

---

# On the Computability of AIXI

---

**Jan Leike**

Australian National University  
jan.leike@anu.edu.au

**Marcus Hutter**

Australian National University  
marcus.hutter@anu.edu.au

## Abstract

How could we solve the machine learning and the artificial intelligence problem if we had infinite computation? Solomonoff induction and the reinforcement learning agent AIXI are proposed answers to this question. Both are known to be incomputable. In this paper, we quantify this using the arithmetical hierarchy, and prove upper and corresponding lower bounds for incomputability. We show that AIXI is not limit computable, thus it cannot be approximated using finite computation. Our main result is a limit-computable  $\varepsilon$ -optimal version of AIXI with infinite horizon that maximizes expected rewards.

**Keywords.** AIXI, Solomonoff induction, general reinforcement learning, computability, complexity, arithmetical hierarchy, universal Turing machine.

## 1 INTRODUCTION

Given infinite computation power, many traditional AI problems become trivial: playing chess, go, or backgammon can be solved by exhaustive expansion of the game tree. Yet other problems seem difficult still; for example, predicting the stock market, driving a car, or babysitting your nephew. How can we solve these problems in theory? A proposed answer to this question is the agent AIXI [Hut00, Hut05]. As a *reinforcement learning agent*, its goal is to maximize cumulative (discounted) rewards obtained from the environment [SB98].

The basis of AIXI is Solomonoff's theory of learning [Sol64, Sol78, LV08], also called *Solomonoff induction*. It arguably solves the induction problem [RH11]: for data drawn from a computable measure  $\mu$ , Solomonoff induction will converge to the correct belief about any hypothesis [BD62, RH11]. Moreover, convergence is extremely fast in the sense that Solomonoff induction will make a total of at most  $E + O(\sqrt{E})$  errors when predicting

the next data points, where  $E$  is the number of errors of the informed predictor that knows  $\mu$  [Hut01]. While learning the environment according to Solomonoff's theory, AIXI selects actions by running an expectimax-search for maximum cumulative discounted rewards. It is clear that AIXI can only serve as an ideal, yet recently it has inspired some impressive applications [VNH<sup>+</sup>11].

Both Solomonoff induction and AIXI are known to be incomputable. But not all incomputabilities are equal. The *arithmetical hierarchy* specifies different levels of computability based on *oracle machines*: each level in the arithmetical hierarchy is computed by a Turing machine which may query a halting oracle for the respective lower level.

We posit that any ideal for a 'perfect agent' needs to be *limit computable* ( $\Delta_2^0$ ). The class of limit computable functions is the class of functions that admit an *anytime algorithm*. It is the highest level of the arithmetical hierarchy which can be approximated using a regular Turing machine. If this criterion is not met, our model would be useless to guide practical research.

For MDPs, planning is already P-complete for finite and infinite horizons [PT87]. In POMDPs, planning is undecidable [MHC99, MHC03]. The existence of a policy whose expected value exceeds a given threshold is PSPACE-complete [MGLA00], even for purely epistemic POMDPs in which actions do not change the hidden state [SLR07]. In this paper we derive hardness results for planning in general semicomputable environments; this environment class is even more general than POMDPs. We show that finding an optimal policy is  $\Pi_2^0$ -hard and finding an  $\varepsilon$ -optimal policy is undecidable.

Moreover, we show that by default, AIXI is not limit computable. The reason is twofold: First, when picking the next action, two or more actions might have the same value (expected future rewards). The choice between them is easy, but determining whether such a tie exists is difficult. Second, in case of an infinite horizon (using discounting), the iterative definition of the value function [Hut05, Def. 5.30] conditions on surviving forever. The first problem

Model	$\gamma$	Optimal	$\varepsilon$ -Optimal
Iterative AINU	DC	$\Delta_4^0, \Sigma_3^0$ -hard	$\Delta_3^0, \Pi_2^0$ -hard
	LT	$\Delta_3^0, \Pi_2^0$ -hard	$\Delta_2^0, \Sigma_1^0$ -hard
Iterative AIXI	DC	$\Delta_4^0, \Pi_2^0$ -hard	$\Delta_3^0, \Pi_2^0$ -hard
	LT	$\Delta_3^0, \Sigma_1^0$ -hard	$\Delta_2^0, \Sigma_1^0$ -hard
Iterative AIMU	DC	$\Delta_2^0$	$\Delta_1^0$
	LT	$\Delta_2^0$	$\Delta_1^0$
Recursive AINU	DC	$\Delta_3^0, \Pi_2^0$ -hard	$\Delta_2^0, \Sigma_1^0$ -hard
	LT	$\Delta_3^0, \Pi_2^0$ -hard	$\Delta_2^0, \Sigma_1^0$ -hard
Recursive AIXI	DC	$\Delta_3^0, \Sigma_1^0$ -hard	$\Delta_2^0, \Sigma_1^0$ -hard
	LT	$\Delta_3^0, \Sigma_1^0$ -hard	$\Delta_2^0, \Sigma_1^0$ -hard
Recursive AIMU	DC	$\Delta_2^0$	$\Delta_1^0$
	LT	$\Delta_2^0$	$\Delta_1^0$

Table 1: Computability results for different agent models derived in Section 3. DC means general discounting, a lower semicomputable discount function  $\gamma$ ; LT means finite lifetime, undiscounted rewards up to a fixed lifetime  $m$ . Hardness results for AIXI are with respect to a specific universal Turing machine; hardness results for AINU are with respect to a specific environment  $\nu \in \mathcal{M}$ .

can be circumvented by settling for an  $\varepsilon$ -optimal agent. We show that the second problem can be solved by using the recursive instead of the iterative definition of the value function. With this we get a limit-computable agent with infinite horizon. Table 1 and Table 3 summarize our computability results.

## 2 PRELIMINARIES

### 2.1 THE ARITHMETICAL HIERARCHY

A set  $A \subseteq \mathbb{N}$  is  $\Sigma_n^0$  iff there is a computable relation  $S$  such that

$$k \in A \iff \exists k_1 \forall k_2 \dots Q_n k_n S(k, k_1, \dots, k_n) \quad (1)$$

where  $Q_n = \forall$  if  $n$  is even,  $Q_n = \exists$  if  $n$  is odd [Nie09, Def. 1.4.10]. A set  $A \subseteq \mathbb{N}$  is  $\Pi_n^0$  iff its complement  $\mathbb{N} \setminus A$  is  $\Sigma_n^0$ . We call the formula on the right hand side of (1) a  $\Sigma_n^0$ -formula, its negation is called  $\Pi_n^0$ -formula. It can be shown that we can add any bounded quantifiers and duplicate quantifiers of the same type without changing the classification of  $A$ . The set  $A$  is  $\Delta_n^0$  iff  $A$  is  $\Sigma_n^0$  and  $A$  is  $\Pi_n^0$ . We get that  $\Sigma_1^0$  as the class of recursively enumerable sets,  $\Pi_1^0$  as the class of co-recursively enumerable sets and  $\Delta_1^0$  as the class of recursive sets.

We say the set  $A \subseteq \mathbb{N}$  is  $\Sigma_n^0$ -hard ( $\Pi_n^0$ -hard,  $\Delta_n^0$ -hard) iff for any set  $B \in \Sigma_n^0$  ( $B \in \Pi_n^0$ ,  $B \in \Delta_n^0$ ),  $B$  is many-one reducible to  $A$ , i.e., there is a computable function  $f$  such that  $k \in B \leftrightarrow f(k) \in A$  [Nie09, Def. 1.2.1]. We get  $\Sigma_n^0 \subset$

$\Delta_{n+1}^0 \subset \Sigma_{n+1}^0 \subset \dots$  and  $\Pi_n^0 \subset \Delta_{n+1}^0 \subset \Pi_{n+1}^0 \subset \dots$ . This hierarchy of subsets of natural numbers is known as the *arithmetical hierarchy*.

By Post's Theorem [Nie09, Thm. 1.4.13], a set is  $\Sigma_n^0$  if and only if it is recursively enumerable on an oracle machine with an oracle for a  $\Sigma_{n-1}^0$ -complete set.

### 2.2 STRINGS

Let  $\mathcal{X}$  be some finite set called *alphabet*. The set  $\mathcal{X}^* := \bigcup_{n=0}^{\infty} \mathcal{X}^n$  is the set of all finite strings over the alphabet  $\mathcal{X}$ , the set  $\mathcal{X}^\infty$  is the set of all infinite strings over the alphabet  $\mathcal{X}$ , and the set  $\mathcal{X}^\# := \mathcal{X}^* \cup \mathcal{X}^\infty$  is their union. The empty string is denoted by  $\epsilon$ , not to be confused with the small positive real number  $\varepsilon$ . Given a string  $x \in \mathcal{X}^*$ , we denote its length by  $|x|$ . For a (finite or infinite) string  $x$  of length  $\geq k$ , we denote with  $x_{1:k}$  the first  $k$  characters of  $x$ , and with  $x_{<k}$  the first  $k-1$  characters of  $x$ . The notation  $x_{1:\infty}$  stresses that  $x$  is an infinite string. We write  $x \sqsubseteq y$  iff  $x$  is a prefix of  $y$ , i.e.,  $x = y_{1:|x|}$ .

### 2.3 COMPUTABILITY OF REAL-VALUED FUNCTIONS

We fix some encoding of rational numbers into binary strings and an encoding of binary strings into natural numbers. From now on, this encoding will be done implicitly wherever necessary.

**Definition 1** ( $\Sigma_n^0$ -,  $\Pi_n^0$ -,  $\Delta_n^0$ -computable). A function  $f : \mathcal{X}^* \rightarrow \mathbb{R}$  is called  $\Sigma_n^0$ -computable ( $\Pi_n^0$ -computable,  $\Delta_n^0$ -computable) iff the set  $\{(x, q) \in \mathcal{X}^* \times \mathbb{Q} \mid f(x) > q\}$  is  $\Sigma_n^0$  ( $\Pi_n^0$ ,  $\Delta_n^0$ ).

A  $\Delta_1^0$ -computable function is called *computable*, a  $\Sigma_1^0$ -computable function is called *lower semicomputable*, and a  $\Pi_1^0$ -computable function is called *upper semicomputable*. A  $\Delta_2^0$ -computable function  $f$  is called *limit computable*, because there is a computable function  $\phi$  such that

$$\lim_{k \rightarrow \infty} \phi(x, k) = f(x).$$

The program  $\phi$  that limit computes  $f$  can be thought of as an *anytime algorithm* for  $f$ : we can stop  $\phi$  at any time  $k$  and get a preliminary answer. If the program  $\phi$  ran long enough (which we do not know), this preliminary answer will be close to the correct one.

Limit-computable sets are the highest level in the arithmetical hierarchy that can be approached by a regular Turing machine. Above limit-computable sets we necessarily need some form of halting oracle. See Table 2 for the definition of lower/upper semicomputable and limit-computable functions in terms of the arithmetical hierarchy.

**Lemma 2** (Computability of Arithmetical Operations). *Let  $n > 0$  and let  $f, g : \mathcal{X}^* \rightarrow \mathbb{R}$  be two  $\Delta_n^0$ -computable functions. Then*

	$f_>$	$f_<$
$f$ is computable	$\Delta_1^0$	$\Delta_1^0$
$f$ is lower semicomputable	$\Sigma_1^0$	$\Pi_1^0$
$f$ is upper semicomputable	$\Pi_1^0$	$\Sigma_1^0$
$f$ is limit computable	$\Delta_2^0$	$\Delta_2^0$
$f$ is $\Delta_n^0$ -computable	$\Delta_n^0$	$\Delta_n^0$
$f$ is $\Sigma_n^0$ -computable	$\Sigma_n^0$	$\Pi_n^0$
$f$ is $\Pi_n^0$ -computable	$\Pi_n^0$	$\Sigma_n^0$

Table 2: Connection between the computability of real-valued functions and the arithmetical hierarchy. We use the shorthand  $f_> := \{(x, q) \mid f(x) > q\}$  and  $f_< := \{(x, q) \mid f(x) < q\}$ .

- (i)  $\{(x, y) \mid f(x) > g(y)\}$  is  $\Sigma_n^0$ ,
- (ii)  $\{(x, y) \mid f(x) \leq g(y)\}$  is  $\Pi_n^0$ ,
- (iii)  $f + g$ ,  $f - g$ , and  $f \cdot g$  are  $\Delta_n^0$ -computable, and
- (iv)  $f/g$  is  $\Delta_n^0$ -computable if  $g(x) \neq 0$  for all  $x$ .

## 2.4 ALGORITHMIC INFORMATION THEORY

A *semimeasure* over the alphabet  $\mathcal{X}$  is a function  $\nu : \mathcal{X}^* \rightarrow [0, 1]$  such that (i)  $\nu(\epsilon) \leq 1$ , and (ii)  $\nu(x) \geq \sum_{a \in \mathcal{X}} \nu(xa)$  for all  $x \in \mathcal{X}^*$ . A semimeasure is called (probability) *measure* iff for all  $x$  equalities hold in (i) and (ii). *Solomonoff's prior*  $M$  [Sol64] assigns to a string  $x$  the probability that the reference universal monotone Turing machine  $U$  [LV08, Ch. 4.5.2] computes a string starting with  $x$  when fed with uniformly random bits as input. The *measure mixture*  $\bar{M}$  [Gá83, p. 74] removes the contribution of programs that do not compute infinite strings; it is a measure except for a constant factor. Formally,

$$M(x) := \sum_{p: x \sqsubseteq U(p)} 2^{-|p|}, \quad \bar{M}(x) := \lim_{n \rightarrow \infty} \sum_{y \in \mathcal{X}^n} M(xy)$$

Equivalently, the Solomonoff prior  $M$  can be defined as a mixture over all lower semicomputable semimeasures [WSH11]. The function  $M$  is a lower semicomputable semimeasure, but not computable and not a measure [LV08, Lem. 4.5.3]. A semimeasure  $\nu$  can be turned into a measure  $\nu_{\text{norm}}$  using *Solomonoff normalization*:  $\nu_{\text{norm}}(\epsilon) := 1$  and for all  $x \in \mathcal{X}^*$  and  $a \in \mathcal{X}$ ,

$$\nu_{\text{norm}}(xa) := \nu_{\text{norm}}(x) \frac{\nu(xa)}{\sum_{b \in \mathcal{X}} \nu(xb)}. \quad (2)$$

## 2.5 GENERAL REINFORCEMENT LEARNING

In general reinforcement learning the agent interacts with an environment in cycles: at time step  $t$  the agent chooses an *action*  $a_t \in \mathcal{A}$  and receives a *percept*  $e_t = (o_t, r_t) \in \mathcal{E}$  consisting of an *observation*  $o_t \in \mathcal{O}$  and a real-valued

*reward*  $r_t \in \mathbb{R}$ ; the cycle then repeats for  $t + 1$ . A *history* is an element of  $(\mathcal{A} \times \mathcal{E})^*$ . We use  $\mathfrak{x} \in \mathcal{A} \times \mathcal{E}$  to denote one interaction cycle, and  $\mathfrak{x}_{1:t}$  to denote a history of length  $t$ . The goal in reinforcement learning is to maximize total discounted rewards. A *policy* is a function  $\pi : (\mathcal{A} \times \mathcal{E})^* \rightarrow \mathcal{A}$  mapping each history to the action taken after seeing this history.

The environment can be stochastic, but is assumed to be semicomputable. In accordance with the AIXI literature [Hut05], we model environments as lower semicomputable *chronological conditional semimeasures* (LSC-CCSs). A *conditional semimeasure*  $\nu$  takes a sequence of actions  $a_{1:t}$  as input and returns a semimeasure  $\nu(\cdot \parallel a_{1:t})$  over  $\mathcal{E}^\#$ . A conditional semimeasure  $\nu$  is *chronological* iff percepts at time  $t$  do not depend on future actions, i.e.,  $\nu(e_{1:t} \parallel a_{1:k}) = \nu(e_{1:t} \parallel a_{1:t})$  for all  $k > t$ . Despite their name, conditional semimeasures do *not* specify conditional probabilities; the environment  $\nu$  is *not* a joint probability distribution on actions and percepts. Here we only care about the computability of the environment  $\nu$ ; for our purposes, chronological conditional semimeasures behave just like semimeasures.

## 2.6 THE UNIVERSAL AGENT AIXI

Our environment class  $\mathcal{M}$  is the class of all LSC-CCSs. Typically, Bayesian agents such as AIXI only function well if the true environment is in their hypothesis class. Since the hypothesis class  $\mathcal{M}$  is extremely large, the assumption that it contains the true environment is rather weak. We fix the *universal prior*  $(w_\nu)_{\nu \in \mathcal{M}}$  with  $w_\nu > 0$  for all  $\nu \in \mathcal{M}$  and  $\sum_{\nu \in \mathcal{M}} w_\nu \leq 1$ , given by the reference machine  $U$ . The universal prior  $w$  gives rise to the *universal mixture*  $\xi$ , which is a convex combination of all LSC-CCSs  $\mathcal{M}$ :

$$\xi(e_{<t} \parallel a_{<t}) := \sum_{\nu \in \mathcal{M}} w_\nu \nu(e_{<t} \parallel a_{<t})$$

It is analogous to the Solomonoff prior  $M$  but defined for reactive environments. Like  $M$ , the universal mixture  $\xi$  is lower semicomputable [Hut05, Sec. 5.10].

We fix a *discount function*  $\gamma : \mathbb{N} \rightarrow \mathbb{R}$  with  $\gamma_t := \gamma(t) \geq 0$  and  $\sum_{t=1}^\infty \gamma_t < \infty$  and make the following assumptions.

- Assumption 3.** (a) *The discount function  $\gamma$  is lower semicomputable.*
- (b) *Rewards are bounded between 0 and 1.*
- (c) *The set of actions  $\mathcal{A}$  and the set of percepts  $\mathcal{E}$  are both finite.*

Assumption 3 (b) could be relaxed to bounded rewards because we can rescale rewards  $r \mapsto cr + d$  for any  $c, d \in \mathbb{R}$  without changing optimal policies if the environment  $\nu$  is a measure. However, for our value-related results, we require that rewards are nonnegative.

We define the *discount normalization factor*  $\Gamma_t := \sum_{i=t}^{\infty} \gamma_i$ . There is no requirement that  $\Gamma_t > 0$ . In fact, we use  $\gamma$  for both, AIXI with discounted infinite horizon ( $\Gamma_t > 0$  for all  $t$ ), and AIXI with finite lifetime  $m$ . In the latter case we set

$$\gamma_{\text{LT}m}(t) := \begin{cases} 1 & \text{if } t \leq m \\ 0 & \text{if } t > m. \end{cases}$$

If we knew the true environment  $\nu \in \mathcal{M}$ , we would choose the  $\nu$ -optimal agent known as AINU that maximizes  $\nu$ -expected value (if  $\nu$  is a measure). Since we do not know the true environment, we use the universal mixture  $\xi$  over all environments in  $\mathcal{M}$  instead. This yields the Bayesian agent AIXI: it weighs every environment  $\nu \in \mathcal{M}$  according to its prior probability  $w_\nu$ .

**Definition 4** (Iterative Value Function [Hut05, Def. 5.30]). The *value* of a policy  $\pi$  in an environment  $\nu$  given history  $\mathbf{x}_{<t}$  is

$$V_\nu^\pi(\mathbf{x}_{<t}) := \frac{1}{\Gamma_t} \lim_{m \rightarrow \infty} \sum_{e_{t:m}} R(e_{t:m}) \nu(e_{1:m} \mid e_{<t} \parallel a_{1:m})$$

if  $\Gamma_t > 0$  and  $V_\nu^\pi(\mathbf{x}_{<t}) := 0$  if  $\Gamma_t = 0$  where  $a_i := \pi(e_{<i})$  for all  $i \geq t$  and  $R(e_{t:m}) := \sum_{k=t}^m \gamma_k r_k$ . The *optimal value* is defined as  $V_\nu^*(h) := \sup_\pi V_\nu^\pi(h)$ .

Let  $\mathbf{x}_{<t} \in (\mathcal{A} \times \mathcal{E})^*$  be some history. We extend the value functions  $V_\nu^\pi$  to include initial interactions (in reinforcement learning literature on MDPs these are called  $Q$ -values),  $V_\nu^\pi(\mathbf{x}_{<t} a_t) := V_\nu^{\pi'}(\mathbf{x}_{<t})$  where  $\pi'$  is the policy  $\pi$  except that it takes action  $a_t$  next, i.e.,  $\pi'(\mathbf{x}_{<t}) := a_t$  and  $\pi'(h) := \pi(h)$  for all  $h \neq \mathbf{x}_{<t}$ . We define  $V_\nu^*(\mathbf{x}_{<t} a_t) := \sup_\pi V_\nu^\pi(\mathbf{x}_{<t} a_t)$  analogously.

**Definition 5** (Optimal Policy [Hut05, Def. 5.19 & 5.30]). A policy  $\pi$  is *optimal in environment  $\nu$*  ( $\nu$ -optimal) iff for all histories the policy  $\pi$  attains the optimal value:  $V_\nu^\pi(h) = V_\nu^*(h)$  for all  $h \in (\mathcal{A} \times \mathcal{E})^*$ .

Since the discount function is summable, rewards are bounded (Assumption 3b), and actions and percepts spaces are both finite (Assumption 3c), an optimal policy exists for every environment  $\nu \in \mathcal{M}$  [LH14, Thm. 10]. For a fixed environment  $\nu$ , an explicit expression for the optimal value function is

$$V_\nu^*(\mathbf{x}_{<t}) = \frac{1}{\Gamma_t} \lim_{m \rightarrow \infty} \mathop{\text{m}}\max_{\mathbf{x}_{t:m}} \sum R(e_{t:m}) \nu(e_{1:m} \mid e_{<t} \parallel a_{1:m}), \quad (3)$$

where  $\mathop{\text{m}}\max$  denotes the expectimax operator:

$$\mathop{\text{m}}\max_{\mathbf{x}_{t:m}} := \max_{a_t \in \mathcal{A}} \sum_{e_t \in \mathcal{E}} \dots \max_{a_m \in \mathcal{A}} \sum_{e_m \in \mathcal{E}}$$

For an environment  $\nu \in \mathcal{M}$  (an LSCCCS), AINU is defined as a  $\nu$ -optimal policy  $\pi_\nu^* = \arg \max_\pi V_\nu^\pi(\epsilon)$ . To

	Plain	Conditional
$M$	$\Sigma_1^0 \setminus \Delta_1^0$	$\Delta_2^0 \setminus (\Sigma_1^0 \cup \Pi_1^0)$
$M_{\text{norm}}$	$\Delta_2^0 \setminus (\Sigma_1^0 \cup \Pi_1^0)$	$\Delta_2^0 \setminus (\Sigma_1^0 \cup \Pi_1^0)$
$\overline{M}$	$\Pi_2^0 \setminus \Delta_2^0$	$\Delta_3^0 \setminus (\Sigma_2^0 \cup \Pi_2^0)$
$\overline{M}_{\text{norm}}$	$\Delta_3^0 \setminus (\Sigma_2^0 \cup \Pi_2^0)$	$\Delta_3^0 \setminus (\Sigma_2^0 \cup \Pi_2^0)$

Table 3: The complexity of the set  $\{(x, q) \in \mathcal{X}^* \times \mathbb{Q} \mid f(x) > q\}$  where  $f \in \{M, M_{\text{norm}}, \overline{M}, \overline{M}_{\text{norm}}\}$  is one of the various versions of Solomonoff’s prior. Lower bounds on the complexity of  $\overline{M}$  and  $\overline{M}_{\text{norm}}$  hold only for specific universal Turing machines.

stress that the environment is given by a measure  $\mu \in \mathcal{M}$  (as opposed to a semimeasure), we use AIMU. AIXI is defined as a  $\xi$ -optimal policy  $\pi_\xi^*$  for the universal mixture  $\xi$  [Hut05, Ch. 5]. Since  $\xi \in \mathcal{M}$  and every measure  $\mu \in \mathcal{M}$  is also a semimeasure, both AIMU and AIXI are a special case of AINU. However, AIXI is not a special case of AIMU since the mixture  $\xi$  is not a measure.

Because there can be more than one optimal policy, the definitions of AINU, AIMU and AIXI are not unique. More specifically, a  $\nu$ -optimal policy maps a history  $h$  to

$$\pi_\nu^*(h) \in \arg \max_{a \in \mathcal{A}} V_\nu^*(ha). \quad (4)$$

If there are multiple actions  $\alpha, \beta \in \mathcal{A}$  that attain the optimal value,  $V_\nu^*(h\alpha) = V_\nu^*(h\beta)$ , we say there is an *argmax tie*. Which action we settle on in case of a tie (how we break the tie) is irrelevant and can be arbitrary.

### 3 THE COMPLEXITY OF AINU, AIMU, AND AIXI

#### 3.1 THE COMPLEXITY OF SOLOMONOFF INDUCTION

AIXI uses an analogue to Solomonoff’s prior on all possible environments  $\mathcal{M}$ . Therefore we first state computability results for *Solomonoff’s prior*  $M$  and the *measure mixture*  $\overline{M}$  in Table 3 [LH15b]. Notably,  $M$  is lower semicomputable and its conditional is limit computable. However, neither the measure mixture  $\overline{M}$  nor any of its variants are limit computable.

#### 3.2 UPPER BOUNDS

In this section, we derive upper bounds on the computability of AINU, AIMU, and AIXI. Except for Corollary 13, all results in this section apply generally to any LSCCCS  $\nu \in \mathcal{M}$ , hence they apply to AIXI even though they are stated for AINU.

For a fixed lifetime  $m$ , only the first  $m$  interactions matter. There is a finite number of policies that are different for

the first  $m$  interactions, and the optimal policy  $\pi_\xi^*$  can be encoded in a finite number of bits and is thus computable. To make a meaningful statement about the computability of  $\text{AINU}_{\text{LT}}$ , we have to consider it as the function that takes the lifetime  $m$  and outputs a policy  $\pi_\xi^*$  that is optimal in the environment  $\xi$  using the discount function  $\gamma_{\text{LT}m}$ . In contrast, for infinite lifetime discounting we just consider the function  $\pi_\xi^* : (\mathcal{A} \times \mathcal{E})^* \rightarrow \mathcal{A}$ .

In order to position AINU in the arithmetical hierarchy, we need to identify these functions with sets of natural numbers. In both cases, finite and infinite lifetime, we represent these functions as relations over  $\mathbb{N} \times (\mathcal{A} \times \mathcal{E})^* \times \mathcal{A}$  and  $(\mathcal{A} \times \mathcal{E})^* \times \mathcal{A}$  respectively. These relations are easily identified with sets of natural numbers by encoding the tuple with their arguments into one natural number. From now on this translation of policies (and  $m$ ) into sets of natural numbers will be done implicitly wherever necessary.

**Lemma 6** (Policies are in  $\Delta_n^0$ ). *If a policy  $\pi$  is  $\Sigma_n^0$  or  $\Pi_n^0$ , then  $\pi$  is  $\Delta_n^0$ .*

*Proof.* Let  $\varphi$  be a  $\Sigma_n^0$ -formula ( $\Pi_n^0$ -formula) defining  $\pi$ , i.e.,  $\varphi(h, a)$  holds iff  $\pi(h) = a$ . We define the formula  $\varphi'$ ,

$$\varphi'(h, a) := \bigwedge_{a' \in \mathcal{A} \setminus \{a\}} \neg \varphi(h, a').$$

The set of actions  $\mathcal{A}$  is finite, hence  $\varphi'$  is a  $\Pi_n^0$ -formula ( $\Sigma_n^0$ -formula). Moreover,  $\varphi'$  is equivalent to  $\varphi$ .  $\square$

To compute the optimal policy, we need to compute the value function. The following lemma gives an upper bound on the computability of the value function for environments in  $\mathcal{M}$ .

**Lemma 7** (Complexity of  $V_\nu^*$ ). *For every LSCCCS  $\nu \in \mathcal{M}$ , the function  $V_\nu^*$  is  $\Pi_2^0$ -computable. For  $\gamma = \gamma_{\text{LT}m}$  the function  $V_\nu^*$  is  $\Delta_2^0$ -computable.*

*Proof.* Multiplying (3) with  $\Gamma_t \nu(e_{<t} \parallel a_{<t})$  yields  $V_\nu^*(\mathbf{x}_{<t}) > q$  if and only if

$$\lim_{m \rightarrow \infty} \bigvee_{\mathbf{x}_{t:m}} \nu(e_{1:m} \parallel a_{1:m}) R(e_{t:m}) > q \Gamma_t \nu(e_{<t} \parallel a_{<t}). \quad (5)$$

The inequality's right side is lower semicomputable, hence there is a computable function  $\psi$  such that  $\psi(\ell) \nearrow q \Gamma_t \nu(e_{<t} \parallel a_{<t}) =: q'$  for  $\ell \rightarrow \infty$ . For a fixed  $m$ , the left side is also lower semicomputable, therefore there is a computable function  $\phi$  such that  $\phi(m, k) \nearrow \bigvee_{\mathbf{x}_{t:m}} \nu(e_{1:m} \parallel a_{1:m}) R(e_{t:m}) =: f(m)$  for  $k \rightarrow \infty$ . We already know that the limit of  $f(m)$  for  $m \rightarrow \infty$  exists (uniquely), hence we

can write (5) as

$$\begin{aligned} & \lim_{m \rightarrow \infty} f(m) > q' \\ \iff & \forall m_0 \exists m \geq m_0. f(m) > q' \\ \iff & \forall m_0 \exists m \geq m_0 \exists k. \phi(m, k) > q' \\ \iff & \forall \ell \forall m_0 \exists m \geq m_0 \exists k. \phi(m, k) > \psi(\ell), \end{aligned}$$

which is a  $\Pi_2^0$ -formula. In the finite lifetime case where  $m$  is fixed, the value function  $V_\nu^*(\mathbf{x}_{<t})$  is  $\Delta_2^0$ -computable by Lemma 2 (iv), since  $V_\nu^*(\mathbf{x}_{<t}) = f(m)q/q'$ .  $\square$

From the optimal value function  $V_\nu^*$  we get the optimal policy  $\pi_\nu^*$  according to (4). However, in cases where there is more than one optimal action, we have to break an argmax tie. This happens iff  $V_\nu^*(h\alpha) = V_\nu^*(h\beta)$  for two potential actions  $\alpha \neq \beta \in \mathcal{A}$ . This equality test is more difficult than determining which is larger in cases where they are unequal. Thus we get the following upper bound.

**Theorem 8** (Complexity of Optimal Policies). *For any environment  $\nu \in \mathcal{M}$ , if  $V_\nu^*$  is  $\Delta_n^0$ -computable, then there is an optimal policy  $\pi_\nu^*$  for the environment  $\nu$  that is  $\Delta_{n+1}^0$ .*

*Proof.* To break potential ties, we pick an (arbitrary) total order  $\succ$  on  $\mathcal{A}$  that specifies which actions should be preferred in case of a tie. We define

$$\begin{aligned} \pi_\nu(h) = a \iff & \bigwedge_{a': a' \succ a} V_\nu^*(ha) > V_\nu^*(ha') \\ & \wedge \bigwedge_{a': a' \succ a'} V_\nu^*(ha) \geq V_\nu^*(ha'). \end{aligned} \quad (6)$$

Then  $\pi_\nu$  is a  $\nu$ -optimal policy according to (4). By assumption,  $V_\nu^*$  is  $\Delta_n^0$ -computable. By Lemma 2 (i) and (ii)  $V_\nu^*(ha) > V_\nu^*(ha')$  is in  $\Sigma_n^0$  and  $V_\nu^*(ha) \geq V_\nu^*(ha')$  is  $\Pi_n^0$ . Therefore the policy  $\pi_\nu$  defined in (6) is a conjunction of a  $\Sigma_n^0$ -formula and a  $\Pi_n^0$ -formula and thus in  $\Delta_{n+1}^0$ .  $\square$

**Corollary 9** (Complexity of AINU).  *$\text{AINU}_{\text{LT}}$  is  $\Delta_3^0$  and  $\text{AINU}_{\text{DC}}$  is  $\Delta_4^0$  for every environment  $\nu \in \mathcal{M}$ .*

*Proof.* From Lemma 7 and Theorem 8.  $\square$

Usually we do not mind taking slightly suboptimal actions. Therefore actually trying to determine if two actions have the exact same value seems like a waste of resources. In the following, we consider policies that attain a value that is always within some  $\varepsilon > 0$  of the optimal value.

**Definition 10** ( $\varepsilon$ -Optimal Policy). *A policy  $\pi$  is  $\varepsilon$ -optimal in environment  $\nu$  iff  $V_\nu^*(h) - V_\nu^\pi(h) < \varepsilon$  for all histories  $h \in (\mathcal{A} \times \mathcal{E})^*$ .*

**Theorem 11** (Complexity of  $\varepsilon$ -Optimal Policies). *For any environment  $\nu \in \mathcal{M}$ , if  $V_\nu^*$  is  $\Delta_n^0$ -computable, then there is an  $\varepsilon$ -optimal policy  $\pi_\nu^\varepsilon$  for the environment  $\nu$  that is  $\Delta_n^0$ .*

*Proof.* Let  $\varepsilon > 0$  be given. Since the value function  $V_\nu^*(h)$  is  $\Delta_n^0$ -computable, the set  $V_\varepsilon := \{(ha, q) \mid |q - V_\nu^*(ha)| < \varepsilon/2\}$  is in  $\Delta_n^0$  according to Definition 1. Hence we compute the values  $V_\nu^*(ha')$  until we get within  $\varepsilon/2$  for every  $a' \in \mathcal{A}$  and then choose the action with the highest value so far. Formally, let  $\succ$  be an arbitrary total order on  $\mathcal{A}$  that specifies which actions should be preferred in case of a tie. Without loss of generality, we assume  $\varepsilon = 1/k$ , and define  $Q$  to be an  $\varepsilon/2$ -grid on  $[0, 1]$ , i.e.,  $Q := \{0, 1/2k, 2/2k, \dots, 1\}$ . We define

$$\begin{aligned} \pi_\nu^\varepsilon(h) = a &: \iff \\ \exists (q_{a'})_{a' \in \mathcal{A}} \in Q^{\mathcal{A}}. & \bigwedge_{a' \in \mathcal{A}} (ha', q_{a'}) \in V_\varepsilon \\ & \wedge \bigwedge_{a': a' \succ a} q_a > q_{a'} \wedge \bigwedge_{a': a \succ a'} q_a \geq q_{a'} \\ & \wedge \text{the tuple } (q_{a'})_{a' \in \mathcal{A}} \text{ is minimal with} \\ & \text{respect to the lex. ordering on } Q^{\mathcal{A}}. \end{aligned} \quad (7)$$

This makes the choice of  $a$  unique. Moreover,  $Q^{\mathcal{A}}$  is finite since  $\mathcal{A}$  is finite, and hence (7) is a  $\Delta_n^0$ -formula.  $\square$

**Corollary 12** (Complexity of  $\varepsilon$ -Optimal AINU). *For any environment  $\nu \in \mathcal{M}$ , there is an  $\varepsilon$ -optimal policy for  $\text{AINU}_{LT}$  that is  $\Delta_2^0$  and there is an  $\varepsilon$ -optimal policy for  $\text{AINU}_{DC}$  that is  $\Delta_3^0$ .*

*Proof.* From Lemma 7 and Theorem 11.  $\square$

If the environment  $\nu \in \mathcal{M}$  is a measure, i.e.,  $\nu$  assigns zero probability to finite strings, then we get computable  $\varepsilon$ -optimal policies.

**Corollary 13** (Complexity of AIMU). *If the environment  $\mu \in \mathcal{M}$  is a measure and the discount function  $\gamma$  is computable, then  $\text{AIMU}_{LT}$  and  $\text{AIMU}_{DC}$  are limit computable ( $\Delta_2^0$ ), and  $\varepsilon$ -optimal  $\text{AIMU}_{LT}$  and  $\text{AIMU}_{DC}$  are computable ( $\Delta_1^0$ ).*

*Proof.* In the discounted case, we can truncate the limit  $m \rightarrow \infty$  in (3) at the  $\varepsilon/2$ -effective horizon  $m_{\text{eff}} := \min\{k \mid \Gamma_k/\Gamma_t < \varepsilon/2\}$ , since everything after  $m_{\text{eff}}$  can contribute at most  $\varepsilon/2$  to the value function. Any lower semicomputable measure is computable [LV08, Lem. 4.5.1]. Therefore  $V_\mu^*$  as given in (3) is composed only of computable functions, hence it is computable according to Lemma 2. The claim now follows from Theorem 8 and Theorem 11.  $\square$

### 3.3 LOWER BOUNDS

We proceed to show that the bounds from the previous section are the best we can hope for. In environment classes where ties have to be broken,  $\text{AIMU}_{DC}$  has to solve  $\Sigma_3^0$ -hard problems (Theorem 15), and  $\text{AIMU}_{LT}$  has to solve

$\Pi_2^0$ -hard problems (Theorem 16). These lower bounds are stated for particular environments  $\nu \in \mathcal{M}$ .

We also construct universal mixtures that yield bounds on  $\varepsilon$ -optimal policies. In the finite lifetime case, there is an  $\varepsilon$ -optimal  $\text{AIXI}_{LT}$  that solves  $\Sigma_1^0$ -hard problems (Theorem 17), and for general discounting, there is an  $\varepsilon$ -optimal  $\text{AIXI}_{DC}$  that solves  $\Pi_2^0$ -hard problems (Theorem 18). For arbitrary universal mixtures, we prove the following weaker statement that only guarantees incomputability.

**Theorem 14** (No AIXI is computable).  *$\text{AIXI}_{LT}$  and  $\text{AIXI}_{DC}$  are not computable for any universal Turing machine  $U$ .*

This theorem follows from the incomputability of Solomonoff induction. Since AIXI uses an analogue of Solomonoff's prior for learning, it succeeds to predict the environment's behavior for its own policy [Hut05, Thm. 5.31]. If AIXI were computable, then there would be computable environments more powerful than AIXI: they can simulate AIXI and anticipate its prediction, which leads to a contradiction.

*Proof.* Assume there is a computable policy  $\pi_\xi^*$  that is optimal in  $\xi$ . We define a deterministic environment  $\mu$ , the adversarial environment to  $\pi_\xi^*$ . The environment  $\mu$  gives rewards 0 as long as the agent follows the policy  $\pi_\xi^*$ , and rewards 1 once the agent deviates. Formally, we ignore observations by setting  $\mathcal{O} := \{0\}$ , and define

$$\mu(r_{1:t} \parallel a_{1:t}) := \begin{cases} 1 & \text{if } \forall k \leq t. a_k = \pi_\xi^*((ar)_{<k}) \text{ and } r_k = 0 \\ 1 & \text{if } \forall k \leq t. r_k = \mathbb{1}_{k \geq i} \\ & \text{where } i := \min\{j \mid a_j \neq \pi_\xi^*((ar)_{<j})\} \\ 0 & \text{otherwise.} \end{cases}$$

The environment  $\mu$  is computable because the policy  $\pi_\xi^*$  was assumed to be computable. Suppose  $\pi_\xi^*$  acts in  $\mu$ , then by [Hut05, Thm. 5.36], AIXI learns to predict perfectly on policy:

$$V_\xi^*(\mathbf{x}_{<t}) = V_\xi^{\pi_\xi^*}(\mathbf{x}_{<t}) \rightarrow V_\mu^{\pi_\xi^*}(\mathbf{x}_{<t}) = 0 \text{ as } t \rightarrow \infty,$$

since both  $\pi_\xi^*$  and  $\mu$  are deterministic. Therefore we find a  $t$  large enough such that  $V_\xi^*(\mathbf{x}_{<t}) < w_\mu$  (in the finite lifetime case we choose  $m > t$ ) where  $\mathbf{x}_{<t}$  is the interaction history of  $\pi_\xi^*$  in  $\mu$ . A policy  $\pi$  with  $\pi(\mathbf{x}_{<t}) \neq \pi_\xi^*(\mathbf{x}_{<t})$ , gets a reward of 1 in environment  $\mu$  for all time steps after  $t$ , hence  $V_\mu^\pi(\mathbf{x}_{<t}) = 1$ . With linearity of  $V_\xi^\pi(\mathbf{x}_{<t})$  in  $\xi$  [Hut05, Thm. 5.31],

$$V_\xi^\pi(\mathbf{x}_{<t}) \geq w_\mu \frac{\mu(e_{1:t} \parallel a_{1:t})}{\xi(e_{1:t} \parallel a_{1:t})} V_\mu^\pi(\mathbf{x}_{<t}) \geq w_\mu,$$

since  $\mu(e_{1:t} \parallel a_{1:t}) = 1$  ( $\mu$  is deterministic),  $V_\mu^\pi(\mathbf{x}_{<t}) = 1$ , and  $\xi(e_{1:t} \parallel a_{1:t}) \leq 1$ . Now we get a contradiction:

$$w_\mu > V_\xi^*(\mathbf{x}_{<t}) = \max_{\pi'} V_\xi^{\pi'}(\mathbf{x}_{<t}) \geq V_\xi^\pi(\mathbf{x}_{<t}) \geq w_\mu \quad \square$$

For the remainder of this section, we fix the action space to be  $\mathcal{A} := \{\alpha, \beta\}$  with action  $\alpha$  favored in ties. The percept space is fixed to a tuple of binary observations and rewards,  $\mathcal{E} := \mathcal{O} \times \{0, 1\}$  with  $\mathcal{O} := \{0, 1\}$ .

**Theorem 15** (AINU<sub>DC</sub> is  $\Sigma_3^0$ -hard). *If  $\Gamma_t > 0$  for all  $t$ , there is an environment  $\nu \in \mathcal{M}$  such that AINU<sub>DC</sub> is  $\Sigma_3^0$ -hard.*

*Proof.* Let  $A$  be any  $\Sigma_3^0$  set, then there is a computable relation  $S$  such that

$$n \in A \iff \exists i \forall t \exists k S(n, i, t, k). \quad (8)$$

We define a class of environments  $\mathcal{M}' = \{\rho_0, \rho_1, \dots\} \subseteq \mathcal{M}$  where each environment  $\rho_i$  is defined by

$$\rho_i((or)_{1:t} \parallel a_{1:t}) := \begin{cases} 2^{-t}, & \text{if } o_{1:t} = 1^t \text{ and } \forall t' \leq t. r_{t'} = 0 \\ 2^{-n-1}, & \text{if } \exists n. 1^n 0 \sqsubseteq o_{1:t} \sqsubseteq 1^n 0^\infty \text{ and } a_{n+2} = \alpha \\ & \text{and } \forall t' \leq t. r_{t'} = 0 \\ 2^{-n-1}, & \text{if } \exists n. 1^n 0 \sqsubseteq o_{1:t} \sqsubseteq 1^n 0^\infty \text{ and } a_{n+2} = \beta \\ & \text{and } \forall t' \leq t. r_{t'} = \mathbb{1}_{t' > n+1} \\ & \text{and } \forall t' \leq t \exists k S(n, i, t', k) \\ 0, & \text{otherwise.} \end{cases}$$

Every  $\rho_i$  is a chronological conditional semimeasure by definition, so  $\mathcal{M}' \subseteq \mathcal{M}$ . Furthermore, every  $\rho_i$  is lower semicomputable since  $S$  is computable.

We define our environment  $\nu$  as a mixture over  $\mathcal{M}'$ ,

$$\nu := \sum_{i \in \mathbb{N}} 2^{-i-1} \rho_i;$$

the choice of the weights on the environments  $\rho_i$  is arbitrary but positive. Let  $\pi_\nu^*$  be an optimal policy for the environment  $\nu$  and recall that the action  $\alpha$  is preferred in ties. We claim that for the  $\nu$ -optimal policy  $\pi_\nu^*$ ,

$$n \in A \iff \pi_\nu^*(1^n 0) = \beta. \quad (9)$$

This enables us to decide whether  $n \in A$  given the policy  $\pi_\nu^*$ , hence proving (9) concludes this proof.

Let  $n, i \in \mathbb{N}$  be given, and suppose we are in environment  $i$  and observe  $1^n 0$ . Taking action  $\alpha$  next yields rewards 0 forever; taking action  $\beta$  next yields a reward of 1 for those time steps  $t \geq n + 2$  for which  $\forall t' \leq t \exists k S(n, i, t', k)$ , after that the semimeasure assigns probability 0 to all next observations. Therefore, if for some  $t_0$  there is no  $k$  such that  $S(n, i, t_0, k)$ , then  $\rho_i(e_{1:t_0} \parallel \dots \beta \dots) = 0$ , and hence

$$V_{\rho_i}^*(1^n 0 \beta) = 0 = V_{\rho_i}^*(1^n 0 \alpha),$$

and otherwise  $\rho_i$  yields reward 1 for every time step after  $n + 1$ , therefore

$$V_{\rho_i}^*(1^n 0 \beta) = \Gamma_{n+2} > 0 = V_{\rho_i}^*(1^n 0 \alpha)$$

(omitting the first  $n + 1$  actions and rewards in the argument of the value function). We can now show (9): By (8),  $n \in A$  if and only if there is an  $i$  such that for all  $t$  there is a  $k$  such that  $S(n, i, t, k)$ , which happens if and only if there is an  $i \in \mathbb{N}$  such that  $V_{\rho_i}^*(1^n 0 \beta) > 0$ , which is equivalent to  $V_\nu^*(1^n 0 \beta) > 0$ , which in turn is equivalent to  $\pi_\nu^*(1^n 0) = \beta$  since  $V_\nu^*(1^n 0 \alpha) = 0$  and action  $\alpha$  is favored in ties.  $\square$

**Theorem 16** (AINU<sub>LT</sub> is  $\Pi_2^0$ -hard). *There is an environment  $\nu \in \mathcal{M}$  such that AINU<sub>LT</sub> is  $\Pi_2^0$ -hard.*

The proof of Theorem 16 is analogous to the proof of Theorem 15. The notable difference is that we replace  $\forall t' \leq t \exists k S(n, i, t', k)$  with  $\exists k S(n, i, k)$ . Moreover, we swap actions  $\alpha$  and  $\beta$ : action  $\alpha$  ‘checks’ the relation  $S$  and action  $\beta$  gives a sure reward of 1.

**Theorem 17** (Some  $\varepsilon$ -optimal AIXI<sub>LT</sub> are  $\Sigma_1^0$ -hard). *There is a universal Turing machine  $U'$  and an  $\varepsilon > 0$  such that any  $\varepsilon$ -optimal policy for AIXI<sub>LT</sub> is  $\Sigma_1^0$ -hard.*

*Proof.* Let  $\xi$  denote the universal mixture derived from the reference universal monotone Turing machine  $U$ . Let  $A$  be a  $\Sigma_1^0$ -set and  $S$  computable relation such that  $n + 1 \in A$  iff  $\exists k S(n, k)$ . We define the environment

$$\nu((or)_{1:t} \parallel a_{1:t}) := \begin{cases} \xi((or)_{1:n} \parallel a_{1:n}), & \text{if } \exists n. o_{1:n} = 1^{n-1} 0 \\ & \text{and } a_n = \alpha \\ & \text{and } \forall t' > n. e_{t'} = (0, \frac{1}{2}) \\ \xi((or)_{1:n} \parallel a_{1:n}), & \text{if } \exists n. o_{1:n} = 1^{n-1} 0 \\ & \text{and } a_n = \beta \\ & \text{and } \forall t' > n. e_{t'} = (0, 1) \\ & \text{and } \exists k S(n-1, k). \\ \xi((or)_{1:t} \parallel a_{1:t}), & \text{if } \nexists n. o_{1:n} = 1^{n-1} 0 \\ 0, & \text{otherwise.} \end{cases}$$

The environment  $\nu$  mimics the universal environment  $\xi$  until the observation history is  $1^{n-1} 0$ . Taking the action  $\alpha$  next gives rewards  $1/2$  forever. Taking the action  $\beta$  next gives rewards 1 forever if  $n \in A$ , otherwise the environment  $\nu$  ends at some future time step. Therefore we want to take action  $\beta$  if and only if  $n \in A$ . We have that  $\nu$  is an LSCCCS since  $\xi$  is an LSCCCS and  $S$  is computable.

We define the universal lower semicomputable semimeasure  $\xi' := \frac{1}{2}\nu + \frac{1}{8}\xi$ . Choose  $\varepsilon := 1/9$ . Let  $n \in A$  be given and define the lifetime  $m := n + 1$ . Let  $h \in (\mathcal{A} \times \mathcal{E})^n$  be any history with observations  $o_{1:n} = 1^{n-1} 0$ . Since  $\nu(1^{n-1} 0 \mid a_{1:n}) = \xi(1^{n-1} 0 \mid a_{1:n})$  by definition, the posterior weights of  $\nu$  and  $\xi$  in  $\xi'$  are equal to the prior weights, analogously to [LH15a, Thm. 7]. In the following, we use the linearity of  $V_\rho^{\pi_{\xi'}}^*$  in  $\rho$  [Hut05, Thm. 5.21], and the fact that values are bounded between 0 and 1. If there is a  $k$

such that  $S(n-1, k)$ ,

$$\begin{aligned} & V_{\xi'}^*(h\beta) - V_{\xi'}^*(h\alpha) \\ &= \frac{1}{2}V_{\nu}^{\pi_{\xi'}^*}(h\beta) - \frac{1}{2}V_{\nu}^{\pi_{\xi'}^*}(h\alpha) + \frac{1}{8}V_{\xi}^{\pi_{\xi'}^*}(h\beta) - \frac{1}{8}V_{\xi}^{\pi_{\xi'}^*}(h\alpha) \\ &\geq \frac{1}{2} - \frac{1}{4} + 0 - \frac{1}{8} = \frac{1}{8}, \end{aligned}$$

and similarly if there is no  $k$  such that  $S(n-1, k)$ , then

$$\begin{aligned} & V_{\xi'}^*(h\alpha) - V_{\xi'}^*(h\beta) \\ &= \frac{1}{2}V_{\nu}^{\pi_{\xi'}^*}(h\alpha) - \frac{1}{2}V_{\nu}^{\pi_{\xi'}^*}(h\beta) + \frac{1}{8}V_{\xi}^{\pi_{\xi'}^*}(h\alpha) - \frac{1}{8}V_{\xi}^{\pi_{\xi'}^*}(h\beta) \\ &\geq \frac{1}{4} - 0 + 0 - \frac{1}{8} = \frac{1}{8}. \end{aligned}$$

In both cases  $|V_{\xi'}^*(h\beta) - V_{\xi'}^*(h\alpha)| > 1/9$ . Hence we pick  $\varepsilon := 1/9$  and get for every  $\varepsilon$ -optimal policy  $\pi_{\xi'}^{\varepsilon}$ , that  $\pi_{\xi'}^{\varepsilon}(h) = \beta$  if and only if  $n \in A$ .  $\square$

**Theorem 18** (Some  $\varepsilon$ -optimal AIXI<sub>DC</sub> are  $\Pi_2^0$ -hard). *There is a universal Turing machine  $U'$  and an  $\varepsilon > 0$  such that any  $\varepsilon$ -optimal policy for AIXI<sub>DC</sub> is  $\Pi_2^0$ -hard.*

The proof of Theorem 18 is analogous to the proof of Theorem 17 except that we choose  $\forall m' \leq m \exists k S(x, m, k)$  as a condition for reward 1 after playing action  $\beta$ .

## 4 ITERATIVE VS. RECURSIVE AINU

Generally, our environment  $\nu \in \mathcal{M}$  is only a semimeasure and not a measure. I.e., there is a history  $\mathbf{x}_{<t}$  such that

$$1 > \sum_{e_t \in \mathcal{E}} \nu(e_t \mid e_{<t} \parallel a_{1:t}).$$

In such cases, with positive probability the environment  $\nu$  does not produce a new percept  $e_t$ . If this occurs, we shall use the informal interpretation that the environment  $\nu$  *ended*, but our formal argument does not rely on this interpretation.

The following proposition shows that for a semimeasure  $\nu \in \mathcal{M}$  that is not a measure, the iterative definition of AINU does not maximize  $\nu$ -expected rewards. Recall that  $\gamma_1$  states the discount of the first reward. In the following, we assume without loss of generality that  $\gamma_1 > 0$ , i.e., we are not indifferent about the reward received in time step 1.

**Proposition 19** (Iterative AINU is not a  $\nu$ -Expected Rewards Maximizer). *For any  $\varepsilon > 0$  there is an environment  $\nu \in \mathcal{M}$  that is not a measure and a policy  $\pi$  that receives a total of  $\gamma_1$  rewards in  $\nu$ , but AINU receives only  $\varepsilon\gamma_1$  rewards in  $\nu$ .*

Informally, the environment  $\nu$  is defined as follows. In the first time step, the agent chooses between the two actions  $\alpha$  and  $\beta$ . Taking action  $\alpha$  gives a reward of 1, and subsequently the environment ends. Action  $\beta$  gives a reward of  $\varepsilon$ , but the environment continues forever. There are no

other rewards in this environment. From the perspective of  $\nu$ -expected reward maximization, it is better to take action  $\alpha$ , however AINU takes action  $\beta$ .

*Proof.* Let  $\varepsilon > 0$ . We ignore observations and set  $\mathcal{E} := \{0, \varepsilon, 1\}$ ,  $\mathcal{A} := \{\alpha, \beta\}$ . The environment  $\nu$  is formally defined by

$$\nu(r_{1:t} \parallel a_{1:t}) := \begin{cases} 1 & \text{if } a_1 = \alpha \text{ and } r_1 = 1 \text{ and } t = 1 \\ 1 & \text{if } a_1 = \beta \text{ and } r_1 = \varepsilon \text{ and } r_k = 0 \forall 1 < k \leq t \\ 0 & \text{otherwise.} \end{cases}$$

Taking action  $\alpha$  first, we have  $\nu(r_{1:t} \parallel \alpha a_{2:t}) = 0$  for  $t > 1$  (the environment  $\nu$  ends in time step 2 given history  $\alpha$ ). Hence we use (3) to conclude

$$V_{\nu}^*(\alpha) = \frac{1}{\Gamma_t} \lim_{m \rightarrow \infty} \sum_{r_{1:m}} \nu(r_{1:m} \parallel \alpha a_{2:m}) \sum_{i=1}^m r_i = 0.$$

Taking action  $\beta$  first we get

$$V_{\nu}^*(\beta) = \frac{1}{\Gamma_t} \lim_{m \rightarrow \infty} \sum_{r_{1:m}} \nu(r_{1:m} \parallel \beta a_{2:m}) \sum_{i=1}^m r_i = \frac{\gamma_1}{\Gamma_1} \varepsilon.$$

Since  $\gamma_1 > 0$  and  $\varepsilon > 0$ , we have  $V_{\nu}^*(\beta) > V_{\nu}^*(\alpha)$ , and thus AINU will use a policy that plays action  $\beta$  first, receiving a total discounted reward of  $\varepsilon\gamma_1$ . In contrast, any policy  $\pi$  that takes action  $\alpha$  first receives a larger total discounted reward of  $\gamma_1$ .  $\square$

Whether it is reasonable to assume that our environment has a nonzero probability of ending is a philosophical debate we do not want to engage in here. Instead, we have a different motivation to use the recursive value function: we get an improved computability result. Concretely, we show that for all environments  $\nu \in \mathcal{M}$ , there is a limit-computable  $\varepsilon$ -optimal policy maximizing  $\nu$ -expected rewards using an infinite horizon. According to Theorem 18, this does not hold for all  $V_{\nu}^*$ -maximizing agents AINU.

In order to maximize  $\nu$ -expected rewards in case  $\nu$  is not a measure, we need the recursive definition of the value function (analogously to [Hut05, Eq. 4.12]). To avoid confusion, we denote it  $W_{\nu}^{\pi}$ :

$$\begin{aligned} W_{\nu}^{\pi}(\mathbf{x}_{<t}) &= \frac{1}{\Gamma_t} \sum_{e_t} (\gamma_t r_t \\ &\quad + \Gamma_{t+1} W_{\nu}^{\pi}(\mathbf{x}_{1:t})) \nu(e_t \mid e_{<t} \parallel a_{1:t}) \end{aligned}$$

where  $a_t := \pi(\mathbf{x}_{<t})$ . In the following we write it in non-recursive form.

**Definition 20** ( $\nu$ -Expected Value Function). The  $\nu$ -expected value of a policy  $\pi$  in an environment  $\nu$  given



history  $\mathbf{x}_{<t}$  is

$$W_\nu^\pi(\mathbf{x}_{<t}) := \frac{1}{\Gamma_t} \sum_{m=t}^{\infty} \sum_{e_{t:m}} \gamma_m r_m \nu(e_{1:m} \parallel a_{1:m})$$

if  $\Gamma_t > 0$  and  $W_\nu^\pi(\mathbf{x}_{<t}) := 0$  if  $\Gamma_t = 0$  where  $a_i := \pi(e_{<i})$  for all  $i \geq t$ . The *optimal  $\nu$ -expected value* is defined as  $W_\nu^*(h) := \sup_\pi W_\nu^\pi(h)$ .

The difference between  $V_\nu^\pi$  and  $W_\nu^\pi$  is that for  $W_\nu^\pi$  all obtained rewards matter, but for  $V_\nu^\pi$  only the rewards in timelines that continue indefinitely. In this sense the value function  $V_\nu^\pi$  conditions on surviving forever. If the environment  $\mu$  is a measure, then the history is infinite with probability one, and so  $V_\nu^\pi$  and  $W_\nu^\pi$  coincide. Hence this distinction is not relevant for AIMU, only for AINU and AIXI.

So why use  $V_\nu^\pi$  in the first place? Historically, this is how infinite-horizon AIXI has been defined [Hut05, Def. 5.30]. This definition is the natural adaptation of (optimal) minimax search in zero-sum games to the (optimal) expectimax algorithm for stochastic environments. It turns out to be problematic only because semimeasures have positive probability of ending prematurely.

**Lemma 21** (Complexity of  $W_\nu^*$ ). *For every LSCCCS  $\nu \in \mathcal{M}$ , and every lower semicomputable discount function  $\gamma$ , the function  $W_\nu^*$  is  $\Delta_2^0$ -computable.*

*Proof.* The proof is analogous to the proof of Lemma 7. We expand Definition 20 using the expectimax operator analogously to (3). This gives a quotient with numerator

$$\lim_{m \rightarrow \infty} \max_{\mathbf{x}_{t:m}} \sum_{i=t}^m \gamma_i r_i \nu(e_{1:i} \parallel a_{1:i}),$$

and denominator  $\nu(e_{<t} \parallel a_{<t}) \cdot \Gamma_t$ . In contrast to the iterative value function, the numerator is now nondecreasing in  $m$  because we assumed rewards to be nonnegative (Assumption 3b). Hence both numerator and denominator are lower semicomputable functions, so Lemma 2 (iv) implies that  $W_\nu^*$  is  $\Delta_2^0$ -computable.  $\square$

Now we can apply our results from Section 3.2 to show that using the recursive value function  $W_\nu^\pi$ , we get a universal AI model with an infinite horizon whose  $\varepsilon$ -approximation is limit computable. Moreover, in contrast to iterative AINU, recursive AINU actually maximizes  $\nu$ -expected rewards.

**Corollary 22** (Complexity of Recursive AINU/AIXI). *For any environment  $\nu \in \mathcal{M}$ , recursive AINU is  $\Delta_3^0$  and there is an  $\varepsilon$ -optimal recursive AINU that is  $\Delta_2^0$ . In particular, for any universal Turing machine, recursive AIXI is  $\Delta_3^0$  and there is an  $\varepsilon$ -optimal recursive AIXI that is limit computable.*

*Proof.* From Theorem 8, Theorem 11, and Lemma 21.  $\square$

Analogously to Theorem 14, Theorem 16, and Theorem 17 we can show that recursive AIXI is not computable, recursive AINU is  $\Pi_2^0$ -hard, and for some universal Turing machines,  $\varepsilon$ -optimal recursive AIXI is  $\Sigma_1^0$ -hard.

## 5 DISCUSSION

We set out with the goal of finding a limit-computable perfect agent. Table 3 on page 4 summarizes our computability results regarding Solomonoff's prior  $M$ : conditional  $M$  and  $M_{\text{norm}}$  are limit computable, while  $\overline{M}$  and  $\overline{M}_{\text{norm}}$  are not. Table 1 on page 2 summarizes our computability results for AINU, AIXI, and AINU: iterative AINU with finite lifetime is  $\Delta_3^0$ . Having an infinite horizon increases the level by one, while restricting to  $\varepsilon$ -optimal policies decreases the level by one. All versions of AINU are situated between  $\Delta_2^0$  and  $\Delta_4^0$  (Corollary 9 and Corollary 12). For environments that almost surely continue forever (semimeasure that are measures), AIMU is limit-computable and  $\varepsilon$ -optimal AIMU is computable. We proved that these computability bounds on iterative AINU are generally unimprovable (Theorem 15 and Theorem 16). Additionally, we proved weaker lower bounds for AIXI independent of the universal Turing machine (Theorem 14) and for  $\varepsilon$ -optimal AIXI for specific choices of the universal Turing machine (Theorem 17 and Theorem 18).

We considered  $\varepsilon$ -optimality in order to avoid having to break argmax ties. This  $\varepsilon$  does not have to be constant over time, instead we may let  $\varepsilon \rightarrow 0$  as  $t \rightarrow \infty$  at any computable rate. With this we retain the computability results of  $\varepsilon$ -optimal policies and get that the value of the  $\varepsilon(t)$ -optimal policy  $\pi_\nu^{\varepsilon(t)}$  converges rapidly to the  $\nu$ -optimal value:  $V_\nu^*(\mathbf{x}_{<t}) - V_\nu^{\pi_\nu^{\varepsilon(t)}}(\mathbf{x}_{<t}) \rightarrow 0$  as  $t \rightarrow \infty$ . Therefore the limitation to  $\varepsilon$ -optimal policies is not very restrictive.

When the environment  $\nu$  has nonzero probability of not producing a new percept, the iterative definition (Definition 4) of AINU fails to maximize  $\nu$ -expected rewards (Proposition 19). We introduced a recursive definition of the value function for infinite horizons (Definition 20), which correctly returns  $\nu$ -expected value. The difference between the iterative value function  $V$  and recursive value function  $W$  is readily exposed in the difference between  $M$  and  $\overline{M}$ . Just like  $V$  conditions on surviving forever, so does  $\overline{M}$  eliminate the weight of programs that do not produce infinite strings. Both  $\overline{M}$  and  $V$  are not limit computable for this reason.

Our main motivation for the introduction of the recursive value function  $W$  is the improvement of the computability of optimal policies. Recursive AINU is  $\Delta_3^0$  and admits a limit-computable  $\varepsilon$ -optimal policy (Corollary 22). In this sense our goal to find a limit-computable perfect agent has been accomplished.

## REFERENCES

- [BD62] David Blackwell and Lester Dubins. Merging of opinions with increasing information. *The Annals of Mathematical Statistics*, pages 882–886, 1962.
- [Gá83] Péter Gács. On the relation between descriptive complexity and algorithmic probability. *Theoretical Computer Science*, 22(1–2):71–93, 1983.
- [Hut00] Marcus Hutter. A theory of universal artificial intelligence based on algorithmic complexity. Technical Report cs.AI/0004001, 2000. <http://arxiv.org/abs/cs.AI/0004001>.
- [Hut01] Marcus Hutter. New error bounds for Solomonoff prediction. *Journal of Computer and System Sciences*, 62(4):653–667, 2001.
- [Hut05] Marcus Hutter. *Universal Artificial Intelligence: Sequential Decisions Based on Algorithmic Probability*. Springer, 2005.
- [LH14] Tor Lattimore and Marcus Hutter. General time consistent discounting. *Theoretical Computer Science*, 519:140–154, 2014.
- [LH15a] Jan Leike and Marcus Hutter. Bad universal priors and notions of optimality. In *Conference on Learning Theory*, 2015.
- [LH15b] Jan Leike and Marcus Hutter. On the computability of Solomonoff induction and knowledge-seeking. 2015. Forthcoming.
- [LV08] Ming Li and Paul M. B. Vitányi. *An Introduction to Kolmogorov Complexity and Its Applications*. Texts in Computer Science. Springer, 3rd edition, 2008.
- [MGLA00] Martin Mundhenk, Judy Goldsmith, Christopher Lusena, and Eric Allender. Complexity of finite-horizon Markov decision process problems. *Journal of the ACM*, 47(4):681–720, 2000.
- [MHC99] Omid Madani, Steve Hanks, and Anne Condon. On the undecidability of probabilistic planning and infinite-horizon partially observable Markov decision problems. In *AAAI/IAAI*, pages 541–548, 1999.
- [MHC03] Omid Madani, Steve Hanks, and Anne Condon. On the undecidability of probabilistic planning and related stochastic optimization problems. *Artificial Intelligence*, 147(1):5–34, 2003.
- [Nie09] André Nies. *Computability and Randomness*. Oxford University Press, 2009.
- [PT87] Christos H Papadimitriou and John N Tsitsiklis. The complexity of Markov decision processes. *Mathematics of Operations Research*, 12(3):441–450, 1987.
- [RH11] Samuel Rathmanner and Marcus Hutter. A philosophical treatise of universal induction. *Entropy*, 13(6):1076–1136, 2011.
- [SB98] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, 1998.
- [SLR07] Régis Sabbadin, Jérôme Lang, and Nasolo Ravoanjanahry. Purely epistemic Markov decision processes. In *AAAI*, volume 22, pages 1057–1062, 2007.
- [Sol64] Ray Solomonoff. A formal theory of inductive inference. Parts 1 and 2. *Information and Control*, 7(1):1–22 and 224–254, 1964.
- [Sol78] Ray Solomonoff. Complexity-based induction systems: Comparisons and convergence theorems. *IEEE Transactions on Information Theory*, 24(4):422–432, 1978.
- [VNH<sup>+</sup>11] Joel Veness, Kee Siong Ng, Marcus Hutter, William Uther, and David Silver. A Monte-Carlo AIXI approximation. *Journal of Artificial Intelligence Research*, 40(1):95–142, 2011.
- [WSH11] Ian Wood, Peter Sunehag, and Marcus Hutter. (Non-)equivalence of universal priors. In *Solomonoff 85th Memorial Conference*, pages 417–425. Springer, 2011.