

A Strongly Asymptotically Optimal Agent in General Environments

Michael K. Cohen^{1*}, Elliot Catt¹ and Marcus Hutter¹

¹Australian National University

{michael.cohen, elliot.carpentercatt, marcus.hutter}@anu.edu.au

Abstract

Reinforcement Learning agents are expected to eventually perform well. Typically, this takes the form of a guarantee about the asymptotic behavior of an algorithm given some assumptions about the environment. We present an algorithm for a policy whose value approaches the optimal value with probability 1 in all computable probabilistic environments, provided the agent has a bounded horizon. This is known as strong asymptotic optimality, and it was previously unknown whether it was possible for a policy to be strongly asymptotically optimal in the class of all computable probabilistic environments. Our agent, Inquisitive Reinforcement Learner (Inq), is more likely to explore the more it expects an exploratory action to reduce its uncertainty about which environment it is in, hence the term inquisitive. Exploring inquisitively is a strategy that can be applied generally; for more manageable environment classes, inquisitiveness is tractable. We conducted experiments in “grid-worlds” to compare the Inquisitive Reinforcement Learner to other weakly asymptotically optimal agents.

1 Introduction

“Efforts to solve [an instance of the exploration-exploitation problem] so sapped the energies and minds of Allied analysts that the suggestion was made that the problem be dropped over Germany, as the ultimate instrument of intellectual sabotage.”
—Peter Whittle [Whittle, 1979]

The Allied analysts were considering the simplest possible problem in which there is a trade-off to be made between exploiting, taking the apparently best option, and exploring, choosing a different option to learn more. We tackle what we consider the most difficult instance of the exploration-exploitation trade-off problem: when the environment could be any computable probability distribution, not just a multi-armed bandit, how can one achieve optimal performance in the limit?

*Contact Author

Our work is within the Reinforcement Learning (RL) paradigm: an agent selects an action, and the environment responds with an observation and a reward. The interaction may end, or it may continue forever. Each interaction cycle is called a timestep. The agent has a discount function that weights its relative concern for the reward it achieves at various future timesteps. The agent’s job is to select actions that maximize the total expected discounted reward it achieves in its lifetime. The “value” of an agent’s policy at a certain point in time is the expected total discounted reward it achieves after that time if it follows that policy. One formal specification of the exploration-exploitation problem is: what policy can an agent follow so that the policy’s value approaches the value of the optimal informed policy with probability 1, even when the agent doesn’t start out knowing the true dynamics of its environment?

Most work in RL makes strong assumptions about the environment—that the environment is Markov, for instance. Impressive recent development in the field of reinforcement learning often makes use of the Markov assumption, including Deep Q Networks [Mnih *et al.*, 2015], A3C [Mnih *et al.*, 2016], Rainbow [Hessel *et al.*, 2018], and AlphaZero [Silver *et al.*, 2017]. Another example of making strong assumptions in RL comes from some model-based algorithms that implicitly assume that the environment is representable by, for example, a fixed-size neural network, or whatever construct is used to model the environment. We do not make any such assumptions.

Many recent developments in RL are largely about tractably learning to exploit; how to explore intelligently is a separate problem. We address the latter problem. Our approach, inquisitiveness, is based on Orseau *et al.*’s [2013] Knowledge Seeking Agent for Stochastic Environments, which selects the actions that best inform the agent about what environment it is in. Our Inquisitive Reinforcement Learner (Inq) explores like a knowledge seeking agent, and is more likely to explore when there is apparently (according to its current beliefs) more to be learned. Sometimes exploring well requires “expeditions,” or many consecutive exploratory actions. Inq entertains expeditions of all lengths, although it follows the longer ones less often, and doesn’t resolutely commit in advance to seeing the expedition through.

This is a very human approach to information acquisition.

When we spot an opportunity to learn something about our natural environment, we feel inquisitive. We get distracted. We are inclined to check it out, even if we don't see directly in advance how this information might help us better achieve our goals. Moreover, if we can tell that the opportunity to learn something requires a longer term project, we may find ourselves less inquisitive.

For the class of computable environments (stochastic environments that follow a computable probability distribution), it was previously unknown whether any policy could achieve strong asymptotic optimality (convergence of the value to optimality with probability 1). Lattimore et al. [2011] showed that no deterministic policy could achieve this. The key advantage that stochastic policies have is that they can let the exploration probability go to 0 while still exploring infinitely often. (For example, an agent that explores with probability $1/t$ at time t still explores infinitely often).

There is a weaker notion of optimality—"weak asymptotic optimality"—for which positive results already exist; this condition requires that the average value over the agent's lifetime approach optimality. Lattimore et al. [2011] identified a weakly asymptotically optimal agent for *deterministic* computable environments; the agent maintains a list of environments consistent with its observations, exploiting as if it is in the first such one, and exploring in bursts. A recent algorithm for a Thompson Sampling Bayesian agent was shown, with an elegant proof, to be weakly asymptotically optimal in all computable environments, but not strongly asymptotically optimal [Leike et al., 2016].

Most work in RL regards (Partially Observable) Markov Decision Processes (PO)MDPs. However, environments that enter completely novel states infinitely often render (PO)MDP algorithms helpless. For example, an RL agent acting as a chatbot, optimizing a function, or proving mathematical theorems would struggle to model the environment as an MDP, and would likely require an exploration mechanism like ours. In the chatbot case, for instance, as a conversation with a person progresses, the person never returns to the same state.

If we formally compare Inq to existing algorithms in MDPs, we find that many achieve asymptotic optimality. Epsilon-greedy, upper confidence bound, and Thompson sampling exploration strategies suffice in MDPs. Our primary motivation is for the sorts of environments described above. To discriminate between exploratory approaches in *ergodic* MDPs, one can formally bound regret, and we would like to do this for Inq in the future.

For comparison, some algorithms which use the MDP formalism also consider information-theoretic approaches to exploration, such as VIME [Houthoof et al., 2016], the agent in [Still, 2009], and TEXPLORE-VANIR [Hester and Stone, 2012].

In Section 2, we formally describe the RL setup and present notation. In Section 3, we present the algorithm for Inq. In Section 4, we prove our main result: that Inq is strongly asymptotically optimal. In Section 5, we present experimental results comparing Inq to weakly asymptotically optimal agents. Finally, we discuss the relevance of this exploration regime to tractable algorithms. Appendix A collates notation

and definitions for quick reference. Appendix B contains the proofs of the lemmas.

2 Notation

We follow the notation of Orseau, et al. [2013]. The reinforcement learning setup is as follows: \mathcal{A} is a finite set of actions available to the agent; \mathcal{O} is a finite set of observations it might observe, and $\mathcal{R} = [0, 1] \cap \mathbb{Q}$ is the set of possible rewards. The set of all possible interactions in a timestep is $\mathcal{H} := \mathcal{A} \times \mathcal{O} \times \mathcal{R}$. At every timestep, one element from this set occurs. A reinforcement learner's policy π is a stochastic function which outputs an action given an interaction history, denoted by $\pi : \mathcal{H}^* \rightsquigarrow \mathcal{A}$. ($\mathcal{H}^* := \bigcup_{i=0}^{\infty} \mathcal{H}^i$ represents all finite strings from an alphabet \mathcal{X}). An environment is a stochastic function which outputs an observation and reward given an interaction history and an action: $\nu : \mathcal{H}^* \times \mathcal{A} \rightsquigarrow \mathcal{O} \times \mathcal{R}$. For a stochastic function $f : \mathcal{X} \rightarrow \mathcal{Y}$, $f(y|x)$ denotes the probability that f outputs $y \in \mathcal{Y}$ when $x \in \mathcal{X}$ is input.

A policy and an environment induce a probability measure over \mathcal{H}^{∞} , the set of all possible infinite histories: for $h \in \mathcal{H}^*$, $P_{\nu}^{\pi}(h)$ denotes the probability that an infinite history begins with h when actions are sampled from the policy π , and observations and rewards are sampled from the environment ν . Formally, we define this inductively: $P_{\nu}^{\pi}(\epsilon) \mapsto 1$, where ϵ is the empty history, and for $h \in \mathcal{H}^*$, $a \in \mathcal{A}$, $o \in \mathcal{O}$, $r \in \mathcal{R}$, we define $P_{\nu}^{\pi}(haor) \mapsto P_{\nu}^{\pi}(h)\pi(a|h)\nu(or|ha)$. In an infinite history $h_{1:\infty} \in \mathcal{H}^{\infty}$, a_t , o_t , and r_t refer to the t th action, observation and reward, and h_t refers to the t th timestep: $a_t o_t r_t$. $h_{<t}$ refers to the first $t - 1$ timesteps, and $h_{t:k}$ refers to the string of timesteps t through k (inclusive). Strings of actions, observations, and rewards are notated similarly.

A Bayesian agent deems a class of environments a priori feasible. Its "beliefs" take the form of a probability distribution over which environment is the true one. We call this the agent's belief distribution. In our formulation, Inq considers any computable environment feasible, and starts with a prior belief distribution based on the environments' Kolmogorov complexities: that is, the length of the shortest program that computes the environment on some reference machine. However, all our results hold as long as the true environment is contained in the class of environments that are considered feasible, and as long as the prior belief distribution assigns nonzero probability to each environment in the class. We take \mathcal{M} to be the class of all computable environments, and $w(\nu) := 2^{-K(\nu)(1+\epsilon)}/\mathcal{N}$ to be the prior probability of the environment ν , where K is the Kolmogorov complexity, $\epsilon > 0$, and \mathcal{N} is a normalization constant. ($\epsilon > 0$ ensures the prior has finite entropy, which facilitates analysis.) A smaller class with a different prior probability could easily be substituted for \mathcal{M} and $w(\nu)$.

We use ξ to denote the agent's beliefs about future observations. Together with a policy π it defines a Bayesian mixture measure: $P_{\xi}^{\pi}(\cdot) := \sum_{\nu \in \mathcal{M}} w(\nu) P_{\nu}^{\pi}(\cdot)$. The posterior belief distribution of the agent after observing a history $h \in \mathcal{H}^*$ is $w(\nu|h) := w(\nu) P_{\nu}^{\pi'}(h) / P_{\xi}^{\pi'}(h)$. This definition is independent of the choice of π' as long as $P_{\xi}^{\pi'}(h) > 0$; we can fix a reference policy π' just for this definition if we

like. We sometimes also refer to the conditional distribution $\xi(or|ha) := \sum_{\nu \in \mathcal{M}} w(\nu|h)\nu(or|ha)$.

The agent’s discount at a timestep is denoted γ_t . To normalize the agent’s policy’s value to $[0, 1]$, we introduce $\Gamma_t := \sum_{k=t}^{\infty} \gamma_k$. (Normalization makes value convergence nontrivial). We consider an agent with a bounded horizon: $\forall \varepsilon > 0 \exists m \forall t : \Gamma_{t+m}/\Gamma_t \leq \varepsilon$. Intuitively, this means that the agent does not become more and more farsighted over time. Note this does not require a finite horizon. A classic discount function giving a bounded horizon is a geometric one: for $0 \leq \gamma < 1$, $\gamma_t = \gamma^t$. The value of a policy π in an environment ν , given a history $h_{<t} \in \mathcal{H}^{t-1}$, is

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

Here, the expectation is with respect to the probability measure P_{ν}^{π} . Reinforcement Learning is the attempt to find a policy that makes this value high, without access to ν .

3 Inquisitive Reinforcement Learner

We first describe how Inq exploits, then how it explores. It exploits by maximizing the discounted sum of its reward in expectation over its current beliefs, and it explores by following maximally informative “exploratory expeditions” of various lengths.

An optimal policy with respect to an environment ν is a policy that maximizes the value.

$$\pi_{\nu}^*(\cdot) := \operatorname{argmax}_{\pi \in \Pi} V_{\nu}^{\pi}(\cdot) \quad (2)$$

where $\Pi = \mathcal{H}^* \rightsquigarrow \mathcal{A}$ is the space of all policies. An optimal deterministic policy always exists [Lattimore and Hutter, 2014b]. When exploiting, Inq simply maximizes the value according to its belief distribution ξ . Since this policy is deterministic, we write $a^*(h_{<t})$ to mean the unique action at time t for which $\pi_{\xi}^*(a|h_{<t}) = 1$. That is the exploitative action.

The most interesting feature of Inq is how it gets distracted by the opportunity to explore. Inq explores to learn. An agent has learned from an observation if its belief distribution w changes significantly after making that observation. If the belief distribution has hardly changed, then the observation was not very informative. The typical information-theoretic measure for how well a distribution Q approximates a distribution P is the KL-divergence, $\text{KL}(P||Q)$. Thus, a principled way to quantify the information that an agent gains in a timestep is the KL-divergence from the belief distribution at time $t+1$ to the belief distribution at time t . This is the rationale behind the construction of Orseau, et al.’s [2013] Knowledge Seeking Agent, which maximizes this expected information gain.

Letting $h_{<t} \in \mathcal{H}^{t-1}$ and $h' \in \mathcal{H}^*$, the information gain at time t is defined:

$$\text{IG}(h'|h_{<t}) := \sum_{\nu \in \mathcal{M}} w(\nu|h_{<t}h') \log \frac{w(\nu|h_{<t}h')}{w(\nu|h_{<t})} \quad (3)$$

Recall that $w(\nu|h)$ is the posterior probability assigned to ν after observing h .

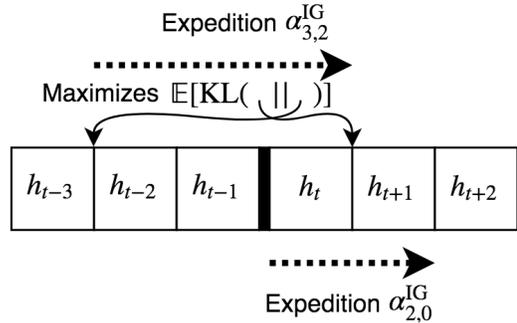


Figure 1: **Example Expeditions.** Expeditions maximize the expected KL-divergence from the posterior at the end to the posterior at the beginning.

An m -step expedition, denoted α^m , represents all contingencies for how an agent will act for the next m timesteps. It is a deterministic policy that takes history-fragments of length less than m and returns an action:

$$\alpha^m : \bigcup_{i=0}^{m-1} \mathcal{H}^i \rightarrow \mathcal{A} \quad (4)$$

$P_{\xi}^{\alpha^m}(h_{<t+k}|h_{<t})$ is a conditional distribution defined for $0 \leq k \leq m$, which represents the conditional probability of observing $h_{<t+k}$ if the expedition α^m is followed starting at time t , after observing $h_{<t}$. Now we can consider the information-gain value of an m -step expedition. It is the expected information gain upon following that expedition:

$$V^{\text{IG}}(\alpha^m, h_{<t}) := \sum_{h_{t:t+m-1} \in \mathcal{H}^m} P_{\xi}^{\alpha^m}(h_{<t+m}|h_{<t}) \text{IG}(h_{t:t+m-1}|h_{<t}) \quad (5)$$

At a time t , one might consider many expeditions: the one-step expedition which maximizes expected information gain, the two-step expedition doing the same, etc. Or one might consider carrying on with an expedition that began three timesteps ago.

Definition 1. At time t , the m - k expedition is the m -step expedition beginning at time $t-k$ which maximized the expected information gain from that point.¹

$$\alpha_{m,k}^{\text{IG}}(h_{<t}) := \operatorname{argmax}_{\alpha^m : \bigcup_{i=0}^{m-1} \mathcal{H}^i \rightarrow \mathcal{A}} V^{\text{IG}}(\alpha^m, h_{<t-k}) \quad (6)$$

Example expeditions are diagrammed in Figure 1.

Expeditions are functions which return an action given what has been seen so far on the expedition. The m - k exploratory action is the action to take at time t according to the m - k expedition:

$$a_{m,k}^{\text{IG}}(h_{<t}) := \alpha_{m,k}^{\text{IG}}(h_{<t})(h_{t-k:t-1}) \quad (7)$$

Naturally, this is only defined for $k < m, t$, since the expedition function can’t accept a history fragment of length $\geq m$, and $t-k$ must be positive. Note also that if $k = 0$, $h_{t-k:t-1}$ evaluates to the empty string, ϵ .

¹Ties in the argmax are broken arbitrarily.

The reason Inq doesn't ignore expeditions that started in the past is that Inq must have some chance of actually executing the whole expedition (for every expedition). If the probability of completing an expedition is 0, one cannot use it for a bound on Inq's belief-accuracy.

Definition 2. Let $\rho(h_{<t}, m, k)$ be the probability of taking the m - k exploratory action after observing a history $h_{<t}$.

$$\rho(h_{<t}, m, k) := \min \left\{ \frac{1}{m^2(m+1)}, \eta V^{\text{IG}}(\alpha_{m,k}^{\text{IG}}(h_{<t}), h_{<t-k}) \right\} \quad (8)$$

where η is an exploration constant.

Note in the definition of $\rho(h_{<t}, m, k)$ that the probability of following an expedition goes to 0 if the expected information gain from that expedition goes to 0. The first term in the min ensures the probabilities will not sum to more than 1. The total probability of exploration is defined:

$$\beta(h_{<t}) := \sum_{m \in \mathbb{N}} \sum_{k < m, t} \rho(h_{<t}, m, k) \leq \sum_{m \in \mathbb{N}} \sum_{k < m} \frac{1}{m^2(m+1)} = 1 \quad (9)$$

The feature that makes Inq inquisitive is that $\rho(h_{<t}, m, k)$ is proportional to the expected information gain from the m - k expedition, $V^{\text{IG}}(\alpha_{m,k}^{\text{IG}}(h_{<t}), h_{<t-k})$. Note that completing an m -step expedition requires randomly deciding to explore in that way on m separate occasions. While this may seem inefficient, if the agent always got boxed into long expeditions, the value of its policy would plummet infinitely often.

Finally, Inq's policy π^\dagger , defined in Algorithm 1, takes the m - k exploratory action with probability $\rho(\cdot, m, k)$, and takes the exploitative action otherwise.²

Algorithm 1 Inquisitive Reinforcement Learner's Policy π^\dagger

- 1: **while** True **do**
 - 2: calculate $\rho(h_{<t}, m, k)$ for all m and for all $k < \min\{m, t\}$
 - 3: take action $a_{m,k}^{\text{IG}}(h_{<t})$ with probability $\rho(h_{<t}, m, k)$
 - 4: take action $a^*(h_{<t})$ with probability $1 - \beta(h_{<t})$
-

4 Strong Asymptotic Optimality

Here we present our central result: that the value of π^\dagger approaches the optimal value. We present the theorem, motivate the result, and proceed to the proof. We recommend the reader have Appendix A at hand for quickly looking up definitions and notation.

Before presenting the theorem, we clarify an assumption, and define the optimal value. We call the true environment

²This algorithm is written in a simplified way that does not halt, but if a real number in $[0, 1]$ is sampled first, the actions can be assigned to disjoint intervals successively until the sampled real number lands in one of them.

μ , and we assume that $\mu \in \mathcal{M}$. For \mathcal{M} the class of computable environments, this is a very unassuming assumption. The optimal value is simply the value of the optimal policy with respect to the true environment:

$$V_\mu^*(h_{<t}) := \sup_{\pi \in \Pi} V_\mu^\pi(h_{<t}) = V_\mu^{\pi^*}(h_{<t}) \quad (10)$$

Recall also that we have assumed the agent has a bounded horizon in the sense that $\forall \varepsilon \exists m \forall t : \Gamma_{t+m}/\Gamma_t \leq \varepsilon$. The Strong Asymptotic Optimality theorem is that under these conditions, the value of Inq's policy approaches the optimal value with probability 1, when actions are sampled from Inq's policy and observations and rewards are sampled from the true environment μ .

Theorem 3 (Strong Asymptotic Optimality). As $t \rightarrow \infty$,

$$V_\mu^*(h_{<t}) - V_\mu^{\pi^\dagger}(h_{<t}) \rightarrow 0 \text{ with } P_\mu^{\pi^\dagger}\text{-prob. } 1$$

where $\mu \in \mathcal{M}$ is the true environment.

For a Bayesian agent, uncertainty about on-policy observations goes to 0. Since "on-policy" for Inq includes, with some probability, all maximally informative expeditions, Inq eventually has little uncertainty about the result of any course of action, and can therefore successfully select the optimal course. For any fixed horizon, Inq's mixture measure ξ approaches the true environment μ .

We use the following notation for a particular KL-divergence that plays a central role in the proof:

$$\text{KL}_{h_{<t}, n}(P_{\nu_1}^\pi \parallel P_{\nu_2}^\pi) := \sum_{h' \in \mathcal{H}^n} P_{\nu_1}^\pi(h'|h_{<t}) \log \frac{P_{\nu_1}^\pi(h'|h_{<t})}{P_{\nu_2}^\pi(h'|h_{<t})} \quad (11)$$

This quantifies the difference between the expected observations of two different environments that would arise in the next n timesteps when following policy π . $\text{KL}_{h_{<t}, \infty}$ denotes the limit of the above as $n \rightarrow \infty$, which exists by [Orseau *et al.*, 2013, proof of Theorem 3].

In dealing with the KL-divergence, we simplify matters by asserting that $0 \log 0 := 0$, and $0 \log \frac{0}{0} := 0$.

We begin with a lemma that equates the information gain value of an expedition with the expected prediction error. The KL-divergence on the right hand side represents how different ν and ξ appear when following the expedition in question.

Lemma 4.

$$V^{\text{IG}}(\alpha^m, h_{<t}) = \sum_{\nu \in \mathcal{M}} w(\nu|h_{<t}) \text{KL}_{h_{<t}, m} \left(P_\nu^{\alpha^m} \parallel P_\xi^{\alpha^m} \right)$$

Proofs of Lemmas appear in Appendix B.

Recall that $w(\nu|h_{<t})$ is the posterior weight that Inq assigns to the environment ν after observing $h_{<t}$. We show that the infimum of this value is strictly positive with probability 1.

Lemma 5. $\inf_t w(\mu|h_{<t}) > 0$ w. $P_\mu^{\pi^\dagger}$ -p. 1

Next, we show that every exploration probability $\rho(h_{<t}, m, k)$ goes to 0. From here, all "w.p.1" statements mean with $P_\mu^{\pi^\dagger}$ -probability 1, if not otherwise specified.

Lemma 6.

$$\rho(h_{<t}, m, k) \xrightarrow{t \rightarrow \infty} 0 \text{ w.p.1}$$

The essence of the proof is that with a finite-entropy prior, there is only a finite amount of information to gain, so the expected information gain (and the exploration probability) goes to 0.

Next, we show that the total exploration probability goes to 0:

Lemma 7.

$$\beta(h_{<t}) \rightarrow 0 \text{ w.p.1}$$

Lemma 8 shows that the probabilities assigned by ξ converge to those of μ .

Lemma 8. $\forall m \in \mathbb{N}, h_{t:t+m-1} \in \mathcal{H}^m, \alpha^m : \bigcup_{i=0}^{m-1} \mathcal{H}^i \rightarrow \mathcal{A}$:

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

The proof of Lemma 8 roughly follows the following argument: if all exploration probabilities go to 0, then the informativeness of the maximally informative expeditions goes to 0, so the informativeness of all expeditions goes to 0, meaning the prediction error goes to 0.

Finally, we prove the Strong Asymptotic Optimality Theorem: $V_\mu^*(h_{<t}) - V_\mu^{\pi^\dagger}(h_{<t}) \rightarrow 0$ with $P_\mu^{\pi^\dagger}$ -prob. 1.

Proof of Theorem 3. 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

$$V_\mu^*(h_{<t}) = \frac{1}{\Gamma_t} \mathbb{E}_\mu^{\pi_\mu^*} \left[\sum_{k=t}^{\infty} \gamma_k r_k \mid h_{<t} \right] \quad (12)$$

Using the m from above, let

$$V_\mu^{*\setminus m}(h_{<t}) := \frac{1}{\Gamma_t} \mathbb{E}_\mu^{\pi_\mu^*} \left[\sum_{k=t}^{t+m-1} \gamma_k r_k \mid h_{<t} \right] \quad (13)$$

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

$$|V_\mu^*(h_{<t}) - V_\mu^{*\setminus m}(h_{<t})| \leq \frac{\Gamma_{t+m}}{\Gamma_t} \leq \varepsilon \quad (14)$$

We continue from there:

$$\begin{aligned} & V_\mu^*(h_{<t}) \\ & \leq V_\mu^{*\setminus m}(h_{<t}) + \varepsilon \\ & = \frac{1}{\Gamma_t} \sum_{h_{t:t+m-1} \in \mathcal{H}^m} P_\mu^{\pi_\mu^*}(h_{t:t+m-1}|h_{<t}) \sum_{k=t}^{t+m-1} \gamma_k r_k + \varepsilon \\ \stackrel{(a)}{\leq} & \frac{1}{\Gamma_t} \sum_{h_{t:t+m-1} \in \mathcal{H}^m} P_\xi^{\pi_\xi^*}(h_{t:t+m-1}|h_{<t}) \sum_{k=t}^{t+m-1} \gamma_k r_k + 2\varepsilon \\ \stackrel{(b)}{\leq} & \frac{1}{\Gamma_t} \mathbb{E}_\xi^{\pi_\xi^*} \left[\sum_{k=t}^{\infty} \gamma_k r_k \mid h_{<t} \right] + 2\varepsilon \end{aligned}$$

$$\begin{aligned} & \stackrel{(c)}{\leq} \frac{1}{\Gamma_t} \mathbb{E}_\xi^{\pi_\xi^*} \left[\sum_{k=t}^{\infty} \gamma_k r_k \mid h_{<t} \right] + 2\varepsilon \\ & \stackrel{(d)}{\leq} \frac{1}{\Gamma_t} \sum_{h_{t:t+m-1} \in \mathcal{H}^m} P_\xi^{\pi_\xi^*}(h_{t:t+m-1}|h_{<t}) \sum_{k=t}^{t+m-1} \gamma_k r_k + 3\varepsilon \\ \stackrel{(e)}{\leq} & \frac{1}{\Gamma_t} \sum_{h_{t:t+m-1} \in \mathcal{H}^m} P_\mu^{\pi_\mu^*}(h_{t:t+m-1}|h_{<t}) \sum_{k=t}^{t+m-1} \gamma_k r_k + 4\varepsilon \\ \stackrel{(f)}{\leq} & \frac{1}{\Gamma_t} \sum_{h_{t:t+m-1} \in \mathcal{H}^m} \frac{P_\mu^{\pi_\mu^\dagger}(h_{t:t+m-1}|h_{<t})}{\prod_{k=t}^{t+m-1} (1 - \beta(h_{<k}))} \sum_{k=t}^{t+m-1} \gamma_k r_k \\ & + 4\varepsilon \\ & \leq \frac{1}{\Gamma_t} \sum_{h_{t:t+m-1} \in \mathcal{H}^m} \frac{P_\mu^{\pi_\mu^\dagger}(h_{t:t+m-1}|h_{<t})}{(1 - \max_{t \leq k < t+m} \beta(h_{<k}))^m} \sum_{k=t}^{t+m-1} \gamma_k r_k \\ & + 4\varepsilon \\ \stackrel{(g)}{\leq} & \frac{1}{\Gamma_t} \sum_{h_{t:t+m-1} \in \mathcal{H}^m} \frac{P_\mu^{\pi_\mu^\dagger}(h_{t:t+m-1}|h_{<t})}{(1 - \varepsilon')^m} \sum_{k=t}^{t+m-1} \gamma_k r_k \\ & + 4\varepsilon \\ & \stackrel{(h)}{\leq} \frac{1}{(1 - \varepsilon')^m \Gamma_t} \mathbb{E}_\mu^{\pi_\mu^\dagger} \left[\sum_{k=t}^{\infty} \gamma_k r_k \mid h_{<t} \right] + 4\varepsilon \\ & = \frac{1}{(1 - \varepsilon')^m} V_\mu^{\pi_\mu^\dagger}(h_{<t}) + 4\varepsilon \\ & = V_\mu^{\pi^\dagger}(h_{<t}) + 4\varepsilon + \left(\frac{1}{(1 - \varepsilon')^m} - 1 \right) V_\mu^{\pi^\dagger}(h_{<t}) \\ \stackrel{(i)}{\leq} & V_\mu^{\pi^\dagger}(h_{<t}) + 4\varepsilon + \left(\frac{1}{(1 - \varepsilon')^m} - 1 \right) \quad (15) \end{aligned}$$

(a), (e), (f), and (g) all hold with probability 1. (a) follows from Lemma 8: for all $m, P_\xi^{\pi_\xi^*}(\cdot|h_{<t}) \rightarrow P_\mu^{\pi_\mu^*}(\cdot|h_{<t})$ for all conditional probabilities of histories of length m , with probability 1, and the countable sum is bounded (by Γ_t). (b) follows from adding more non-negative terms to the sum. (c) follows π_ξ^* being the ξ -optimal policy, and therefore it accrues at least as much expected reward in environment ξ as π_μ^* does. (d) follows from $\sum_{k=t+m}^{\infty} \gamma_k / \Gamma_t = \Gamma_{t+m} / \Gamma_t \leq \varepsilon$, and $r_t \in [0, 1]$. (e) follows from Lemma 8 just as (a) did. (f) follows because the product in the denominator is the probability that π^\dagger mimics π_ξ^* for m consecutive timesteps, and by Lemma 7 there is a time after which this probability is uniformly strictly positive. (g) follows from Lemma 7: $\beta(h_{<k}) \rightarrow 0$ with probability 1. (h) follows from adding more non-negative terms to the sum. Finally, (i) follows from the value being normalized to $[0, 1]$ by Γ_t .

$\forall \delta > 0 \exists \varepsilon > 0, \varepsilon' > 0 : 4\varepsilon + \left(\frac{1}{(1 - \varepsilon')^m} - 1 \right) < \delta$. Letting $T = \max\{T_1, T_2, T_3, T_4\}$, we can combine the equations above to give

$$\forall \delta > 0 \exists T \forall t > T : V_\mu^*(h_{<t}) - V_\mu^{\pi^\dagger}(h_{<t}) < \delta \text{ w.p.1} \quad (16)$$

Since $V_\mu^*(h_{<t}) \geq V_\mu^{\pi^\dagger}(h_{<t})$,

$$V_\mu^*(h_{<t}) - V_\mu^{\pi^\dagger}(h_{<t}) \rightarrow 0 \text{ w.p.1} \quad (17)$$

□

Strong Asymptotic Optimality is not a guarantee of efficacy; consider an agent that “commits suicide” on the first timestep, and thereafter receives a reward of 0 no matter what it does. This agent is asymptotically optimal, but not very useful. In general, when considering many environments with many different “traps,” bounded regret is impossible to guarantee [Hutter, 2005], but one can still demand from a reinforcement learner that it make the best of whatever situation it finds itself in by correctly identifying (in the limit) the optimal policy.

We suspect that strong asymptotic optimality would not hold if Inq had an unbounded horizon, since its horizon of concern may grow faster than it can learn about progressively more long-term dynamics of the environment. Going more into the technical details, let Δ_{kt} be, roughly “at time t , how much does ξ differ from μ regarding predictions about the next k timesteps?” A lemma in our proof is that $\forall k \lim_{t \rightarrow \infty} \Delta_{kt} = 0$, but this does not imply, for example, that $\lim_{z \rightarrow \infty} \Delta_{zz} = 0$. If the horizon which is necessary to predict is growing over time, Inq might not be strongly asymptotically optimal.

Indeed, we tenuously suspect that it is impossible for an agent with an unbounded time horizon to be strongly asymptotically optimal in the class of all computable environments. If that is true, then the assumptions that our result relies on (namely that the true environment is computable, and the agent has a bounded horizon) are the bare minimum for strong asymptotic optimality to be possible.

Inq is not computable; in fact, no computable policy can be strongly asymptotically optimal in the class of all computable environments (Lattimore, et al. [2011] show this for deterministic policies, but a simple modification extends this to stochastic policies). For many smaller environment classes, however, Inq would be computable, for example if \mathcal{M} is finite, and perhaps for decidable \mathcal{M} in general. The central result, that inquisitiveness is an effective exploration strategy, applies to any Bayesian agent.

5 Experimental Results

We compared Inq with other known weakly asymptotically optimal agents, Thompson sampling and BayesExp [Lattimore and Hutter, 2014a], in the grid-world environment using AIXIjs [Aslanides, 2017] which has previously been used to compare asymptotically optimal agents [Aslanides *et al.*, 2017]. We tested in 10×10 grid-worlds, and 20×20 grid-worlds, both with a single dispenser with probability of dispensing reward 0.75; that is, if the agent enters that cell, the probability of a reward of 1 is 0.75. Following the conventions of [Aslanides *et al.*, 2017] we averaged over 50 simulations, used discount factor $\gamma = 0.99$, 600 MCTS samples, and planning horizon of 6. The planning horizon restricts m , and the number of MCTS samples is an input to ρ UCT [Silver and Veness, 2010], which we use instead of expectimax. The algorithm for the approximate version of Inq is in Appendix C.

The code used for this experiment is available online at <https://github.com/ejcatt/aixijs>, and this version of Inq can be run in the browser at <https://ejcatt.github.io/aixijs/demo.html#inq>. We found that using small values for η , specifically $\eta \leq 1$ worked well. For our experiments we chose $\eta = 1$.

In the 10×10 grid-worlds Inq performed comparably to both BayesExp and Thompson sampling. However in the 20×20 grid-worlds Inq performed comparably to BayesExp, and outperformed Thompson sampling. This is likely because when the Thompson Sampling Agent samples an environment with a reward dispenser that is inaccessible within its planning horizon, the agent acts randomly rather than seeking new cells. This is contrast to Inq and BayesExp which always have an incentive to explore the frontier of cells that have not been visited. This is especially relevant in the larger grid where the Thompson sampling agent is more likely to act as if the dispenser is deep in uncharted territory, rather than nearby. In a grid-world, good exploration is just about visiting new states, which both Inq and BayesExp successfully seek.

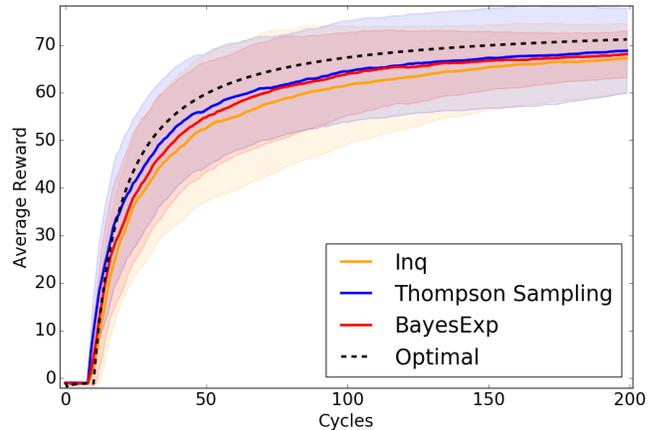


Figure 2: 10×10 Grid-worlds

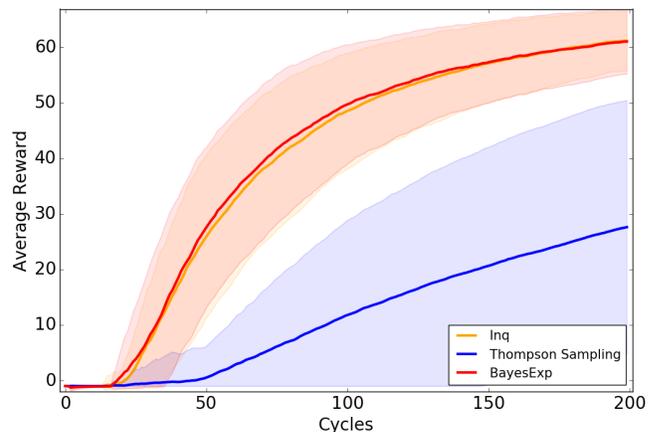


Figure 3: 20×20 Grid-worlds

6 Conclusion

We have shown that it is possible for an agent with a bounded horizon to be strongly asymptotically optimal in the class of all computable environments. No existing RL agent has as strong an optimality guarantee as Inq. The nature of the exploration regime that accomplishes this is perhaps of wider interest. We formalize an agent that gets distracted from reward maximization by its inquisitiveness: the more it expects to learn from an expedition, the more inclined it is to take it.

We have confirmed experimentally that inquisitiveness is a practical and effective exploration strategy for Bayesian agents with manageable model classes.

There are two main avenues for future work we would like to see. The first regards possible extensions of inquisitiveness: we have defined inquisitiveness for Bayesian agents with countable model-classes, but inquisitiveness could also be defined for a Bayesian agent with a continuous model class, such as a Q-learner using a Bayesian Neural Network. The second avenue regards the theory of strong asymptotic optimality itself: is Inq strongly asymptotically optimal for more farsighted discounters? If not, can it be modified to accomplish that? Or is it indeed impossible for an agent with an unbounded horizon to be strongly asymptotically optimal in the class of computable environments? Answers to these questions, besides being interesting in their own right, will likely inform the design of tractable exploration strategies, in the same way that this work has done.

Acknowledgements

This work was supported by the Open Philanthropy Project AI Scholarship and the Australian Research Council Discovery Projects DP150104590.

References

- [Aslanides *et al.*, 2017] John Aslanides, Jan Leike, and Marcus Hutter. Universal reinforcement learning algorithms: survey and experiments. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 1403–1410. AAAI Press, 2017.
- [Aslanides, 2017] John Aslanides. AIXIjs: A software demo for general reinforcement learning. *arXiv preprint arXiv:1705.07615*, 2017.
- [Hessel *et al.*, 2018] Matteo Hessel, Joseph Modayil, Hado van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning. In *Proc. of AAAI Conference on Artificial Intelligence*, 2018.
- [Hester and Stone, 2012] Todd Hester and Peter Stone. Intrinsically motivated model learning for a developing curious agent. In *2012 IEEE international conference on development and learning and epigenetic robotics (ICDL)*, pages 1–6. IEEE, 2012.
- [Houthoofd *et al.*, 2016] Rein Houthoofd, Xi Chen, Yan Duan, John Schulman, Filip De Turck, and Pieter Abbeel. Curiosity-driven exploration in deep reinforcement learning via bayesian neural networks. *arXiv preprint arXiv:1605.09674*, 2016.
- [Hutter, 2005] Marcus Hutter. *Universal Artificial Intelligence: Sequential Decisions based on Algorithmic Probability*. Springer, Berlin, 2005.
- [Lattimore and Hutter, 2011] Tor Lattimore and Marcus Hutter. Asymptotically optimal agents. In *Proc. 22nd International Conf. on Algorithmic Learning Theory (ALT'11)*, volume 6925 of *LNAI*, pages 368–382, Espoo, Finland, 2011. Springer.
- [Lattimore and Hutter, 2014a] Tor Lattimore and Marcus Hutter. Bayesian reinforcement learning with exploration. In *International Conference on Algorithmic Learning Theory*, pages 170–184. Springer, 2014.
- [Lattimore and Hutter, 2014b] Tor Lattimore and Marcus Hutter. General time consistent discounting. *Theoretical Computer Science*, 519:140–154, 2014.
- [Leike *et al.*, 2016] Jan Leike, Tor Lattimore, Laurent Orseau, and Marcus Hutter. Thompson sampling is asymptotically optimal in general environments. In *Proc. 32nd International Conf. on Uncertainty in Artificial Intelligence (UAI'16)*, pages 417–426, New Jersey, USA, 2016. AUAI Press.
- [Mnih *et al.*, 2015] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellefleur, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- [Mnih *et al.*, 2016] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937, 2016.
- [Orseau *et al.*, 2013] Laurent Orseau, Tor Lattimore, and Marcus Hutter. Universal knowledge-seeking agents for stochastic environments. In *Proc. 24th International Conf. on Algorithmic Learning Theory (ALT'13)*, volume 8139 of *LNAI*, pages 158–172, Singapore, 2013. Springer.
- [Silver and Veness, 2010] David Silver and Joel Veness. Monte-carlo planning in large pomdps. In *Advances in neural information processing systems*, pages 2164–2172, 2010.
- [Silver *et al.*, 2017] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, et al. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *arXiv preprint arXiv:1712.01815*, 2017.
- [Still, 2009] Susanne Still. Information-theoretic approach to interactive learning. *EPL (Europhysics Letters)*, 85(2):28005, 2009.
- [Whittle, 1979] Peter Whittle. Discussion of Dr Gittins' paper. *Journal of the Royal Statistical Society*, 41:164–177, 1979.

Appendices

A Definitions and Notation – Quick Reference

$$\begin{aligned}
 \mathcal{H} &:= \mathcal{A} \times \mathcal{O} \times \mathcal{R} \\
 \left. \begin{aligned}
 h_{<t} &\in \mathcal{H}^{t-1} \\
 h_{t:k} &\in \mathcal{H}^{k-t+1} \\
 \mu, \nu &\in \mathcal{M} \\
 \mu, \nu &: \mathcal{H}^* \times \mathcal{A} \rightsquigarrow \mathcal{O} \times \mathcal{R} \\
 \pi &: \mathcal{H}^* \rightsquigarrow \mathcal{A}
 \end{aligned} \right\} \text{typical meaning of certain notation} \\
 P_\nu^\pi(\epsilon) &:= 1; \quad P_\nu^\pi(haor) := P_\nu^\pi(h)\pi(a|h)\nu(or|ha) \\
 w(\nu) &:= 2^{-K(\nu)} \\
 P_\xi^\pi(\cdot) &:= \sum_{\nu \in \mathcal{M}} w(\nu) P_\nu^\pi(\cdot) \\
 w(\nu|h) &:= w(\nu) \frac{P_\nu^\pi(h)}{P_\xi^\pi(h)} \\
 w_{\text{inf}}(\mu|h_{1:\infty}) &:= \inf_{k \in \mathbb{N}} w(\mu|h_{<k}) \\
 \xi(or|ha) &:= \sum_{\nu \in \mathcal{M}} w(\nu|h)\nu(or|ha) \\
 \text{IG}(h_{t:t+k-1}|h_{<t}) &:= \sum_{\nu \in \mathcal{M}} w(\nu|h_{<t+k-1}) \log \frac{w(\nu|h_{<t+k-1})}{w(\nu|h_{<t})} \quad \text{for } h_{<t} \in \mathcal{H}^{t-1}, h' \in \mathcal{H}^k \\
 V^{\text{IG}}(\alpha^m, h_{<t}) &:= \sum_{h_{t:t+m-1} \in \mathcal{H}^m} P_\xi^{\alpha^m}(h_{<t+m}|h_{<t}) \text{IG}(h_{t:t+m-1}|h_{<t}) \\
 \alpha_{m,k}^{\text{IG}}(h_{<t}) &:= \operatorname{argmax}_{\alpha^m : \bigcup_{i=0}^{m-1} \mathcal{H}^i \rightarrow \mathcal{A}} V^{\text{IG}}(\alpha^m, h_{<t-k}) \\
 a_{m,k}^{\text{IG}}(h_{<t}) &:= \alpha_{m,k}^{\text{IG}}(h_{<t})(h_{t-k:t-1}) \\
 V_\nu^\pi(h_{<t}) &:= \frac{1}{\Gamma_t} \mathbb{E}_\nu^\pi \left[\sum_{k=t}^{\infty} \gamma_k r_k | h_{<t} \right] \\
 a^*(h_{<t}) &:= \pi_\xi^*(h_{<t}) \\
 V_\nu^*(h_{<t}) &:= \sup_{\pi \in \Pi} V_\nu^\pi(h_{<t}) = V_\nu^{\pi^*}(h_{<t}) \\
 \rho(h_{<t}, m, k) &:= \max \left\{ \frac{1}{m^2(m+1)}, \eta V^{\text{IG}}(\alpha_{m,k}^{\text{IG}}(h_{<t}), h_{<t-k}) \right\} \\
 \pi^\dagger(a|h_{<t}) &:= \sum_{m \in \mathbb{N}} \sum_{k < m, t} \rho(h_{<t}, m, k) [[a = a_{m,k}^{\text{IG}}(h_{<t})]] + (1 - \beta(h_{<t}, m, k)) [[a = a^*(h_{<t})]] \\
 \text{KL}_{h_{<t}, n}(P_{\nu_1}^\pi || P_{\nu_2}^\pi) &:= \sum_{h' \in \mathcal{H}^n} P_{\nu_1}^\pi(h'|h_{<t}) \log \frac{P_{\nu_1}^\pi(h'|h_{<t})}{P_{\nu_2}^\pi(h'|h_{<t})}
 \end{aligned}$$

B Proofs of Lemmas

We begin with a lemma that equates the information gain value of an expedition with the expected prediction error. The KL-divergence on the right hand side represents how different ν and ξ appear when following the expedition in question.

Lemma 4.

$$V^{\text{IG}}(\alpha^m, h_{<t}) = \sum_{\nu \in \mathcal{M}} w(\nu|h_{<t}) \text{KL}_{h_{<t}, m} \left(\mathbb{P}_\nu^{\alpha^m} \parallel \mathbb{P}_\xi^{\alpha^m} \right)$$

Proof. This result is shown in [Orseau *et al.*, 2013, Equation 4]. \square

Recall that $w(\nu|h_{<t})$ is the posterior weight that Inq assigns to the environment ν after observing $h_{<t}$. We show that the infimum of this value is strictly positive with probability 1.

Lemma 5. $\inf_t w(\mu|h_{<t}) > 0$ w. \mathbb{P}_μ^π -p. 1

Proof. Suppose $\inf_t w(\mu|h_{<t}) = 0$. $w(\mu|h_{<t}) > 0$ for all histories generated by \mathbb{P}_μ^π . Therefore, $\inf_t w(\mu|h_{<t}) = 0 \implies \liminf_{t \rightarrow \infty} w(\mu|h_{<t}) = 0$, and $\limsup_{t \rightarrow \infty} w(\mu|h_{<t})^{-1} = \infty$. We show that this has probability 0.

Let

$$z_t := w(\mu|h_{\leq t})^{-1} = \frac{\mathbb{P}_\xi^\pi(h_{\leq t})}{\mathbb{P}_\mu^\pi(h_{\leq t})} w(\mu)^{-1} \quad (18)$$

I first show that z_t is a μ -martingale.

$$\begin{aligned} \mathbb{E}_\mu^\pi[z_t|h_{<t}] &= w(\mu)^{-1} \sum_{h_t \in \mathcal{H}} \mathbb{P}_\mu^\pi(h_t|h_{<t}) \frac{\mathbb{P}_\xi^\pi(h_{\leq t})}{\mathbb{P}_\mu^\pi(h_{\leq t})} \\ &= w(\mu)^{-1} \frac{\sum_{h_t \in \mathcal{H}} \mathbb{P}_\xi^\pi(h_{\leq t})}{\mathbb{P}_\mu^\pi(h_{<t})} \\ &= w(\mu)^{-1} \frac{\mathbb{P}_\xi^\pi(h_{<t})}{\mathbb{P}_\mu^\pi(h_{<t})} \\ &= z_{t-1} \end{aligned} \quad (19)$$

By the martingale convergence theorem $z_t \rightarrow f(\omega) < \infty$ w.p.1, for $\omega \in \Omega$, the sample space, and some $f : \Omega \rightarrow \mathbb{R}$. Therefore, $\inf_t w(\mu|h_{<t}) > 0$ w.p.1. \square

Next we show that every exploration probability $\rho(h_{<t}, m, k)$ goes to 0. From here, all “w.p.1” statements mean with \mathbb{P}_μ^π -probability 1, if not otherwise specified.

Lemma 6.

$$\rho(h_{<t}, m, k) \xrightarrow{t \rightarrow \infty} 0 \text{ w.p.1}$$

Proof. $\rho(h_{<t}, m, k) = \rho(h_{<t-k}, m, 0)$, so we need only show that $\rho(h_{<t}, m, 0) \rightarrow 0$ w.p.1. We do this by showing that the expectation of $\rho(h_{<t}, m, 0)^{m+1}$ is summable. (This is a stronger result, since it implies that it is summable with probability 1, so the probability that it is greater than ε infinitely often is 0.) A bit of notational background: $0 \in \mathbb{N}$, and $m\mathbb{N} + i = \{i, i + m, i + 2m, \dots\}$. Each equation and inequality is explained below.

$$\begin{aligned} & w(\mu) \mathbb{E}_\mu^{\pi^\dagger} \sum_{t \in m\mathbb{N}+i} \rho(h_{<t}, m, 0)^{m+1} \\ & \stackrel{(a)}{\leq} \sum_{\nu \in \mathcal{M}} w(\nu) \mathbb{E}_\nu^{\pi^\dagger} \sum_{t \in m\mathbb{N}+i} \rho(h_{<t}, m, 0)^{m+1} \\ & \stackrel{(b)}{=} \mathbb{E}_\xi^{\pi^\dagger} \sum_{t \in m\mathbb{N}+i} \rho(h_{<t}, m, 0)^{m+1} \\ & \stackrel{(c)}{\leq} \mathbb{E}_\xi^{\pi^\dagger} \sum_{t \in m\mathbb{N}+i} \rho(h_{<t}, m, 0)^m \eta V^{\text{IG}}(\alpha_{m,0}^{\text{IG}}, h_{<t}) \end{aligned}$$

$$\begin{aligned}
&\stackrel{(d)}{=} \eta \sum_{t \in m\mathbb{N}+i} \mathbb{E}_{h_{<t} \sim P_{\xi}^{\pi^\dagger}} \left[\rho(h_{<t}, m, 0)^m \mathbb{E}_{h_{t:t+m-1} \sim P_{\xi}^{\alpha_{m,0}^{\text{IG}}}} [\text{IG}(h_{t:t+m-1}|h_{<t})] \right] \\
&\stackrel{(e)}{=} \eta \sum_{t \in m\mathbb{N}+i} \mathbb{E}_{h_{<t} \sim P_{\xi}^{\pi^\dagger}} \left[\sum_{h_{t:t+m-1} \in \mathcal{H}^m} \rho(h_{<t}, m, 0)^m P_{\xi}^{\alpha_{m,0}^{\text{IG}}}(h_{t:t+m-1}) [\text{IG}(h_{t:t+m-1}|h_{<t})] \right] \\
&\stackrel{(f)}{\leq} \eta \sum_{t \in m\mathbb{N}+i} \mathbb{E}_{h_{<t} \sim P_{\xi}^{\pi^\dagger}} \left[\sum_{h_{t:t+m-1} \in \mathcal{H}^m} P_{\xi}^{\pi^\dagger}(h_{t:t+m-1}) [\text{IG}(h_{t:t+m-1}|h_{<t})] \right] \\
&\stackrel{(g)}{=} \eta \sum_{t \in m\mathbb{N}+i} \mathbb{E}_{\xi}^{\pi^\dagger} \text{IG}(h_{t:t+m-1}|h_{<t}) \\
&\stackrel{(h)}{=} \eta \mathbb{E}_{\xi}^{\pi^\dagger} \sum_{t \in m\mathbb{N}+i} \sum_{\nu \in \mathcal{M}} w(\nu|h_{<t+m}) \log \frac{w(\nu|h_{<t+m})}{w(\nu|h_{<t})} \\
&\stackrel{(i)}{=} \eta \sum_{t \in m\mathbb{N}+i} \sum_{\nu \in \mathcal{M}} \mathbb{E}_{\xi}^{\pi^\dagger} \frac{w(\nu)\nu(\alpha_{<t+m}|a_{<t+m})}{\xi(\alpha_{<t+m}|a_{<t+m})} \log \frac{w(\nu|h_{<t+m})}{w(\nu|h_{<t})} \\
&\stackrel{(j)}{=} \eta \sum_{t \in m\mathbb{N}+i} \sum_{\nu \in \mathcal{M}} \mathbb{E}_{\nu}^{\pi^\dagger} w(\nu) \log \frac{w(\nu|h_{<t+m})}{w(\nu|h_{<t})} \\
&\stackrel{(k)}{=} \lim_{N \rightarrow \infty} \eta \sum_{k=0}^{N-1} \sum_{\nu \in \mathcal{M}} \mathbb{E}_{\nu}^{\pi^\dagger} w(\nu) \log \frac{w(\nu|h_{<mk+i+m})}{w(\nu|h_{<mk+i})} \\
&\stackrel{(l)}{=} \lim_{N \rightarrow \infty} \eta \sum_{\nu \in \mathcal{M}} \mathbb{E}_{\nu}^{\pi^\dagger} w(\nu) \log \prod_{k=0}^{N-1} \frac{w(\nu|h_{<m(k+1)+i})}{w(\nu|h_{<mk+i})} \\
&\stackrel{(m)}{=} \lim_{N \rightarrow \infty} \eta \sum_{\nu \in \mathcal{M}} \mathbb{E}_{\nu}^{\pi^\dagger} w(\nu) \log \frac{w(\nu|h_{<mN+i})}{w(\nu|h_{<i})} \\
&\stackrel{(n)}{\leq} \eta \sum_{\nu \in \mathcal{M}} \mathbb{E}_{\nu}^{\pi^\dagger} w(\nu) \log \frac{1}{w(\nu|h_{<i})} \\
&\stackrel{(o)}{=} \eta \sum_{\nu \in \mathcal{M}} \mathbb{E}_{\nu}^{\pi^\dagger} w(\nu) \log \frac{1}{w(\nu)} \frac{w(\nu)}{w(\nu|h_{<i})} \\
&\stackrel{(p)}{=} \eta \sum_{\nu \in \mathcal{M}} w(\nu) \log \frac{1}{w(\nu)} + \eta \sum_{\nu \in \mathcal{M}} \mathbb{E}_{\nu}^{\pi^\dagger} w(\nu) \log \frac{w(\nu)}{w(\nu|h_{<i})} \\
&\stackrel{(q)}{=} \eta \text{Ent}(w) + \eta \sum_{h_{<i} \in \mathcal{H}^i} \sum_{\nu \in \mathcal{M}} w(\nu) P_{\nu}^{\pi^\dagger}(h_{<i}) \log \frac{w(\nu)}{w(\nu|h_{<i})} \\
&\stackrel{(r)}{=} \eta \text{Ent}(w) + \eta \sum_{h_{<i} \in \mathcal{H}^i} \sum_{\nu \in \mathcal{M}} w(\nu|h_{<i}) P_{\xi}^{\pi^\dagger}(h_{<i}) \log \frac{w(\nu)}{w(\nu|h_{<i})} \\
&\stackrel{(s)}{=} \eta \text{Ent}(w) - \eta \mathbb{E}_{\xi}^{\pi^\dagger} [\text{IG}(h_{<i}|\epsilon)] \stackrel{(t)}{\leq} \eta \text{Ent}(w) \stackrel{(u)}{<} \infty
\end{aligned} \tag{20}$$

For multiple steps in this derivation, note that the information gain is non-negative; this is a property of the KL-divergence. (a) follows from the l.h.s. being one of the non-negative summands of the r.h.s. (b) follows from the definition of ξ . (c) follows from the definition of ρ . (d) substitutes V^{IG} for its definition. (e) expands the definition of the expectation. (f) follows because π^\dagger mimics $\alpha_{m,0}^{\text{IG}}$ for m consecutive timesteps with probability $\prod_{i=0}^{m-1} \rho(h_{<t+i}, m, i) = \rho(h_{<t}, m, 0)^m$, so the probability of any history under $P_{\xi}^{\pi^\dagger}$ is at least the probability of that history under $P_{\xi}^{\alpha_{m,0}^{\text{IG}}}$ times $\rho(h_{<t}, m, 0)^m$. (g) combines the two expectations, which are now with respect to the same probability measure. (h) expands the definition of the information gain. (i) rearranges the expectations and the sums, and expands $w(\nu|h_{<t+m})$ according to Bayes' rule. (j) converts the expectation to an 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. And (u) is shown in [Orseau *et al.*, 2013, Proposition 13].

Finally,

$$\begin{aligned} \mathbb{E}_\mu^{\pi^\dagger} \sum_{t=0}^{\infty} \rho(h_{<t}, m, 0)^{m+1} &= \sum_{i=0}^{m-1} \mathbb{E}_\mu^{\pi^\dagger} \sum_{t \in m\mathbb{N}+i} \rho(h_{<t}, m, 0)^{m+1} \\ &\stackrel{(20)}{\leq} \sum_{i=0}^{m-1} \frac{\eta \text{Ent}(w)}{w(\mu)} = \frac{m\eta \text{Ent}(w)}{w(\mu)} < \infty \end{aligned} \quad (21)$$

□

Now, we show that the total exploration probability goes to 0:

Lemma 7.

$$\beta(h_{<t}) \rightarrow 0 \text{ w.p.1}$$

Proof.

$$\begin{aligned} \beta(h_{<t}) &= \sum_{m \in \mathbb{N}} \sum_{k=0}^{\min\{m-1, t\}} \rho(h_{<t}, m, k) \\ &= \sum_{m \in \mathbb{N}} \sum_{k=0}^{\min\{m-1, t\}} \min \left\{ \frac{1}{m^2(m+1)}, V_{m,k}^{\text{IG}}(h_{<t}) \right\} \end{aligned} \quad (22)$$

Each of the terms in the sum approaches 0 with probability 1 by Lemma 6, and because $\rho(h_{<t}, m, k) = \rho(h_{<t-k}, m, 0)$. Suppose by contradiction $\beta(h_{<t}) > \varepsilon > 0$ infinitely often. There exists an M such that

$$\sum_{m=M}^{\infty} \sum_{k=0}^{\min\{m-1, t\}} \rho(h_{<t}, m, k) < \sum_{m=M}^{\infty} \sum_{k=0}^{m-1} \frac{1}{m^2(m+1)} < \varepsilon/2 \quad (23)$$

for all t . With that M , then if $\beta(h_{<t}) > \varepsilon$ infinitely often, it must be the case that $\sum_{m=0}^{M-1} \sum_{k=0}^{m-1} \rho(h_{<t}, m, k) > \varepsilon/2$ infinitely often, but this is a finite sum of terms that all approach 0, a contradiction. □

Lemma 8 shows that the probabilities assigned by ξ converge to those of μ .

Lemma 8. $\forall m \in \mathbb{N}, h_{t:t+m-1} \in \mathcal{H}^m, \alpha^m : \bigcup_{i=0}^{m-1} \mathcal{H}^i \rightarrow \mathcal{A}$:

$$\mathbb{P}_\mu^{\alpha^m}(h_{t:t+m-1}|h_{<t}) - \mathbb{P}_\xi^{\alpha^m}(h_{t:t+m-1}|h_{<t}) \xrightarrow{t \rightarrow \infty} 0 \text{ w.p.1}$$

Proof. Suppose that $0 < \varepsilon \leq (\mathbb{P}_\mu^{\alpha^m}(h_{t:t+m-1}|h_{<t}) - \mathbb{P}_\xi^{\alpha^m}(h_{t:t+m-1}|h_{<t}))^2$ for some $h_{t:t+m-1}$.

$$\begin{aligned} \varepsilon &\leq (\mathbb{P}_\mu^{\alpha^m}(h_{t:t+m-1}|h_{<t}) - \mathbb{P}_\xi^{\alpha^m}(h_{t:t+m-1}|h_{<t}))^2 \\ &\stackrel{(a)}{\leq} \text{KL}_{h_{<t}, m}(\mathbb{P}_\mu^{\alpha^m} \parallel \mathbb{P}_\xi^{\alpha^m}) \\ &\stackrel{(b)}{\leq} \sum_{\nu \in \mathcal{M}} \frac{w(\nu|h_{<t})}{w(\mu|h_{<t})} \text{KL}_{h_{<t}, m}(\mathbb{P}_\nu^{\alpha^m} \parallel \mathbb{P}_\xi^{\alpha^m}) \\ &\stackrel{(c)}{=} \frac{1}{w(\mu|h_{<t})} V^{\text{IG}}(\alpha^m, h_{<t}) \\ &\stackrel{(d)}{\leq} \frac{1}{\inf_k w(\mu|h_{<k})} V^{\text{IG}}(\alpha^m, h_{<t}) \\ &\stackrel{(e)}{\leq} \frac{1}{\inf_k w(\mu|h_{<k})} V^{\text{IG}}(\alpha_{m,0}^{\text{IG}}(h_{<t}), h_{<t}) \end{aligned} \quad (24)$$

(a) is a result from information theory known as the entropy inequality. (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 Lemma 4. (d) follows from the definition of the infimum. And (e) follows from the fact that $\alpha_{m,0}^{\text{IG}}(h_{<t})$ maximizes $V^{\text{IG}}(\cdot, h_{<t})$, by definition.

Therefore,

$$\begin{aligned}
& (\mathbb{P}_\mu^{\alpha^m}(h_{t:t+m-1}|h_{<t}) - \mathbb{P}_\xi^{\alpha^m}(h_{t:t+m-1}|h_{<t}))^2 \geq \varepsilon \text{ i.o.} \\
& \text{implies } V^{\text{IG}}(\alpha_{m,0}^{\text{IG}}(h_{<t}), h_{<t}) \geq \varepsilon \inf_k w(\mu|h_{<k}) \text{ i.o.} \\
& \text{which implies } \rho(h_{<t}, m, 0) \geq \min\left\{\frac{1}{m^2(m+1)}, \varepsilon \inf_k w(\mu|h_{<k})\right\} \text{ i.o.} \\
& \text{which implies } \sum_{t=0}^{\infty} \rho(h_{<t}, m, 0)^{m+1} = \infty \text{ or } \inf_k w(\mu|h_{<k}) = 0
\end{aligned}$$

This has probability 0 by Lemmas 6 and 5. Thus, with probability 1, $\mathbb{P}_\mu^{\alpha^m}(h_{t:t+m-1}|h_{<t}) - \mathbb{P}_\xi^{\alpha^m}(h_{t:t+m-1}|h_{<t}) \rightarrow 0$. \square

C Approximation of Inq

Following Aslanides [2017], our approximation of Inq calls ρUCT [Silver and Veness, 2010] as a subroutine in place of expectimax.

Algorithm 2 Approximation of Inquisitive Reinforcement Learner's Policy

Require: MCTS Samples, horizon, γ

Initialize: uniform prior over model class

- 1: **while** True **do**
 - 2: **for all** $m \leq \text{horizon}$ and $k < \min\{m, t\}$ **do**
 - 3: using information gain as reward, $a_{m,k}^{\text{IG}} \sim \rho\text{UCT}(h_{<t-k}, \text{MCTS samples}, m, \gamma)$
 - 4: $\rho(m, k) = \min\{\text{information-gain-value of } a_{m,k}^{\text{IG}}, 1/(m^2(m+1))\}$
 - 5: using the actual reward, $a^* \sim \rho\text{UCT}(h_{<t}, \text{MCTS samples}, \text{horizon}, \gamma)$
 - 6: take action $a_{m,k}^{\text{IG}}$ with probability $\rho(m, k)$ for all $m \leq \text{horizon}$ and $k < \min\{m, t\}$ else take action a^*
 - 7: update posterior from observation and reward
-