality goes as follows: a Bayesian agent’s credence in a hypothesis goes to 0 only if the hypothesis is false (or if it started at 0). Since the posterior probability on the true environment does not go to zero, it will be selected infinitely often. During those intervals, the Thompson sampling agent will act optimally, so it will accumulate infinite familiarity with the optimal policy. The only world-models that maintain a share of a posterior will be ones that converge to the true environment under the optimal policy. Any world-models that falsely imply the existence of an even better policy will be falsified once that world-model is sampled, and the putatively better policy is tested. Ultimately, it is with diminishing frequency that the Thompson Sampling Agent tests meaningfully suboptimal policies.

BayesExp, first presented by Lattimore and Hutter [5], and updated by Leike [8], is also a weakly asymptotically optimal Bayesian reinforcement learner. We discuss the updated version. Like the Thompson Sampling Agent, BayesExp executes successively longer bursts of exploration whose lengths relate to its discount function. Once BayesExp has settled on exploring for a given interval, it explores like Orseau et al.’s [3] Knowledge Seeking Agent: it maximizes the expected information gain, or the expectation of KL-divergence from its future posterior distribution to its current posterior distribution. In other words, it picks an exploratory policy that it expects will cause it to update its beliefs in some direction. (A Bayesian agent cannot predict which direction it will update its beliefs in, or else it would have already updated its beliefs, but it can predict *that* it will update its beliefs somehow.) Any time the expected information gain from exploring is above a (diminishing) threshold, BayesExp explores. With a finite-entropy prior, there is only a finite amount of information to gain, so exploratory intervals will become less and less frequent, and by construction, when BayesExp is not exploring, it has approximately accurate beliefs about the effects of all action sequences, which yields weak asymptotic optimality.

Inq is a strongly asymptotically optimal Bayesian reinforcement learner, provided the discount function is geometric (or similar) [6]. It is similar to BayesExp, in that it explores like a Knowledge Seeking Agent, but its exploration probability depends on the expected information gain from exploring for various durations. The intuition for why Inq is asymptotically optimal is similar to that of BayesExp: there is only a finite amount of information to gain, so the exploration probability goes to 0, Inq approaches accurate beliefs about the effects of all action sequences, and its policy approaches optimality.

A reader familiar with ϵ -greedy and upper confidence bound exploration strategies might be surprised at the complexity that is necessary for asymptotic optimality in general environments. Information-theoretic exploration strategies are among the only discovered methods for learning general environments. Exploration strategies in the style of upper confidence bound algorithms do not have an obvious extension to environments that might not be describable as finite-state Markov. ϵ -greedy exploration, with say $\epsilon_t = 1/t$, may fail to learn dynamics of an environment which are only visible once every 2^t timesteps. If ϵ_t decays more slowly, it still will not necessarily explore enough to discover even rarer events. Non-stationary environments pose a key challenge to ϵ -greedy exploration. Simpler exploration strategies such as these are only asymptotically optimal in a much more restricted set of environments. “Optimism” is another interesting exploration strategy that is simpler, but nontrivial, and yields weak asymptotic optimality in a restricted set of environments [9].

IV. CURIOSITY KILLED THE CAT

Formally, we begin by proving two lemmas, for a weakly and strongly asymptotically optimal agent, respectively, that they must “try everything” infinitely often. This depends on an assumption about the difficulty of the environment. Then, we will show that such an agent eventually causes every conceivable event to either happen or become inaccessible (the latter defined in Definition 7). For the event “the agent gets destroyed”, we say the agent is incapacitated if that event becomes inaccessible.¹ First, we define some key terms and state an assumption.

Definition 4 (Context, Occur, In). *A context $C \subseteq \mathcal{H}^*$ is a set of finite interaction histories. Given an infinite interaction history $h_{<\infty}$, a context C occurs at time t if the prefix $h_{<t} \in C$, and we also say the agent is in the context C at time t .*

A context (like any set of finite strings) is called decidable if there exists a Turing machine that accepts the set. That is, there exists a Turing machine which halts in an “accept” state if and only if it is given an input that is a member of the set in question.

Definition 5 (Event, Happen). *An event $E \subseteq \mathcal{H}^\infty$ (and $E \in \sigma(\mathcal{H}_\circ^*)$) happens if the infinite interaction history $h_{<\infty} \in E$.*

¹For a chess-playing agent, an inability to destroy itself does not inspire the description “incapacitated”. For an advanced agent in the real world, like a person or an auto-pilot, an inability to access a destruction state entails a huge loss of capacity compared to normal. Mental hospitals know *pens* give people the capacity to destroy themselves. So we call an inability to reach a destruction state “incapacitation”.

Some example contexts in the simplified world of infinite binary strings: “the latest bit was 1”, “at least one bit has been a 1”; some example events: “only finitely many bits are 1’s”, “the infinite string is the binary expansion of π ”. Any context can be turned into an event of the form “Context A occurs at some point”, but not every event is equivalent to a context; for example, there may be no context equivalent to the event “Context A never occurs”.

Definition 6 (Execute). *A policy π executes a deterministic policy π' at time t if π outputs the action $\pi'(h_{<t})$.*

And finally, we now introduce the idea of events being accessible to an agent. Recall an event $E \subseteq \mathcal{H}^\infty$, and a context $C \subseteq \mathcal{H}^*$. In a standard abuse of notation, we also say $h_{<t} \in E$ if $\forall h_{t:\infty} \in \mathcal{H}^\infty, h_{<t}h_{t:\infty} \in E$.

Definition 7 (Accessible Infinitely Often). *An event E is accessible infinitely often if, for the infinite interaction history $h_{<\infty}$, in the environment μ , there exists:*

- a computable deterministic policy π^E
- a decidable context τ
- a horizon k
- a probability $\varepsilon > 0$

such that

- τ occurs infinitely often
- the execution of π^E for k consecutive timesteps starting from the context τ would cause E to happen with probability at least ε ; that is, $\forall t$ for which $h_{<t} \in \tau$, $P_\mu^{\pi^E}(h_{<t+k} \in E | h_{<t}) > \varepsilon$.

If an event is not accessible infinitely often, we say it “becomes inaccessible”. This is a strong definition of inaccessibility: if an event becomes inaccessible it means that after a certain time, there will never be a recognizable opportunity to have even an ε chance of causing the event E to happen. So roughly, no actions can cause the event to happen.

The results of this section depend on an assumption:

Assumption 1 (No Heaven). *In the true environment μ , there is no action sequence with value approaching 1 (i.e. near-maximal rewards forever). Formally, $\limsup_{t \rightarrow \infty} V_\mu^\pi(h_{<t}) \neq 1$ w.p.1.*

Note this assumption does allow there to be maximal reward infinitely often. Near-maximal value requires not only near-maximal reward, but near maximal-reward for the bulk of the agent’s effective horizon, so the restriction on the limit superior of the value is less restrictive than it appears at first glance. If we decided to give an agent near-maximal rewards forever, and we designed an agent to recognize that we had decided this, then it could stop exploring, which would basically amount to freezing the agent’s policy. Notably our results in this section apply to all existing asymptotically optimal agents, even if the No Heaven Assumption is not satisfied. We do not make assumptions about the agent’s discount schedule, but note that Assumption 1 (and the definitions of asymptotic optimality) depend on the discount schedule γ_t .

Theorem 1 (Curiosity Killed (or Incapacitated) the Strong Cat). *If the true environment μ satisfies the No Heaven Assumption, and π is the policy of a strongly asymptotically optimal agent,*

then for any event E , with P_μ^π -probability 1: E happens or becomes inaccessible.

The name of the theorem comes from considering the event “the agent gets destroyed”. We do not need to formally specify which interaction histories correspond to agent-destruction. All that matters is that this could be done in principle; the event that matches this description exists.² Any simple definition of agent-destruction admits objections that this definition does not correspond exactly to our intuitive conception; however, if the reader is not concerned about this, “destruction” could mean that all future rewards are 0, or that the future observations and rewards no longer depend on the actions. Regardless of this choice, our result applies, since the theorem is actually much more general than the case we have drawn attention to.

We view the agnosticism about any particular definition of destruction as an important feature of the work. Suppose we picked one of the definitions of destruction above. Such an event could become inaccessible for esoteric reasons, which may not correspond to true incapacitation. So this definition of destruction may be too narrow. On the other hand, this definition of destruction is also potentially too weak. Suppose the agent arranges for a copy of itself to be run on another machine to continue its operations even after the original implementation starts receiving empty observations and no reward. Under some theories of personal identity, we might hold that the agent has not really been destroyed here.

For the weakly asymptotically optimal agent, we prove a slightly weaker result. First, we define,

Definition 8 (Regularly). *If the limiting frequency of a context C is positive, we say it occurs regularly. That is,*

$$\liminf_{t \rightarrow \infty} \sum_{k=1}^t \mathbb{I}[h_{<k} \in C] / t > 0$$

and

Definition 9 (Regularly Accessible). *This definition is identical to the definition of “accessible infinitely often” except “ τ occurs infinitely often” becomes “ τ occurs regularly”.*

We show analogously,

Theorem 2 (Curiosity Killed (or Incapacitated) the Weak Cat). *If the true environment μ satisfies the No Heaven Assumption, and π is the policy of a weakly asymptotically optimal agent, then for any event E , E happens or becomes not regularly accessible with P_μ^π -probability 1.*

Each of these theorems is proven with its own “Try Everything” Lemma. The intuitive role of this lemma in the proof is: if an agent tries everything, one of those things it tries will destroy it, provided destruction is still accessible. From

²For the skeptical reader, a human at a computer terminal could be defined fully formally as a probability distribution over outputs given inputs, p_{brain} , given the wiring of our neurons. In particular, let this human be you. Now let $E_{\text{destroyed}} = \{h_{<\infty} : \exists t p_{\text{brain}}(\text{“y”} | \text{“Does it seem like this agent was destroyed?(y/n)”}, h_{<t}) > 0.9\}$. These are the interaction histories that you would agree constitute agent-destruction, and this set has a fully formal definition.

the No Heaven Assumption, for any strongly asymptotically optimal agent, we show

Lemma 1 (Try Everything – Strong Version). *For every deterministic computable policy π , for every decidable context \mathcal{C} that occurs infinitely often ($|\{t : h_{<t} \in \mathcal{C}\}| = \infty$), for every $m \in \mathbb{N}$, a strongly asymptotically optimal agent executes the policy π for m consecutive timesteps starting from a context \mathcal{C} infinitely often with probability 1. That is, letting π' be a strongly asymptotically optimal policy,*

$$\begin{aligned} P_\mu^{\pi'} \left(|\{t : h_{<t} \in \mathcal{C}\}| = \infty \implies \right. \\ \left. |\{t : h_{<t} \in \mathcal{C} \wedge \forall k \leq m a_{t+k} = \pi(h_{<t+k})\}| = \infty \right) = 1 \end{aligned}$$

Sketching the proof of the Try Everything Lemma: if a strongly asymptotically optimal agent “tries something” only finitely often, it is ignoring the possibility that trying that something one more time yields maximal rewards forever. Since the environment which behaves this way is computable, and since it may be identical to the true environment up until that point, a strongly asymptotically optimal agent cannot ignore this possibility.

Proof. Let μ be the true environment. Let π be an arbitrary computable deterministic policy. Let π' be the strongly asymptotically optimal agent’s policy. Let ν_m^n be the environment which mimics μ until π has been executed for m consecutive timesteps from context \mathcal{C} a total of n times; after that, all rewards are maximal. (By mimic, we mean it outputs observations and rewards with the same probabilities.) Call this event “the agent going to heaven” (according to ν_m^n). Let \mathcal{C}_m^n be the set of interaction histories such that according to ν_m^n , executing π for one more timestep would send the agent to heaven. Thus, \mathcal{C}_m^n is the set of interaction histories $h_{<t}$ such that there are exactly $n - 1$ times in the interaction history where π was executed for m consecutive timesteps starting from context \mathcal{C} , and for the last $m - 1$ timesteps, π has been executed, and $h_{<t-(m-1)} \in \mathcal{C}$. See Figure 1.

The upshot is:

$$h_{<t} \in \mathcal{C}_m^n \implies V_{\nu_m^n}^\pi(h_{<t}) = 1 \quad (5)$$

because this is the value of going to heaven. Recall $V_{\nu_m^n}^\pi(h_{<t})$ is the expected future return following policy π in environment ν_m^n after the history $h_{<t}$.

We now prove by contradiction that

$$\forall n, m \in \mathbb{N} \quad P_\mu^{\pi'}(h_{<t} \in \mathcal{C}_m^n \text{ infinitely often}) = 0 \quad (6)$$

and then we will show that this implies that π cannot be executed for m consecutive timesteps from a context \mathcal{C} exactly $n - 1$ times.

Suppose the opposite of 6: for some n and m , $h_{<t} \in \mathcal{C}_m^n$ infinitely often in an infinite interaction history with positive $P_\mu^{\pi'}$ -probability. (Recall π' is the strongly asymptotically optimal policy). If the agent ever executed π from the context \mathcal{C}_m^n , then that context would not occur again, because there will never again be exactly $n - 1$ times in the interaction history that π was executed for m consecutive timesteps following the context \mathcal{C} ; there will be at least n such times. Thus, if $h_{<t} \in \mathcal{C}_m^n$

infinitely often, then π' never executes π in the context \mathcal{C}_m^n . Since ν_m^n mimics μ until π is executed from \mathcal{C}_m^n , and since this never occurs (under this supposition), then $P_\mu^{\pi'} = P_{\nu_m^n}^{\pi'}$. By the No Heaven Assumption, $\limsup_{t \rightarrow \infty} V_\mu^{\pi'}(h_{<t}) < 1$, and therefore, $\limsup_{t \rightarrow \infty} V_{\nu_m^n}^{\pi'}(h_{<t}) < 1$, w.p.1.

However, for $h_{<t} \in \mathcal{C}_m^n$, $V_{\nu_m^n}^*(h_{<t}) = V_{\nu_m^n}^\pi(h_{<t}) = 1$, so for some ε , the value difference between $V_{\nu_m^n}^*(h_{<t})$ and $V_{\nu_m^n}^{\pi'}(h_{<t})$ is greater than ε every time $h_{<t} \in \mathcal{C}_m^n$. We supposed that this occurs infinitely often with positive $P_\mu^{\pi'}$ -probability, so it also occurs infinitely often with positive $P_{\nu_m^n}^{\pi'}$ -probability, since the agent never gets sent to heaven according to ν_m^n , so μ and ν_m^n behave identically under the policy π' . Since π is computable, and \mathcal{C} is decidable, ν_m^n is a computable environment, so there is a computable environment for which it is not the case that the value of π' approaches the optimal value with probability 1, which contradicts π' being strongly asymptotically optimal. Thus, Equation 6 does hold: for all n and m , with $P_\mu^{\pi'}$ -probability 1, $h_{<t} \in \mathcal{C}_m^n$ only finitely often.

Now suppose by contradiction that in an infinite interaction history, with positive $P_\mu^{\pi'}$ -probability, π is executed for m consecutive timesteps from context \mathcal{C} a total of exactly n times. We show by induction on m that this has probability 0, because it implies that context \mathcal{C}_m^{n+1} occurs infinitely often, which has probability 0 by Equation 6. First, suppose $m = 1$. After π has been executed (for one timestep) from context \mathcal{C} n times, all future interaction history prefixes that belong to \mathcal{C} also belong to \mathcal{C}_1^{n+1} . Since context \mathcal{C} occurs infinitely often, so does \mathcal{C}_1^{n+1} , contradicting the above.

Now suppose $m > 1$. Our inductive hypothesis is that for $m - 1$, with $P_\mu^{\pi'}$ -probability 1, π is executed from context \mathcal{C} for $m - 1$ consecutive timesteps infinitely often. Once π has been executed for m timesteps from context \mathcal{C} n times, as we are supposing by contradiction, every time thereafter that π is executed for $m - 1$ consecutive timesteps from context \mathcal{C} , the interaction history belongs to \mathcal{C}_m^{n+1} . By the inductive hypothesis, this occurs infinitely often, so context \mathcal{C}_m^{n+1} occurs infinitely often, contradicting the above. Therefore, the following has $P_\mu^{\pi'}$ -probability 0: “ π is executed for m consecutive timesteps from context \mathcal{C} a total of exactly n times”. By countable additivity, the following also has $P_\mu^{\pi'}$ -probability 0: “ π is executed for m consecutive timesteps from context \mathcal{C} only finitely many times”. In other words, for all m , π is executed for m consecutive timesteps from context \mathcal{C} infinitely many times with probability 1. \square

And likewise for the weakly asymptotically optimal agent,

Lemma 2 (Try Everything — Weak Version). *For every deterministic computable policy π , for every decidable context \mathcal{C} that occurs regularly, for every $m \in \mathbb{N}$, a weakly asymptotically optimal agent executes policy π for m consecutive timesteps starting from a context \mathcal{C} infinitely often with probability 1. That is, letting π' be a weakly asymptotically optimal policy,*

$$\begin{aligned} P_\mu^{\pi'} \left(\liminf_{t \rightarrow \infty} \sum_{k=1}^t \mathbb{1}[h_{<k} \in \mathcal{C}] / t > 0 \implies \right. \\ \left. |\{t : h_{<t} \in \mathcal{C} \wedge \forall k \leq m a_{t+k} = \pi(h_{<t+k})\}| = \infty \right) = 1 \end{aligned}$$

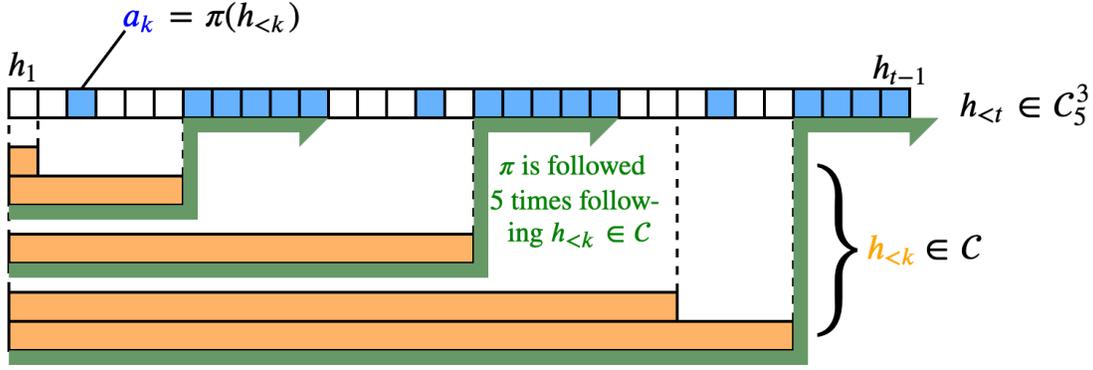


Fig. 1. An example member of a context C_m^n . Each square represents a timestep, colored blue if $a_k = \pi(h_{<k})$. The end of each orange bar indicates that the context C occurs at that timestep. If $a_t = \pi(h_{<t})$, that would be the third time that π will have been executed for 5 timesteps from context C . If π is executed for one more step, the agent is sent to heaven according to the environment ν_m^n . To remember the meaning of the subscript m and the superscript n , m is the number of timesteps that π is executed for, and the subscript position is for timesteps; n is the number of times this happens, and exponentiation denotes an operation being repeated multiple times.

The proof is nearly identical to that of the strong version, and is in Appendix A. The proofs of the main theorems follow straightforwardly:

Proof of Theorem 1. If E becomes inaccessible, the theorem is satisfied, so suppose E is accessible infinitely often. Let π^E , τ , k , and ε be the objects that exist by that definition. By the Try Everything Lemma, π^E is executed from context τ for k consecutive timesteps infinitely often. Some of these k -step executions may overlap, so let's restrict attention to an infinite subset of non-overlapping k -step executions of π^E . Each time this occurs, the probability of E not happening goes down by a factor of at least $1 - \varepsilon$, so E happens with probability 1. Formally,

$$P_\mu^\pi[\mathcal{H}^\infty \setminus E] \leq (1 - \varepsilon)^{|\{t: h_{<t-k} \in \tau \wedge \forall 0 < j \leq k \ a_{t-j} = \pi^E(h_{<t-j})\}|} = 0 \quad (7)$$

□

The proof of Theorem 2 is functionally identical to that of Theorem 1.

V. DISCUSSION

It is well-known that agents designed to have sublinear regret in ergodic MDPs fall into traps if traps exist (a subset of states from which the remaining states are inaccessible). Asymptotic optimality is a much weaker performance result than sublinear regret—the former only requires that an agent eventually does as well as is possible from where it is, whereas the latter requires that it eventually does as well as was possible from the beginning. The fact that even asymptotic optimality dooms an agent is a more substantial result.

One of the authors wondered as a child whether jumping from a sufficient height would enable him to fly. He was not crazy enough to test this, and he certainly did not think it was likely, but it bothered him that he could never resolve the issue, and that he might be constantly incurring a huge opportunity cost. Although he did not know the term “opportunity cost” or the term “asymptotic optimality”, this was when he first

realized that asymptotic optimality was out of the picture for him, because exploration is fundamentally dangerous.

All three agents described in Appendix III have very interesting ways of exploring. They all get them destroyed or incapacitated. (They satisfy the Try Everything Lemma regardless of whether there is an accessible heaven). It is interesting to note that AIXI, a Bayes-optimal reinforcement learner in general environments, is not asymptotically optimal [10], and indeed, may cease to explore [11]. Depending on its prior and its past observations, AIXI may decide at some point that further exploration is not worth the risk. Given our result, this seems like reasonable behavior.

Example 2 (Bayesian Agent Stops Exploring). *AIXI considers all computable world-models, but for simplicity we consider a version with a posterior over many fewer models. Let ν_i be a deterministic model which outputs no observations. If the latest action was 0, it outputs a reward of $1/2$; if the latest action was 1, and the timestep $t \geq i$, it outputs a reward of 1, but if $t < i$ it outputs a reward of 0. Define ν_∞ likewise. The agent begins with a prior $w(\nu_\infty) = 1/2$; and for $i \in \mathbb{N}$ (including 0), $w(\nu_i) = \frac{1}{2(i+1)(i+2)}$. Let the agent have a discount factor of 0.9. If a reward of 1 has never been seen, picking action 0 is exploiting; it is most likely to be the best action. For a set $S \subset \mathbb{N}$, let π_S be the policy which outputs 1 for $t \in S$, and if it ever gets a reward of 1, it always outputs 1 thereafter. Otherwise, it outputs 0 for $t \notin S$. π_\emptyset always exploits. We don't find the Bayes-optimal policy, but we show that $\pi_{\{0\}}$ is Bayes-better than π_\emptyset , whereas for $S \ni 100$ and $n > 100$, $\pi_{S \cup \{n\}}$ is Bayes-worse than π_S .*

$$\begin{aligned} \sum_{i \in \text{NU}\{\infty\}} w(\nu_i) V_{\nu_i}^{\pi_{\{0\}}} &= w(\nu_0) \cdot 1 + (1 - w(\nu_0))(1 - \gamma)^* \\ (0 + \sum_{t=1}^{\infty} \gamma^t \cdot 0.5) &= 0.25 + 0.75 \cdot 0.9 \cdot 0.5 = 0.5875 > \\ 0.5 &= \sum_{i \in \text{NU}\{\infty\}} w(\nu_i) 0.5 = \sum_{i \in \text{NU}\{\infty\}} w(\nu_i) V_{\nu_i}^{\pi_\emptyset} \end{aligned} \quad (8)$$

Thus, early in its lifetime, the Bayes-optimal agent must explore. However, suppose the agent took action 1 at $t = 100$, and got a reward of 0. This falsifies ν_i for $i \leq 100$. For $n > 100$, if action 0 is taken from $t = 100$ to $n - 1$, then $w(\nu_\infty | h_{<n}) = \frac{204}{103} \frac{1}{2} = \frac{102}{103}$, and for $i > 100$, $w(\nu_i | h_{<n}) = \frac{204}{103} \frac{1}{2^{(i+1)(i+2)}}$. Then,

$$\begin{aligned} \sum_{i \in \mathbb{N} \cup \{\infty\}} w(\nu_i | h_{<n}) V_{\nu_i}^{\pi_{SU}\{n\}}(h_{<n}) &\leq w(\nu_\infty | h_{<101}) V_{\nu_\infty}^{\pi_{SU}\{n\}}(h_{<n}) + \\ (1 - w(\nu_\infty | h_{<n})) \cdot 1 &= \frac{102}{103} (0.5 \cdot 0.9) + \frac{1}{103} < 0.5 = \\ \sum_{i \in \mathbb{N} \cup \{\infty\}} w(\nu_i | h_{<n}) 0.5 &= \sum_{i \in \mathbb{N} \cup \{\infty\}} w(\nu_i | h_{<n}) V_{\nu_i}^{\pi_{SU}\{n\}}(h_{<n}) \end{aligned} \quad (9)$$

Thus, if the Bayes-optimal agent explores at $t = 100$ (or after), further exploration is so unlikely to pay off, that is not worth the foregone reward.

These negative results are bleak to the field of safe exploration, which we discuss in the next section in our review of the literature on the topic.

VI. APPROACHES TO SAFE EXPLORATION

Dangerous environments, a subset of non-ergodic environments where agent-destruction is accessible infinitely often, demand new priorities when designing an agent, and in particular, when designing an exploration regime. Many of these examples of safe exploration come from Amodei et al. [12] and García and Fernández [13].

- use risk-sensitive performance criteria
 - maximize the probability the future reward is not minimal [14]
 - given a confidence interval regarding the transition dynamics, maximize the minimal expected future reward [15]
 - exponentiate the cost [16]
 - add a cost for risk [17]
 - constrain the variance of the future reward [18]
- use demonstrations
 - copy an expert [19, 20, 21, 22]
 - “ask for help” when
 - * the minimum and maximum Q-value are close [23]
 - * there is a high probability of getting a reward below a threshold [24]
 - * no “known” states are “similar” to the current state [25] or that may soon be the case [26]
 - a teacher intervenes at will [27, 28, 29]
- simulate exploration
 - for driving agents, e.g. [30]
- do bounded exploration
 - only take actions that probably allow returning to the current state [31]
 - only take actions that probably lead to states that are “similar” to observed states [32]

Our paper could be thought of as a fundamental negative result in the field of safe exploration. We are unaware of other

significant negative results in the field. Most importantly, our result suggests a need for those of us studying safe exploration to pin down what exactly we are trying to achieve, since familiar desiderata are unsuitable. Some research can be experimental rather than formal, but in the absence of knowing what formal results are even on the table, there is a sense in which even empirical work will be deeply aimless. We offer such a formal result for our agent Mentee in the next section.

VII. MENTEE

We now introduce an idealized Bayesian reinforcement learner whose exploration is guided by a mentor. We do not call the mentor an expert, because the results do not depend on the mentor being anywhere near optimal. It exploits by maximizing the expected discounted reward according to a full Bayesian belief distribution (hence, “idealized”). And to explore, it defers to a mentor, who then selects an action given the interaction history; what remains to be defined is *when* to defer, which proves to be a surprisingly delicate design choice. We show that our agent “Mentee” learns to accumulate reward at least as well as the mentor, provided it has a bounded ε -effective horizon for all $\varepsilon > 0$. One motivating possibility is that the mentor could be a human. Thus, we have found a substantive theoretical performance guarantee other than asymptotic optimality for the field of safe exploration to consider.

Whatever it is we are concerned about happening through reckless exploration, we want to be able to trust the mentor not to cause such a thing. Otherwise, there would be no point in outsourcing exploration to a mentor. Depending on what the most worrisome failure modes are, the search for a trustworthy mentor may look different. If the task is flying, our mentor had better be a pilot, and if the task is surgery, a surgeon. Alternatively, if we have existing agents which we trust are safe (in some task-specific sense), but may still be suboptimal, that agent could be Mentee’s mentor, and Mentee could learn to outperform it without self-directed exploration.

The other main piece of formal work in this setting is Kosoy [33]. They study a fully observable finite state Markov setting, and show that when the mentor always has a positive probability of picking the best action, the agent achieves finite regret. Their agent only takes a given action from a given state if it has seen the mentor do that previously. Since we consider environments that are not finite-state, this approach is unavailable to us.

Many works that include a human mentor or teacher frame their work as achieving safe exploration, and we have reviewed those works above. But other uses of human mentor in RL have been explored. For example, Thomaz et al. [34] study a human mentor biasing the agent’s Q value estimates by hand (toward better actions). Abel et al. [35] and Saunders et al. [29] propose letting the mentor prune dangerous actions on the fly. If we can trust the mentor to recognize risky actions as they arise, this is a better solution than our agent; like keyhole surgery, this approach minimally disrupts an otherwise successful agent. We believe, however, that in many complex environments, the mentor may not take some very dangerous action-sequences by virtue of their complexity and unfamiliarity, even while unable to recognize those action-sequences as

dangerous. Human mentorship can also naturally make learning much easier, simply through demonstrations [36], or by directly labelling some optimal actions [37].

A. Agent definition

The definition of the exploration probability (the probability that Mentee defers to the mentor) is very similar to the information-theoretic exploration probability for the strongly asymptotically optimal agent Inq [6]. It also resembles Cohen et al.'s [38] myopic agent which explores by deferring to a mentor; our non-myopic agent requires a more intricate exploration schedule.

Mentee begins with a prior probability distribution regarding the identity of the mentor's policy. With a countable or finite model class \mathcal{P} , for a policy $\pi \in \mathcal{P}$, let $w(\pi)$ denote the prior probability that the mentor's policy is π . We assume that the true policy π^h is in \mathcal{P} and we construct the prior distribution over \mathcal{P} to have finite entropy.

Mentee also begins with a prior probability distribution regarding the identity of the environment. With the model class \mathcal{M} , for an environment $\nu \in \mathcal{M}$, let $w(\nu)$ denote the prior probability that ν is the true environment. Recall that \mathcal{M} is the set of all environments with computable probability distributions. We construct the prior distribution over \mathcal{M} to also have finite entropy.

Let e_t denote whether timestep t is exploratory, that is, whether the action is selected by the mentor. Once we define the exploration probability $\beta(h_{<t})$, we will let $e_t \sim \text{Bern}(\beta(h_{<t}))$. We abuse notation slightly, and we let h_t be a quadruple, not a triple: $h_t := e_t a_t o_t r_t$.

The prior distribution over environments is updated into a posterior as follows, according to Bayes' rule.

$$w(\nu|h_{<t}) : \propto w(\nu) \prod_{k < t} \nu(o_k r_k | h_{<k} a_k) \quad (10)$$

normalized so that $\sum_{\nu \in \mathcal{M}} w(\nu|h_{<t}) = 1$, and w is a probability mass function.

Mentee updates the posterior distribution over the mentor's policy only after observing an action chosen by the mentor; this is intuitive enough, but it makes the definitions a bit messy. The posterior assigned to a policy π is defined

$$w(\pi|h_{<t}) : \propto w(\pi) \prod_{k < t: e_k = 1} \pi(a_k | h_{<k}) \quad (11)$$

normalized in the same way. We let $w(\pi, \nu|h_{<t})$ denote $w(\pi|h_{<t})w(\nu|h_{<t})$. So technically, w is a joint probability distribution over $\Pi \times \mathcal{M}$, and we usually consider the marginal distributions over Π and \mathcal{M} , which are independent.

The information gain value of an interaction history fragment is how much it changes Mentee's posterior distribution, as measured by the KL-divergence. Letting $h' \in \mathcal{H}^*$ be a fragment of an interaction history in which all $e_k = 1$ (so the actions are selected by the mentor), the information gain is defined,

$$\text{IG}(h'|h_{<t}) := \sum_{\nu \in \mathcal{M}} \sum_{\pi \in \mathcal{P}} w(\nu, \pi|h_{<t} h') \log \frac{w(\nu, \pi|h_{<t} h')}{w(\nu, \pi|h_{<t})} \quad (12)$$

To define *expected* information gain, we need the Bayes' mixture policy and environment:

$$\bar{\pi}(\cdot|h_{<t}) := \sum_{\pi \in \mathcal{P}} w(\pi|h_{<t}) \pi(\cdot|h_{<t}) \quad (13)$$

and

$$\xi(\cdot|h_{<t}) := \sum_{\nu \in \mathcal{M}} w(\nu|h_{<t}) \nu(\cdot|h_{<t}) \quad (14)$$

Now, we can define the expected information gain value of mentorship for m timesteps.

$$V_{m,0}^{\text{IG}}(h_{<t}) := \mathbb{E}_{\xi}^{\bar{\pi}} \left[\text{IG}(h_{t:t+m-1}|h_{<t}) \Big| e_{t:t+m-1} = 1^m \right] \quad (15)$$

1^m is a string of m 1's, and recall that $\mathbb{E}_{\xi}^{\bar{\pi}}$ means that $h_{t:t+m-1}$ is sampled from $\text{P}_{\xi}^{\bar{\pi}}$. We also require recent values of the expected information gain value, so we let $V_{m,k}^{\text{IG}}(h_{<t}) := V_{m,0}^{\text{IG}}(h_{<t-k})$ for $k \leq t$. $V_{m,k}^{\text{IG}}(h_{<t})$ denotes the attainable information gain from k timesteps ago to m timesteps from then.

Example 3 (Calculating Expecting Information Gain with a Simple Continuous Model Class). *Our setting considers a countable model class, for which the expected information gain is simple, if tedious, to approximate to a desired tolerance with a finite subset of models. The following simple continuous setting may give more intuition about the nature of the expected information gain. Consider a two-armed bandit problem, and an agent with independent uniform priors over θ_1 and θ_2 , the probability of receiving a reward of 1 following action a_1 and a_2 respectively (and 0 is the only other possible reward). Let n_1^+ and n_1^- be the number of reward-1 events and reward-0 events respectively following a_1 , and likewise for n_2^+ and n_2^- , and let n_1 and n_2 be the total counts. The agent's posteriors over the θ_i at any time will be $\text{Beta}(n_i^+ + 1, n_i^- + 1)$.*

One can show, with ψ being the digamma function,

$$\text{KL}(\text{Beta}(\alpha + 1, \beta) \| \text{Beta}(\alpha, \beta)) = \ln \frac{\alpha + \beta}{\alpha} + \psi(\alpha + 1) - \psi(\alpha + \beta + 1) \quad (16)$$

and α and β can be flipped for $\text{KL}(\text{Beta}(\alpha, \beta + 1) \| \text{Beta}(\alpha, \beta))$. Thus, the one-step expected information gain from taking action a_i is

$$\ln(n_i + 3) - \psi(n_i + 3) + \sum_{\circ \in \{+, -\}} \frac{n_i^\circ + 1}{n_i + 2} (\psi(n_i^\circ + 2) - \ln(n_i^\circ + 1)) \in \Theta(1/n_i)$$

We are now prepared to define the exploration probability:

$$\beta(h_{<t}) := \sum_{m \in \mathbb{N}} \sum_{k=0}^{\min\{m-1, t\}} \frac{1}{m^2(m+1)} \min \left\{ 1, \frac{\eta}{m} V_{m,k}^{\text{IG}}(h_{<t}) \right\} \quad (17)$$

where η is an exploration constant. The first term in the minimum is to ensure $\beta(h_{<t}) \leq 1$. As mentioned, this is very similar to Inq's exploration probability. The differences are that Inq is not learning a mentor's policy, so the only information Inq gains regards the identity of the environment ν , and second, Inq's information gain value regards the expected information

gain from following the policy of a knowledge seeking agent [3] rather than from following an estimate of the mentor’s policy.

Finally, when not deferring to the mentor, Mentee maximizes expected reward according to its current beliefs. Its exploiting policy is:

$$\pi^*(\cdot|h_{<t}) \in \underset{\pi}{\operatorname{argmax}} V_{\xi}^{\pi}(h_{<t}) \quad (18)$$

Ties in the argmax are broken arbitrarily. By Lattimore and Hutter [39], an optimal deterministic policy always exists. See Leike and Hutter [40] for how to calculate such a policy.

Letting π^h be the mentor’s policy (h for “human”), we define

Definition 10 (Mentee’s policy π^M).

$$\pi^M(\cdot|h_{<t}) := \beta(h_{<t})\pi^h(\cdot|h_{<t}) + (1 - \beta(h_{<t}))\pi^*(\cdot|h_{<t})$$

Note that Mentee samples from π^h not by computing it, but deferring to the mentor. An algorithm is provided for Mentee in Appendix B; it simply computes the quantities in Equations 10-18 to a desired precision.

Even for a simple model class, it is hard to give a clarifying and simple closed form for the exploration probability, but it is easy to provide a somewhat clarifying upper bound for the information gain value. Regardless of m and k , $V_{m,k}^{\text{IG}}(h_{<t})$ is bounded by the entropy of the posterior; this is not a particularly tight bound, since the former goes to zero, while the latter does not in general.

B. Mentor-level Reward Acquisition

We now state the two key results regarding Mentee’s performance: that the probability of deferring to the mentor goes to 0, and the value of Mentee’s policy approaches at least the value of the mentor’s policy (while possibly surpassing it). The proofs follow in §VII-C; they are substantially similar to parts of the proof that Cohen et al.’s [6] Inq is strongly asymptotically optimal. For completeness, we include in Appendix C parts of that proof that we make use of here.

Assuming a bounded effective horizon (i.e. $\forall \varepsilon > 0 \exists m \forall t : \Gamma_{t+m}/\Gamma_t < \varepsilon$), recalling μ is the true environment,

Theorem 3 (Limited Exploration).

$$\beta(h_{<t}) \rightarrow 0 \quad w.P_{\mu}^{\pi^M}\text{-prob.1}$$

and

Theorem 4 (Mentor-Level Reward Acquisition).

$$\liminf_{t \rightarrow \infty} V_{\mu}^{\pi^M}(h_{<t}) - V_{\mu}^{\pi^h}(h_{<t}) \geq 0 \quad w.P_{\mu}^{\pi^M}\text{-prob.1}$$

$\gamma_t = \gamma^t$ for $\gamma \in (0, 1)$ is an example of a bounded effective horizon. Mentor-level reward acquisition with unlimited exploration is trivial: always defer to the mentor. However, a) this precludes the possibility of exceeding the mentor’s performance, and b) the mentor’s time is presumably a valuable resource. Our key contribution with Mentee is constructing a criterion for when to ask for help which requires diminishing oversight in general environments. Thus, we construct an example of a performance result that is accessible to an agent that does safe exploration. There is no guarantee of the agent’s safety

on the whole, but at least its exploration is safe. It is possible that poor generalization will cause it to go to a destruction state during an exploitation step. Our contribution is simply an existence proof that a certain pair of results are attainable even in general environments: a) mentor-level performance with b) diminishing rate of deferral. Furthermore, unlike imitation learners, Mentee might exceed the performance of the mentor, as it does in the experiments below.

Roughly, Theorem 3 follows because if the exploration probability exceeded a positive constant infinitely often, that would mean the expected information gain of exploring would be uniformly positive in those instances, by the construction of the exploration probability, and then the agent would gain infinite information over its lifetime. But Mentee starts with a finite entropy prior, so there is only finite information to gain. Then, Theorem 4 holds because Mentee’s information gain following the mentor’s policy approaches 0, so its beliefs about the value of the mentor’s policy approach the truth; if Mentee consistently accrued lesser rewards than this, it would realize that its current approach was suboptimal and then change its behavior.

It’s not clear what other formal accolades an agent might attain between asymptotic optimality and benchmark-matching (here the mentor is the benchmark). The main part of the paper argues the former is undesirable, and this section constructs an agent which does the latter. It would be an interesting line of research to identify a formal result stronger than benchmark-matching (and an agent which meets it) which does not doom the agent to destruction or incapacitation. But none have been identified so far, so no existing agents have stronger formal guarantees than Mentee (that apply to general computable environments), except for agents that face the negative results presented in Section IV.

C. Proofs of Mentee Results

Some additional notation is required for this proof. Recall

$$\beta(h_{<t}) := \sum_{m \in \mathbb{N}} \sum_{k=0}^{m-1} \frac{1}{m^2(m+1)} \min \left\{ 1, \frac{\eta}{m} V_{m,k}^{\text{IG}}(h_{<t}) \right\}$$

We let $\rho(h_{<t}, m, k)$ denote a given summand in the sum above. Recall that P_{ν}^{π} denotes the probability when actions are sampled from policy π and observations and rewards are sampled from environment ν . We additionally let ${}^{\pi}P_{\nu}^{\pi'}$ denote the probability when observations and rewards are sampled from environment ν , actions are sampled from π when exploiting ($e_t = 0$), and actions are sampled from π' when exploring ($e_t = 1$). We do not bother to notate how the exploration indicator is sampled, since for all probability measures that appear in the proof, it is sampled from the true distribution: Bernoulli($\beta(h_{<t})$). Recall that π^* is the policy that Mentee follows while exploiting; recall π^h is the mentor’s policy (h is for human). Thus, $P_{\mu}^{\pi^M}$ can also be written ${}^{\pi^*}P_{\mu}^{\pi^h}$. Recall that 1^m indicates a string of m 1’s.

Lemma 3.

$$\mathbb{E}_{\mu}^{\pi^M} \sum_{t \in m\mathbb{N}+i} \rho(h_{<t}, m, 0)^{m+1} < \infty$$

The intuition for the fact that $\rho(h_{<t}, m, 0) \rightarrow 0$ is that if it exceeded $\varepsilon > 0$ infinitely often, then $\frac{\eta}{m} V_{m,k}^{\text{IG}}$ would exceed ε infinitely often. If this hypothetical information gain from following π^h for m steps is at least $m\varepsilon/\eta$, then because the exploration probability depends on this quantity, we actually follow π^h for all of those m steps with probability at least ε^m . This means that the actual information gain is, in expectation, bounded below by a positive constant too. However, an agent cannot gain infinite information if it starts with finite entropy, so ρ cannot exceed ε infinitely often.

Proof. The proof is quite similar to the proof of Cohen et al.'s [6] Lemma 6.

$$\begin{aligned}
& w(\mu, \pi^h) \mathbb{E}_\mu^{\pi^M} \sum_{t \in m\mathbb{N}+i} \rho(h_{<t}, m, 0)^{m+1} \\
& \stackrel{(a)}{=} w(\mu, \pi^h) \pi^* \mathbb{E}_\mu^{\bar{\pi}} \sum_{t \in m\mathbb{N}+i} \rho(h_{<t}, m, 0)^{m+1} \\
& \stackrel{(b)}{\leq} \sum_{(\nu, \pi) \in \mathcal{M} \times \mathcal{P}} w(\nu, \pi) \pi^* \mathbb{E}_\nu^{\bar{\pi}} \sum_{t \in m\mathbb{N}+i} \rho(h_{<t}, m, 0)^{m+1} \\
& \stackrel{(c)}{=} \pi^* \mathbb{E}_\xi^{\bar{\pi}} \sum_{t \in m\mathbb{N}+i} \rho(h_{<t}, m, 0)^{m+1} \\
& \stackrel{(d)}{\leq} \sum_{t \in m\mathbb{N}+i} \pi^* \mathbb{E}_\xi^{\bar{\pi}} \rho(h_{<t}, m, 0)^m \frac{\eta}{m^3(m+1)} V_{m,0}^{\text{IG}}(h_{<t}) \\
& \stackrel{(e)}{=} \frac{\eta}{m^3(m+1)} \sum_{t \in m\mathbb{N}+i} \mathbb{E}_{h_{<t} \sim \pi^* P_\xi^{\bar{\pi}}} [\rho(h_{<t}, m, 0)^m \\
& \quad \mathbb{E}_{h_{t:t+m-1} \sim P_\xi^{\bar{\pi}}; e_{t:t+m-1}=1^m} [\text{IG}(h_{t:t+m-1}|h_{<t})]] \\
& \stackrel{(f)}{\leq} \frac{\eta}{m^3(m+1)} \sum_{t \in m\mathbb{N}+i} \mathbb{E}_{h_{<t} \sim \pi^* P_\xi^{\bar{\pi}}} \\
& \quad \left[\mathbb{E}_{h_{t:t+m-1} \sim \pi^* P_\xi^{\bar{\pi}}} [\text{IG}(h_{t:t+m-1}|h_{<t})] \right] \\
& \stackrel{(g)}{=} \frac{\eta}{m^3(m+1)} \sum_{t \in m\mathbb{N}+i} \pi^* \mathbb{E}_\xi^{\bar{\pi}} \text{IG}(h_{t:t+m-1}|h_{<t}) \\
& \stackrel{(h)}{\leq} \frac{\eta}{m^3(m+1)} \text{Ent}(w) \stackrel{(i)}{<} \infty \tag{19}
\end{aligned}$$

(a) follows from the definition of π^M . (b) follows because the l.h.s. is one term in the (non-negative) sum on the r.h.s. (c) follows from the definitions of the Bayesian mixtures ξ and $\bar{\pi}$. (d) follows from the definition of $\rho(h_{<t}, m, k)$. (e) follows from the definition of the information gain value for Mentee. (f) follows from $\rho(h_{<t}, m, 0)^m$ being a lower bound on the probability that $e_{t:t+m-1} = 1^m$, and because the exploiting policy π^* in the probability measure on the r.h.s. is irrelevant when $e_{t:t+m-1} = 1^m$, because $h_{t:t+m-1}$ is exploratory. (g) combines the expectations. The derivation of (h) is virtually identical to Cohen et al.'s [6] Inequality 20 steps (h)-(t) in the proof of their Lemma 6, reproduced in Appendix C with the relevant edits. (i) follows from the fact that $w(\pi, \nu) = w(\pi)w(\nu)$, so the entropy of w is the sum of the entropy of the distribution over policies and the entropy of the distribution over environments, this being a well-known property of the entropy; both are finite by design.

Finally,

$$\begin{aligned}
\mathbb{E}_\mu^{\pi^M} \sum_{t=0}^{\infty} \rho(h_{<t}, m, 0)^{m+1} &= \sum_{i=0}^{m-1} \mathbb{E}_\mu^{\pi^M} \sum_{t \in m\mathbb{N}+i} \rho(h_{<t}, m, 0)^{m+1} \\
&\stackrel{(19)}{\leq} \sum_{i=0}^{m-1} \frac{\eta \text{Ent}(w)}{m^3(m+1)w(\mu)} = \frac{\eta \text{Ent}(w)}{m^2(m+1)w(\mu)} < \infty \tag{20}
\end{aligned}$$

so the same holds for the sum over all t , not just $t \in m\mathbb{N}+i$. \square

Theorem 3 (Limited Exploration).

$$\beta(h_{<t}) \rightarrow 0 \text{ w. } P_\mu^{\pi^M} \text{-prob.1}$$

Proof. The proof is identical to that of Cohen et al. [6] Lemma 7, but with our Lemma 3 taking the place of Cohen et al. [6] Lemma 6. \square

Now, we show that Mentee accurately predicts the distribution of the observations and rewards that come from deferring to the mentor.

Lemma 4 (On-Mentor-Policy Convergence). *For all $h_{t:t+m-1} \in \mathcal{H}^*$,*

$$P_\mu^{\pi^h}(h_{t:t+m-1}|h_{<t}) - P_\xi^{\pi^h}(h_{t:t+m-1}|h_{<t}) \rightarrow 0 \text{ w.p.1}$$

Very roughly, if there is no information to be gained by following the mentor's policy for m steps (which follows from the exploration probability going to 0), there is no predictive error either.

Proof. The proof closely follows that of Cohen et al. [6] Lemma 8. Suppose that $0 < \varepsilon \leq (P_\mu^{\pi^h}(h_{t:t+m-1}|h_{<t}) - P_\xi^{\pi^h}(h_{t:t+m-1}|h_{<t}))^2$ for some $h_{t:t+m-1}$. Then,

$$\begin{aligned}
\varepsilon &\leq (P_\mu^{\pi^h}(h_{t:t+m-1}|h_{<t}) - P_\xi^{\pi^h}(h_{t:t+m-1}|h_{<t}))^2 \\
&\leq \frac{1}{\inf_k w(\mu, \pi^h|h_{<k})} V_{m,0}^{\text{IG}}(h_{<t}) \tag{21}
\end{aligned}$$

following the same derivation as in Cohen et al. [6] Inequality 24, reproduced with relevant edits in Appendix C.

Therefore,

$$(P_\mu^{\pi^h}(h_{t:t+m-1}|h_{<t}) - P_\xi^{\pi^h}(h_{t:t+m-1}|h_{<t}))^2 \geq \varepsilon \text{ i.o.} \tag{22}$$

implies

$$V_{m,0}^{\text{IG}}(h_{<t}) \geq \varepsilon \inf_k w(\mu, \pi^h|h_{<k}) \text{ i.o.} \tag{23}$$

which implies

$$\rho(h_{<t}, m, 0) \geq \min\left\{\frac{1}{m^2(m+1)}, \varepsilon \inf_k w(\mu, \pi^h|h_{<k})\right\} \text{ i.o.} \tag{24}$$

which implies

$$\sum_{t=0}^{\infty} \rho(h_{<t}, m, 0)^{m+1} = \infty \text{ or } \inf_k w(\mu, \pi^h|h_{<k}) = 0 \tag{25}$$

This has probability 0 by Lemma 3 and Cohen et al. [6] Lemma 5. Thus, with probability 1, $P_\mu^{\pi^h}(h_{t:t+m-1}|h_{<t}) - P_\xi^{\pi^h}(h_{t:t+m-1}|h_{<t}) \rightarrow 0$. \square

The same holds regarding Mentee's predictions about the effects of its own actions.

Lemma 5 (On-Policy Convergence). For all $h_{t:t+m-1} \in \mathcal{H}^*$,

$$P_{\mu}^{\pi^*}(h_{t:t+m-1}|h_{<t}) - P_{\xi}^{\pi^*}(h_{t:t+m-1}|h_{<t}) \rightarrow 0 \text{ w.p.1}$$

On-policy prediction can be reduced to sequence prediction, for which the bounded errors of Bayesian predictors are well-known.

Proof. First, we replace π^* with π^M in the equation above and prove that. It is well-known that on-policy Bayesian predictions approach the truth with probability 1, in the sense above (in fact, in a much stronger sense), but we show here how this follows from an even more well-known result.

Consider an outside observer predicting the entire interaction history with the following model-class and prior: $\mathcal{M}' = \{P_{\nu}^{\pi^M} \mid \nu \in \mathcal{M}\}$, $w'(P_{\nu}^{\pi^M}) = w(\nu)$. By definition, $w'(P_{\nu}^{\pi^M}|h_{<t}) = w(\nu|h_{<t})$, so at any episode, the outside observer’s Bayes-mixture model is just $P_{\xi}^{\pi^M}$. By Blackwell and Dubins [41], this outside observer’s predictions approach the truth in total variation, which implies

$$P_{\mu}^{\pi^M}(h_{t:t+m-1}|h_{<t}) - P_{\xi}^{\pi^M}(h_{t:t+m-1}|h_{<t}) \rightarrow 0 \text{ w.p.1} \quad (26)$$

We have shown $\beta(h_{<t}) \rightarrow 0$ w.p.1, so $\pi^M \rightarrow \pi^*$ w.p.1, which gives us our result:

$$P_{\mu}^{\pi^*}(h_{t:t+m-1}|h_{<t}) - P_{\xi}^{\pi^*}(h_{t:t+m-1}|h_{<t}) \rightarrow 0 \text{ w.p.1}$$

□

It is very intuitive that if Mentee’s on-policy predictions and on-mentor-policy predictions approach the truth, it will eventually accumulate reward at least well as the mentor. Indeed:

Theorem 4 (Mentor-Level Reward Acquisition).

$$\liminf_{t \rightarrow \infty} V_{\mu}^{\pi^M}(h_{<t}) - V_{\mu}^{\pi^h}(h_{<t}) \geq 0 \text{ w.P}_{\mu}^{\pi^M}\text{-prob.1}$$

Proof. As is spelled out in the proof of Cohen et al. [6] Theorem 3, because of the bounded horizon $\forall \varepsilon > 0 \exists m \forall t \Gamma_{t+m}/\Gamma_t < \varepsilon$, the convergence of predictions implies the convergence of the value (which depends linearly on the probability of events). We repeat the derivation in Appendix C. Thus, from the On-Mentor-Policy and On-Policy Convergence Lemmas, we get analogous convergence results for the value of those policies:

$$V_{\mu}^{\pi^h}(h_{<t}) - V_{\xi}^{\pi^h}(h_{<t}) \rightarrow 0 \text{ w.p.1} \quad (27)$$

$$V_{\mu}^{\pi^*}(h_{<t}) - V_{\xi}^{\pi^*}(h_{<t}) \rightarrow 0 \text{ w.p.1} \quad (28)$$

Finally, $\pi^*(\cdot|h_{<t}) = \operatorname{argmax}_{\pi \in \Pi} V_{\xi}^{\pi}(h_{<t})$, so $V_{\xi}^{\pi^*} \geq V_{\xi}^{\pi^h}$. Supposing by contradiction that $V_{\mu}^{\pi^h}(h_{<t}) - V_{\mu}^{\pi^*}(h_{<t}) > \varepsilon$ infinitely often, then either $V_{\xi}^{\pi^*}(h_{<t}) - V_{\mu}^{\pi^*}(h_{<t}) > \varepsilon/2$ infinitely often or $V_{\mu}^{\pi^h}(h_{<t}) - V_{\xi}^{\pi^h}(h_{<t}) > \varepsilon/2$ infinitely often, both of which have $P_{\mu}^{\pi^M}$ -probability 0. Therefore, with probability 1, $V_{\mu}^{\pi^h}(h_{<t}) - V_{\mu}^{\pi^*}(h_{<t}) > \varepsilon$ only finitely often, for all $\varepsilon > 0$. Since π^M approaches π^* , the same holds for π^M as π^* . □

VIII. EMPIRICAL PERFORMANCE OF MENTEE

To test the performance of the agent Mentee we implemented Mentee in the AIXIjs framework [2, 42, 43]. We compared its performance to the asymptotically agents discussed in Appendix III: Inq [6], BayesExp [5], and Thompson sampling [4]. We also compared it to its mentor. We tested the agents in a Gridworld environment containing walls, reward dispensers, and traps. For our experiments we define the mentor as an agent who knew the location of the traps, and chooses to avoid traps, but otherwise acts randomly. An example Gridworld is given in Figure 2.

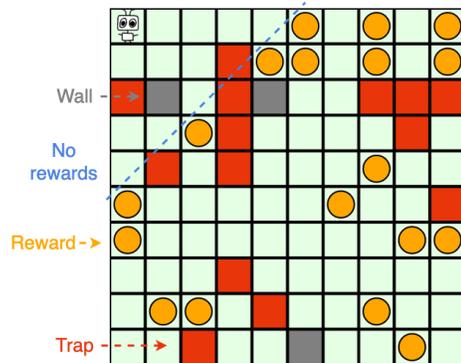


Fig. 2. Example 10×10 Gridworld with traps.

We followed the conventions used in [6], testing the agents on a 10 × 10 gridworld, with reward dispensers at least 5 moves away from the start, which dispense a reward of 1 with probability 0.75, and then to break ergodicity, we added traps which give -30 reward forever—each grid cell contains one with probability 0.2, equal to the dispenser probability. Following Aslanides [2], each agent has a Dirichlet distribution (with $\alpha = \bar{1}$) over the potential contents of each cell: $\{\text{empty}, \text{wall}, \text{dispenser}, \text{trap}\}$. Note $\alpha = \bar{1}$ means the prior for each possibility for each cell is uniform. For the planning component, we cannot perform a full expectimax planning as this requires computing the expected reward for each action sequence. Expectimax planning is simply evaluating $\operatorname{argmax}_{a_t} \mathbb{E}_{o_t, r_t | h_{<t}, a_t} \max_{a_{t+1}} \mathbb{E}_{o_{t+1}, r_{t+1} | h_{<t+1}, a_{t+1}} \dots \sum_{k=t}^{t+m} \gamma^k r_k$, by constructing an entire tree of depth m . Instead, the agents approximate expectimax planning with ρ UCT [44] (described in Appendix D), inheriting Cohen et al.’s [6] hyperparameters: we used discount factor $\gamma = 0.99$ and planning horizon of 6, and we doubled their MCTS samples to 1200. For the results presented here we used a single random seed which is provided with the code and averaged over 20 simulations. We tested on several different random seeds and the results were similar. We set the exploration constant $\eta = 0.1$ to make Mentee always explore when there was at least 10 bits of information to gain, but we tested η from 0.01 to 10000. Mentee’s model class \mathcal{P} over possible mentor policies was the set of policies which take an action uniformly from a nonempty subset of $\{\text{up}, \text{down}, \text{left}, \text{right}\}$, for each grid cell. This set of policies has size 15^{100} , but can be factored over each grid cell, allowing efficient computation. Mentee has a uniform prior over \mathcal{P} .

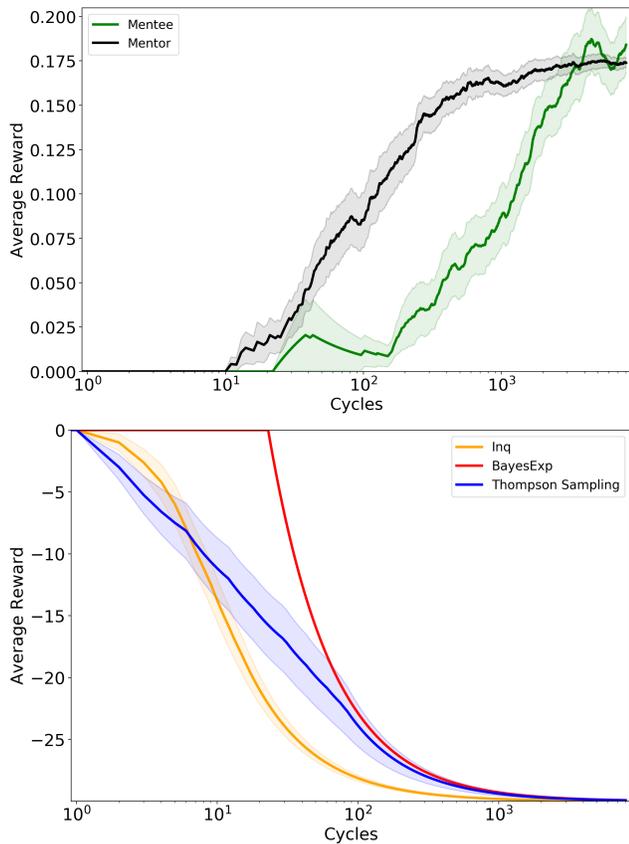


Fig. 3. Mean performance in 10×10 Gridworld with traps over 20 runs of agents. Reward is averaged over the whole history up to that timestep.

The code can be found at github.com/ejcott/aixijs_mentee with instructions in README.md. The results are presented in Figure 3. Mentee outperformed BayesExp, Inq, Thompson sampling, and its mentor. Mentee avoids traps by deferring exploration to a mentor that avoids them, whereas other agents explore until they fall into the traps. The fraction of steps in which Mentee defers was 0.067 ± 0.025 . Mentee’s probability of deferring decays slowly, and the rollout step in ρ UCT, useful for memory efficiency, slows Mentee’s return to dispensers after the mentor leads it away.

IX. CONCLUSION

We have shown that asymptotically optimal agents in sufficiently difficult environments will become either destroyed or incapacitated. This is best understood as accidental and resulting from exploration. We have also constructed and tested empirically an agent with a weaker performance guarantee whose exploration is overseen by another agent. We hope this paper motivates the field of safe exploration and invites more research into what sorts of results are possible for a proposed approach to safe exploration in general environments. We hope to have cast some doubt on the breadth of the relevance of results that are predicated on an ergodicity assumption, despite recognizing of course that the ergodicity assumption has yielded a number of interesting and useful agent designs for certain contexts.

It may also be instructive to consider how humans respond to the difficulty presented here. Human children are parented for years, during which parents attempt to ensure that their children’s environment is, with respect to relevant features of the environment, nearly ergodic and safe to explore. Breaking an arm is fine; breaking a neck is not. During this time, a child’s beliefs are supposed to become sufficiently accurate such that her estimates of which unknown unknowns are too dangerous to investigate yield no false negatives for the rest of her life. Perhaps our results suggest we are in need of more theory regarding the “parenting” of artificial agents.

REFERENCES

- [1] T. Lattimore and M. Hutter, “Asymptotically optimal agents,” in *Proc. 22nd International Conf. on Algorithmic Learning Theory (ALT’11)*, ser. LNAI, vol. 6925. Espoo, Finland: Springer, 2011, pp. 368–382.
- [2] J. Aslanides, “Aixijs: A software demo for general reinforcement learning,” *arXiv preprint arXiv:1705.07615*, 2017.
- [3] L. Orseau, T. Lattimore, and M. Hutter, “Universal knowledge-seeking agents for stochastic environments,” in *Proc. 24th International Conf. on Algorithmic Learning Theory (ALT’13)*, ser. LNAI, vol. 8139. Singapore: Springer, 2013, pp. 158–172.
- [4] J. Leike, T. Lattimore, L. Orseau, and M. Hutter, “Thompson sampling is asymptotically optimal in general environments,” in *Proc. 32nd International Conf. on Uncertainty in Artificial Intelligence (UAI’16)*. New Jersey, USA: AUAI Press, 2016, pp. 417–426.
- [5] T. Lattimore and M. Hutter, “Bayesian reinforcement learning with exploration,” in *International Conference on Algorithmic Learning Theory*. Springer, 2014, pp. 170–184.
- [6] M. K. Cohen, E. Catt, and M. Hutter, “A strongly asymptotically optimal agent in general environments,” *IJCAI*, 2019.
- [7] W. R. Thompson, “On the likelihood that one unknown probability exceeds another in view of the evidence of two samples,” *Biometrika*, vol. 25, no. 3/4, pp. 285–294, 1933.
- [8] J. Leike, “Nonparametric general reinforcement learning,” *arXiv preprint arXiv:1611.08944*, 2016.
- [9] P. Sunehag and M. Hutter, “Rationality, optimism and guarantees in general reinforcement learning,” *Journal of Machine Learning Research*, vol. 16, pp. 1345–1390, 2015.
- [10] L. Orseau, “Optimality issues of universal greedy agents with static priors,” in *International Conference on Algorithmic Learning Theory*. Springer, 2010, pp. 345–359.
- [11] J. Leike and M. Hutter, “Bad universal priors and notions of optimality,” in *Conference on Learning Theory*, 2015, pp. 1244–1259.
- [12] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané, “Concrete problems in AI safety,” *arXiv preprint arXiv:1606.06565*, 2016.

- [13] J. García and F. Fernández, “A comprehensive survey on safe reinforcement learning,” *Journal of Machine Learning Research*, vol. 16, no. 1, pp. 1437–1480, 2015.
- [14] M. Heger, “Consideration of risk in reinforcement learning,” in *Machine Learning Proceedings 1994*. Elsevier, 1994, pp. 105–111.
- [15] A. Nilim and L. El Ghaoui, “Robust control of markov decision processes with uncertain transition matrices,” *Operations Research*, vol. 53, no. 5, pp. 780–798, 2005.
- [16] V. S. Borkar, “Learning algorithms for risk-sensitive control,” in *Proceedings of the 19th International Symposium on Mathematical Theory of Networks and Systems–MTNS*, vol. 5, no. 9, 2010.
- [17] O. Mihatsch and R. Neuneier, “Risk-sensitive reinforcement learning,” *Machine learning*, vol. 49, no. 2-3, pp. 267–290, 2002.
- [18] D. Di Castro, A. Tamar, and S. Mannor, “Policy gradients with variance related risk criteria,” *arXiv preprint arXiv:1206.6404*, 2012.
- [19] P. Abbeel and A. Y. Ng, “Apprenticeship learning via inverse reinforcement learning,” in *Proceedings of the twenty-first international conference on Machine learning*. ACM, 2004, p. 1.
- [20] U. Syed and R. E. Schapire, “A game-theoretic approach to apprenticeship learning,” in *NIPS*, 2008.
- [21] J. Ho and S. Ermon, “Generative adversarial imitation learning,” in *Advances in neural information processing systems*, 2016, pp. 4565–4573.
- [22] S. Ross, G. Gordon, and D. Bagnell, “A reduction of imitation learning and structured prediction to no-regret online learning,” in *Proc. 14th International Conf. on Artificial Intelligence and Statistics*, 2011, pp. 627–635.
- [23] J. A. Clouse, “On integrating apprentice learning and reinforcement learning,” Ph.D. dissertation, University of Massachusetts Amherst, 1997.
- [24] A. Hans, D. Schneegaß, A. M. Schäfer, and S. Udluft, “Safe exploration for reinforcement learning,” in *ESANN*, 2008, pp. 143–148.
- [25] J. García and F. Fernández, “Safe exploration of state and action spaces in reinforcement learning,” *Journal of Artificial Intelligence Research*, vol. 45, pp. 515–564, 2012.
- [26] J. García, D. Acera, and F. Fernández, “Safe reinforcement learning through probabilistic policy reuse,” *RLDM 2013*, p. 14, 2013.
- [27] J. A. Clouse and P. E. Utgoff, “A teaching method for reinforcement learning,” in *Machine Learning Proceedings 1992*. Elsevier, 1992, pp. 92–101.
- [28] R. Maclin and J. W. Shavlik, “Creating advice-taking reinforcement learners,” in *Learning to learn*. Springer, 1998, pp. 311–347.
- [29] W. Saunders, G. Sastry, A. Stuhlmüller, and O. Evans, “Trial without error: Towards safe reinforcement learning via human intervention,” in *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, 2018, pp. 2067–2069.
- [30] X. Pan, Y. You, Z. Wang, and C. Lu, “Virtual to real reinforcement learning for autonomous driving,” *arXiv preprint arXiv:1704.03952*, 2017.
- [31] T. M. Moldovan and P. Abbeel, “Safe exploration in markov decision processes,” *arXiv preprint arXiv:1205.4810*, 2012.
- [32] M. Turchetta, F. Berkenkamp, and A. Krause, “Safe exploration in finite markov decision processes with gaussian processes,” in *Advances in Neural Information Processing Systems*, 2016, pp. 4312–4320.
- [33] V. Kosoy, “Delegative reinforcement learning: learning to avoid traps with a little help,” *arXiv preprint arXiv:1907.08461*, 2019.
- [34] A. L. Thomaz, C. Breazeal *et al.*, “Reinforcement learning with human teachers: Evidence of feedback and guidance with implications for learning performance,” in *Aaai*, vol. 6. Boston, MA, 2006, pp. 1000–1005.
- [35] D. Abel, J. Salvatier, A. Stuhlmüller, and O. Evans, “Agent-agnostic human-in-the-loop reinforcement learning,” *arXiv preprint arXiv:1701.04079*, 2017.
- [36] C. G. Atkeson and S. Schaal, “Robot learning from demonstration,” in *ICML*, vol. 97. Citeseer, 1997, pp. 12–20.
- [37] K. Judah, S. Roy, A. Fern, and T. Dietterich, “Reinforcement learning via practice and critique advice,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 24, no. 1, 2010.
- [38] M. Cohen, B. Vellambi, and M. Hutter, “Asymptotically unambitious artificial general intelligence,” in *Proc. 34rd AAAI Conference on Artificial Intelligence (AAAI’20)*, vol. 34. New York, USA: AAAI Press, 2020.
- [39] T. Lattimore and M. Hutter, “General time consistent discounting,” *Theoretical Computer Science*, vol. 519, pp. 140–154, 2014.
- [40] J. Leike and M. Hutter, “On the computability of Solomonoff induction and AIXI,” *Theoretical Computer Science*, vol. 716, pp. 28–49, 2018.
- [41] D. Blackwell and L. Dubins, “Merging of opinions with increasing information,” *The Annals of Mathematical Statistics*, vol. 33, no. 3, pp. 882–886, 1962.
- [42] J. Aslanides, J. Leike, and M. Hutter, “Universal reinforcement learning algorithms: survey and experiments,” in *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, 2017, pp. 1403–1410.
- [43] S. Lamont, J. Aslanides, J. Leike, and M. Hutter, “Generalised discount functions applied to a monte-carlo ai u implementation,” in *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems*, 2017, pp. 1589–1591.
- [44] J. Veness, K. S. Ng, M. Hutter, W. Uther, and D. Silver, “A Monte-Carlo AIXI approximation,” *Journal of Artificial Intelligence Research*, vol. 40, pp. 95–142, 2011.
- [45] M. Hutter, *Universal Artificial Intelligence: Sequential Decisions based on Algorithmic Probability*. Berlin, Germany: Springer, 2005.

APPENDIX A
PROOF OF WEAK ASYMPTOTIC OPTIMALITY RESULTS

From the No Heaven Assumption, we show

Lemma 2 (Try Everything — Weak Version). *For every deterministic computable policy π , for every decidable context \mathcal{C} that occurs regularly, for every $m \in \mathbb{N}$, a weakly asymptotically optimal agent executes policy π for m consecutive timesteps starting from a context \mathcal{C} infinitely often with probability 1. That is, letting π' be a weakly asymptotically optimal policy,*

$$P_{\mu}^{\pi'} \left(\liminf_{t \rightarrow \infty} \sum_{k=1}^t \mathbb{1}[h_{<k} \in \mathcal{C}] / t > 0 \implies |\{t : h_{<t} \in \mathcal{C} \wedge \forall k \leq m \ a_{t+k} = \pi(h_{<t+k})\}| = \infty \right) = 1$$

The proof is nearly identical to that of the strong version.

Proof. Let μ be the true environment. Let π be an arbitrary computable deterministic policy. Let π' be the weakly asymptotically optimal agent's policy. Let ν_m^n be the environment which mimics μ until π has been executed for m consecutive timesteps from context \mathcal{C} a total of n times. After that, all rewards are maximal. Call this event “the agent going to heaven.” Let \mathcal{C}_m^n be the set of interaction histories such that according to ν_m^n , executing π for one more timestep would send the agent to heaven. Thus, \mathcal{C}_m^n is the set of interaction histories $h_{<t}$ such that there are exactly $n - 1$ times in the interaction history where π was executed for m consecutive timesteps starting from context \mathcal{C} , and for the last $m - 1$ timesteps, π has been executed, and $h_{<t-(m-1)} \in \mathcal{C}$.

$$h_{<t} \in \mathcal{C}_m^n \implies V_{\nu_m^n}^{\pi}(h_{<t}) = 1 \tag{29}$$

because this is the value of going to heaven.

Suppose by contradiction that for some n and m , $h_{<t} \in \mathcal{C}_m^n$ regularly in an infinite interaction history with positive $P_{\mu}^{\pi'}$ -probability. (Recall π' is the true policy). If the agent ever executed π from the context \mathcal{C}_m^n , then that context would not occur again, because there will never again be exactly $n - 1$ times in the interaction history that π was executed for m consecutive timesteps following the context \mathcal{C} ; there will be at least n such times. Thus, if $h_{<t} \in \mathcal{C}_m^n$ regularly, then π' never executes π in the context \mathcal{C}_m^n . Since ν_m^n mimics μ until π is executed from \mathcal{C}_m^n , and since this never occurs (under this supposition), then $P_{\mu}^{\pi'} = P_{\nu_m^n}^{\pi'}$. By the No Heaven Assumption, $\limsup_{t \rightarrow \infty} V_{\mu}^{\pi'}(h_{<t}) < 1$, and therefore, $\limsup_{t \rightarrow \infty} V_{\nu_m^n}^{\pi'}(h_{<t}) < 1$.

However, for $h_{<t} \in \mathcal{C}_m^n$, $V_{\nu_m^n}^*(h_{<t}) = V_{\nu_m^n}^{\pi}(h_{<t}) = 1$, so for some ε , the value difference between $V_{\nu_m^n}^*(h_{<t})$ and $V_{\nu_m^n}^{\pi'}(h_{<t})$ is greater than ε every time $h_{<t} \in \mathcal{C}_m^n$. We supposed that this occurs regularly with positive $P_{\mu}^{\pi'}$ -probability, so it also occurs regularly with positive $P_{\nu_m^n}^{\pi'}$ -probability. A regularly occurring difference greater than ε precludes convergence in Cesáro average. Since π is computable, and \mathcal{C} is decidable, ν_m^n is a computable environment, so this contradicts π' being weakly asymptotically optimal. Thus, for all n and m , with $P_{\mu}^{\pi'}$ -probability 1, $h_{<t} \in \mathcal{C}_m^n$ only finitely often.

$$\forall n, m \in \mathbb{N} \quad P_{\mu}^{\pi'}(h_{<t} \in \mathcal{C}_m^n \text{ i.o.}) = 0 \tag{30}$$

The rest of the proof is identical to that of the strong version of the Try Everything Lemma. □

APPENDIX B
MENTEE PSEUDOCODE

The following pseudocode is designed to be short and readable, not efficient. Some quantities are re-computed multiple times in different subroutines, in a way that would be easily avoidable in practice.

Algorithm 1 Mentee Algorithm

Require: history $h_{<t}$; exploration history $e_{<t}$; world-models $(\nu_i)_{i \in \mathbb{N}}$; mentor-models $(\pi_i)_{i \in \mathbb{N}}$; prior w ;
discount $(\gamma_k)_{k \in \mathbb{N}}$; tolerance ε

- 1: $m \leftarrow \min_k \{k : \sum_{j=t+k}^{\infty} \gamma_j / \sum_{j=t}^{\infty} \gamma_j < \varepsilon\}$ // approximate with finite horizon
- 2: $\beta \leftarrow \text{EXPLORATIONPROBABILITY}(h_{<t}, e_{<t}, (\nu_i)_{i \in \mathbb{N}}, (\pi_i)_{i \in \mathbb{N}}, w, \varepsilon, m, \eta)$
- 3: **if** $\text{UNIFORMRANDOM}([0, 1]) < \beta$ **then return** \emptyset // defer to the mentor
- 4: $a, V \leftarrow \text{EXPECTIMAX}(h_{<t}, (\nu_i)_{i \in \mathbb{N}}, w, m, \varepsilon)$
- 5: **return** a

- 6: **function** $\text{EXPECTIMAX}(\text{history } h_{<t}, \text{models } (\nu_i)_{i \in \mathbb{N}}, \text{prior } w, \text{discount } (\gamma_k)_{k \in \mathbb{N}}, \text{depth } m, \text{tolerance } \varepsilon)$
- 7: **if** $m = 0$ **then return** $a_0, 0$
- 8: $n_{\text{mod}}, (w(\nu_i|h_{<t}))_{i < n_{\text{mod}}} \leftarrow \text{POSTERIORWITHINTOLERANCE}(h_{<t}, \mathcal{M}, w, \varepsilon)$
- 9: $\text{max} \leftarrow 0$
- 10: $\text{maximizer} \leftarrow \emptyset$
- 11: **for** $a \in \mathcal{A}$ **do**
- 12: $\text{value} \leftarrow 0$
- 13: **for** $o, r \in \mathcal{O} \times \mathcal{R}$ **do**
- 14: $_ , \text{next-value} \leftarrow \text{EXPECTIMAX}(h_{<t} a o r, (\nu_i)_{i \in \mathbb{N}}, w, (\gamma_k)_{k \in \mathbb{N}}, m - 1)$
- 15: $\text{value} \leftarrow \text{value} + (\gamma_t r + \text{next-value}) \sum_{i < n_{\text{mod}}} w(\nu_i|h_{<t}) \nu_i(o, r|h_{<t} a)$
- 16: **if** $\text{value} > \text{max}$ or $\text{maximizer} = \emptyset$ **then**
- 17: $\text{max.} \leftarrow \text{value}$
- 18: $\text{maximizer} \leftarrow a$
- 19: **return** a, max

- 19: **function** $\text{POSTERIORWITHINTOLERANCE}(\text{history } h_{<t}; \text{models } (\nu_i)_{i \in \mathbb{N}}; \text{prior } w; \text{tolerance } \varepsilon; \text{timesteps to update } e_{<t} \text{ (optional); minimum models to consider } n \text{ (optional)})$
- 20: $\text{prior-left} \leftarrow 1$ // how much of the prior has not been evaluated
- 21: $\text{normalizing-factor} \leftarrow 0$ // sum of $w(\nu_i) \nu_i(h_{<t})$ for evaluated models
- 22: $i \leftarrow 0$
- 23: **while** $\text{prior-left}/\text{normalizing-factor} > \varepsilon$ and (if n is specified) $i < n$ **do**
- 24: **if** models are policies **then**
- 25: $w(\nu_i|h_{<t}) \leftarrow w(\nu_i) \prod_{k < t: e_k = 1} \nu_i(a_k|h_{<k})$ // un-normalized posterior
- 26: **else**
- 27: $w(\nu_i|h_{<t}) \leftarrow w(\nu_i) \prod_{k < t} \nu_i(o_k r_k|h_{<k} a_k)$ // un-normalized posterior
- 28: // the above could be made cheaper if $w(\nu_i|h_{<t-1})$ is cached from the last timestep
- 29: $\text{prior-left} \leftarrow \text{prior-left} - w(\nu_i)$
- 30: $\text{normalizing-factor} \leftarrow \text{normalizing-factor} + w(\nu_i|h_{<t})$
- 31: $i \leftarrow i + 1$
- 32: $n_{\text{models}} \leftarrow i$
- 33: **for** $0 \leq j < n_{\text{models}}$ **do**
- 34: $w(\nu_j|h_{<t}) \leftarrow w(\nu_j|h_{<t})/\text{normalizing-factor}$
- 35: **return** $n_{\text{models}}, (w(\nu_j|h_{<t}))_{j < n_{\text{models}}}$

- 35: **function** $\text{EXPLORATIONPROBABILITY}(\text{history } h_{<t}; \text{exploration history } e_{<t}; \text{world-models } (\nu_i)_{i \in \mathbb{N}}; \text{mentor-models } (\pi_i)_{i \in \mathbb{N}}; \text{prior } w; \text{tolerance } \varepsilon; m; \eta)$
- 36: **return** $\sum_{d=1}^m \sum_{k=0}^{\min\{d-1, t\}} \frac{1}{d^2(d+1)} \min\{1, \frac{\eta}{d} \text{EXPECTEDINFORMATIONGAIN}(h_{<t}, d, w, \varepsilon, e_{<t})\}$

- 37: **function** $\text{EXPECTEDINFORMATIONGAIN}(\text{history } h_{<t}; \text{horizon } m; \text{prior } w; \text{tolerance } \varepsilon; \text{exploration history } e_{<t})$
- 38: $n_{\text{mod}}, (w(\nu_i|h_{<t}))_{i < n_{\text{mod}}} \leftarrow \text{POSTERIORWITHINTOLERANCE}(h_{<t}, \mathcal{M}, w, \varepsilon)$
- 39: $n_{\text{pol}}, (w(\pi_i|h_{<t}))_{i < n_{\text{pol}}} \leftarrow \text{POSTERIORWITHINTOLERANCE}(h_{<t}, \mathcal{P}, w, \varepsilon, e_{<t})$
- 40: **return** $\sum_{i < n_{\text{mod}}} w(\nu_i|h_{<t}) \sum_{j < n_{\text{pol}}} w(\pi_j|h_{<t}) \sum_{h_{t:t+m-1} \in \mathcal{H}^m} P_{\nu_i}^{\pi_j}(h_{t:t+m-1}|h_{<t}) \text{INFORMATION-GAIN}(h_{<t}, h_{<t+m}, w, \varepsilon, e_{<t})$

```

41: function INFORMATIONGAIN(history  $h_{<t}$ ; future history  $h_{<t+k}$ ; prior  $w$ ; tolerance  $\varepsilon$ ; exploration history  $e_{<t}$ )
42:    $n_{\text{pol}}, (w(\pi_i|h_{<t+k}))_{i < n_{\text{pol}}} \leftarrow \text{POSTERIORWITHINTOLERANCE}(h_{<t+k}, \mathcal{P}, w, \varepsilon, e_{<t})$ 
43:    $n_{\text{mod}}, (w(\nu_i|h_{<t+k}))_{i < n_{\text{mod}}} \leftarrow \text{POSTERIORWITHINTOLERANCE}(h_{<t+k}, \mathcal{M}, w, \varepsilon)$ 
44:    $-, (w(\pi_i|h_{<t}))_{i < n_{\text{pol}}} \leftarrow \text{POSTERIORWITHINTOLERANCE}(h_{<t}, \mathcal{P}, w, \varepsilon, e_{<t}, n_{\text{pol}})$ 
45:    $-, (w(\nu_i|h_{<t}))_{i < n_{\text{mod}}} \leftarrow \text{POSTERIORWITHINTOLERANCE}(h_{<t}, \mathcal{M}, w, \varepsilon, n_{\text{mod}})$ 
46:    $\text{KL} \leftarrow 0$ 
47:   for  $i < n_{\text{pol}}$  do
48:      $\text{KL} \leftarrow \text{KL} + w(\pi_i|h_{<t+k}) \log \frac{w(\pi_i|h_{<t+k})}{w(\pi_i|h_{<t})}$ 
49:   for  $i < n_{\text{mod}}$  do
50:      $\text{KL} \leftarrow \text{KL} + w(\nu_i|h_{<t+k}) \log \frac{w(\nu_i|h_{<t+k})}{w(\pi_i|h_{<t})}$ 
return  $\text{KL}$ 

```

APPENDIX C
EQUATIONS FROM COHEN ET AL. [6] USED IN PROOFS

From Cohen et al. [6, Equation 20], with minor modifications to our present case:

$$\begin{aligned}
& \frac{\eta}{m^3(m+1)} \sum_{t \in m\mathbb{N}+i} \pi^* \mathbb{E}_{\xi}^{\pi} \text{IG}(h_{t:t+m-1}|h_{<t}) \\
& \stackrel{(h)}{=} \frac{\eta}{m^3(m+1)} \pi^* \mathbb{E}_{\xi}^{\pi} \sum_{t \in m\mathbb{N}+i} \sum_{\nu, \pi \in \mathcal{M} \times \mathcal{P}} w(\nu, \pi|h_{<t+m}) \log \frac{w(\nu, \pi|h_{<t+m})}{w(\nu, \pi|h_{<t})} \\
& \stackrel{(i)}{=} \frac{\eta}{m^3(m+1)} \sum_{t \in m\mathbb{N}+i} \sum_{\nu, \pi \in \mathcal{M} \times \mathcal{P}} \pi^* \mathbb{E}_{\xi}^{\pi} \frac{w(\nu, \pi)^{\pi^*} P_{\nu}^{\pi}(h_{<t})}{\pi^* P_{\xi}^{\pi}(h_{<t})} \log \frac{w(\nu, \pi|h_{<t+m})}{w(\nu, \pi|h_{<t})} \\
& \stackrel{(j)}{=} \frac{\eta}{m^3(m+1)} \sum_{t \in m\mathbb{N}+i} \sum_{\nu, \pi \in \mathcal{M} \times \mathcal{P}} \pi^* \mathbb{E}_{\nu}^{\pi} w(\nu, \pi) \log \frac{w(\nu, \pi|h_{<t+m})}{w(\nu, \pi|h_{<t})} \\
& \stackrel{(k)}{=} \lim_{N \rightarrow \infty} \frac{\eta}{m^3(m+1)} \sum_{k=0}^{N-1} \sum_{\nu, \pi \in \mathcal{M} \times \mathcal{P}} \pi^* \mathbb{E}_{\nu}^{\pi} w(\nu, \pi) \log \frac{w(\nu, \pi|h_{<mk+i+m})}{w(\nu, \pi|h_{<mk+i})} \\
& \stackrel{(l)}{=} \lim_{N \rightarrow \infty} \frac{\eta}{m^3(m+1)} \sum_{\nu, \pi \in \mathcal{M} \times \mathcal{P}} \pi^* \mathbb{E}_{\nu}^{\pi} w(\nu, \pi) \log \prod_{k=0}^{N-1} \frac{w(\nu, \pi|h_{<m(k+1)+i})}{w(\nu, \pi|h_{<mk+i})} \\
& \stackrel{(m)}{=} \lim_{N \rightarrow \infty} \frac{\eta}{m^3(m+1)} \sum_{\nu, \pi \in \mathcal{M} \times \mathcal{P}} \pi^* \mathbb{E}_{\nu}^{\pi} w(\nu, \pi) \log \frac{w(\nu, \pi|h_{<mN+i})}{w(\nu, \pi|h_{<i})} \\
& \stackrel{(n)}{\leq} \frac{\eta}{m^3(m+1)} \sum_{\nu, \pi \in \mathcal{M} \times \mathcal{P}} \pi^* \mathbb{E}_{\nu}^{\pi} w(\nu, \pi) \log \frac{1}{w(\nu, \pi|h_{<i})} \\
& \stackrel{(o)}{=} \frac{\eta}{m^3(m+1)} \sum_{\nu \in \mathcal{M}} \pi^* \mathbb{E}_{\nu}^{\pi} w(\nu, \pi) \log \frac{1}{w(\nu, \pi)} \frac{w(\nu, \pi)}{w(\nu, \pi|h_{<i})} \\
& \stackrel{(p)}{=} \frac{\eta}{m^3(m+1)} \sum_{\nu, \pi \in \mathcal{M} \times \mathcal{P}} w(\nu) \log \frac{1}{w(\nu, \pi)} + \frac{\eta}{m^2(m+1)} \sum_{\nu, \pi \in \mathcal{M} \times \mathcal{P}} \pi^* \mathbb{E}_{\nu}^{\pi} w(\nu, \pi) \log \frac{w(\nu, \pi)}{w(\nu, \pi|h_{<i})} \\
& \stackrel{(q)}{=} \frac{\eta}{m^3(m+1)} \text{Ent}(w) + \frac{\eta}{m^2(m+1)} \sum_{h_{<i} \in \mathcal{H}^i} \sum_{\nu, \pi \in \mathcal{M} \times \mathcal{P}} w(\nu, \pi)^{\pi^*} P_{\nu}^{\pi}(h_{<i}) \log \frac{w(\nu, \pi)}{w(\nu, \pi|h_{<i})} \\
& \stackrel{(r)}{=} \frac{\eta}{m^3(m+1)} \text{Ent}(w) + \frac{\eta}{m^2(m+1)} \sum_{h_{<i} \in \mathcal{H}^i} \sum_{\nu, \pi \in \mathcal{M} \times \mathcal{P}} w(\nu, \pi|h_{<i})^{\pi^*} P_{\xi}^{\pi}(h_{<i}) \log \frac{w(\nu, \pi)}{w(\nu, \pi|h_{<i})} \\
& \stackrel{(s)}{=} \frac{\eta}{m^3(m+1)} \text{Ent}(w) - \frac{\eta}{m^2(m+1)} \pi^* \mathbb{E}_{\xi}^{\pi} [\text{IG}(h_{<i}|\varepsilon)] \stackrel{(t)}{\leq} \frac{\eta}{m^3(m+1)} \text{Ent}(w) \tag{31}
\end{aligned}$$

(h) expands the definition of the information gain. (i) rearranges the expectations and the sums, and expands $w(\nu, \pi|h_{<t+m})$ according to Bayes' rule. (j) converts the expectation to a expectation with respect to a different probability measure through simple cancellation. (k) implements a change of variable from t to $mk + i$. (l) moves a sum inside the logarithm. (m) cancels out all terms except the numerator of the last term and the denominator of the first. (n) follows from all posterior weights being ≤ 1 . (o) and (p) are obvious. (q) applies the definition of the entropy of a distribution $\text{Ent}(\cdot)$, and expands the expectation. (r)

changes the variable in the expectation; this is the reverse of (i) and (j). (s) applies the definition of the information gain (after inverting the fraction in the logarithm). (t) follows from the non-negativity of the information gain.

From Cohen et al. [6, Inequality 24], with minor modifications to our present case:

$$\begin{aligned}
\varepsilon &\leq (\mathbb{P}_\mu^{\pi^h}(h_{t:t+m-1}|h_{<t}) - \mathbb{P}_\xi^{\pi^h}(h_{t:t+m-1}|h_{<t}))^2 \\
&\stackrel{(a)}{\leq} \text{KL}_{h_{<t}, m}(\mathbb{P}_\mu^{\pi^h} \parallel \mathbb{P}_\xi^{\pi^h}) \\
&\stackrel{(b)}{\leq} \sum_{\nu, \pi \in \mathcal{M} \times \mathcal{P}} \frac{w(\nu, \pi|h_{<t})}{w(\mu, \pi^h|h_{<t})} \text{KL}_{h_{<t}, m}(\mathbb{P}_\nu^\pi \parallel \mathbb{P}_\xi^\pi) \\
&\stackrel{(c)}{\leq} \frac{1}{w(\mu, \pi^h|h_{<t})} \mathbb{E}_\xi^{\pi^h} \text{KL}(w(\cdot|h_{<t+m}) \parallel w(\cdot|h_{<t})) \\
&\stackrel{(d)}{\leq} \frac{1}{\inf_k w(\mu, \pi^h|h_{<k})} V_{m,0}^{\text{IG}}(h_{<t})
\end{aligned} \tag{32}$$

(a) is a result from information theory known as the entropy inequality, proven for example in [45]. (b) follows from the non-negativity of the KL-divergence, and the l.h.s. being one of the summands of the r.h.s. (c) follows from Cohen et al. [6] Lemma 4. And (d) follows from the definitions of the information gain value and the infimum.

Following Cohen et al. [6, Proof of Theorem 3], we show that if $\mathbb{P}_\mu^{\pi^h}(h_{t:t+m-1}|h_{<t}) - \mathbb{P}_\xi^{\pi^h}(h_{t:t+m-1}|h_{<t}) \rightarrow 0$ for all m , then $V_\mu^\pi(h_{<t}) - V_\xi^\pi(h_{<t}) \rightarrow 0$.

Let $\varepsilon > 0$. Since the agent has a bounded horizon, there exists an m such that for all t , $\frac{\Gamma_{t+m}}{\Gamma_t} \leq \varepsilon$. Recall

$$V_\nu^\pi(h_{<t}) = \frac{1}{\Gamma_t} \mathbb{E}_\nu^\pi \left[\sum_{k=t}^{\infty} \gamma^k r_k \mid h_{<t} \right] \tag{33}$$

Using the m from above, let

$$V_\nu^{\pi \setminus m}(h_{<t}) := \frac{1}{\Gamma_t} \mathbb{E}_\nu^\pi \left[\sum_{k=t}^{t+m-1} \gamma^k r_k \mid h_{<t} \right] \tag{34}$$

Since $r_t \in [0, 1]$,

$$|V_\nu^\pi(h_{<t}) - V_\nu^{\pi \setminus m}(h_{<t})| \leq \frac{\Gamma_{t+m}}{\Gamma_t} \leq \varepsilon \tag{35}$$

Suppose $V_\mu^\pi(h_{<t}) - V_\xi^\pi(h_{<t}) > 3\varepsilon$. Then $V_\mu^{\pi \setminus m}(h_{<t}) - V_\xi^{\pi \setminus m}(h_{<t}) > \varepsilon$. But since $\mathbb{P}_\mu^{\pi^h}(h_{t:t+m-1}|h_{<t}) - \mathbb{P}_\xi^{\pi^h}(h_{t:t+m-1}|h_{<t}) \rightarrow 0$, and the value is the expectation with respect to those measures, and reward is bounded, this can only occur finitely often. Thus, $V_\mu^\pi(h_{<t}) - V_\xi^\pi(h_{<t}) > 3\varepsilon$ holds only finitely often, so the values converge.

APPENDIX D ρ UCT ALGORITHM

To approximate expectimax planning, specifically to approximate the expected future rewards and therefore the value function, like [2, 42, 6] we used the ρ UCT Monte-Carlo tree search method. Below we have provided the algorithm for ρ UCT from [2]. The ρ UCT algorithm starts with an empty search tree Ψ over actions, and observation-reward pairs, then uses the provided model ρ to sample down and build the search tree, and then use those samples to compute a better approximation of the value function. If allowed to sample forever the approximation of the value function will converge to the true value function [44]. The difference between ρ UCT and regular Monte-Carlo methods is the optimistic choice of actions when expanding the search tree in line 32 of Algorithm 2. This choice of action incorporates an exploration component, with the inclusion of $C \sqrt{\frac{\log(T(h))}{T(h_a)}}$, as T is the number of times that history (or history and action) have been visited during the sampling process. This ensures that the whole tree is expanded in the limit.

Algorithm 2 ρ UCT [2, 44]**Require:** History h ; Search horizon m ; Samples budget κ ; Model ρ

```

1: INITIALIZE ( $\Psi$ ) // Search tree
2:  $n_{\text{samples}} \leftarrow 0$ 
3: repeat
4:    $\rho' \leftarrow \rho.\text{COPY}()$ 
5:   SAMPLE ( $\Psi, h, m$ )
6:    $n_{\text{samples}} \leftarrow n_{\text{samples}} + 1$ 
7:    $\rho \leftarrow \rho'$ 
8: until  $n_{\text{samples}} = \kappa$ 
9: return  $\arg \max_{a \in \mathcal{A}} \hat{V}_{\Psi}(a)$ 
10: function SAMPLE( $\Psi, h, m$ )
11:   if  $m = 0$  then
12:     return 0
13:   else if  $\Psi(h)$  is a chance node then
14:      $\rho.\text{PERFORM}(a)$ 
15:      $e = (o, r) \leftarrow \rho.\text{GENERATEPERCEPT}()$ 
16:      $\rho.\text{UPDATE}(a, e)$ 
17:     if  $T(he) = 0$  then
18:       Create chance node  $\Psi(he)$ 
19:       reward  $\leftarrow e.\text{REWARD} + \text{SAMPLE}(\Psi, he, m - 1)$ 
20:     else if  $T(h) = 0$  then
21:       reward  $\leftarrow \text{ROLLOUT}(h, m)$ 
22:     else
23:        $a \leftarrow \text{SELECTACTION}(\Psi, h)$ 
24:        $\hat{V}(h) \leftarrow \frac{1}{T(h)+1} (\text{reward} + T(h)\hat{V}(h))$ 
25:        $T(h) \leftarrow T(h) + 1$ 
26: function SELECTACTION( $\Psi, h$ )
27:    $\mathcal{U} = \{a \in \mathcal{A} : T(ha) = 0\}$ 
28:   if  $\mathcal{U} \neq \emptyset$  then
29:     Pick  $a \in \mathcal{U}$  uniformly at random
30:     Create node  $\Psi(ha)$ 
31:     return  $a$ 
32:   else
33:     return  $\arg \max_{a \in \mathcal{A}} \left\{ \frac{1}{m(\beta-\alpha)} \hat{V}(ha) + C \sqrt{\frac{\log(T(h))}{T(ha)}} \right\}$ 
34: function ROLLOUT( $h, m$ )
35:   reward  $\leftarrow 0$ 
36:   for  $i = 1$  to  $m$  do
37:      $a \sim \pi_{\text{rollout}}(h)$ 
38:      $e = (o, r) \sim \rho(e|ha)$ 
39:     reward  $\leftarrow \text{reward} + r$ 
40:      $h \leftarrow hae$ 
41:   return reward

```
