

# Universal Reinforcement Learning Algorithms: Survey and Experiments

John Aslanides<sup>†</sup>, Jan Leike<sup>‡</sup>, Marcus Hutter<sup>†</sup>  
{john.aslanides,marcus.hutter}@anu.edu.au, leike@google.com

<sup>†</sup> Australian National University  
<sup>‡</sup> Future of Humanity Institute

IJCAI 2017

August 24, 2017

# Overview

- Introduction
- Algorithms
- Experiments

# Motivation

- **Universal RL (URL):** *making very weak assumptions about its environment, what can an agent achieve, in principle?*
  - What is intelligent behavior?
  - What is a useful optimality criterion?
- Theoretically studied, but few to no experiments/reference implementations to date

# Motivation

- **Universal RL (URL):** *making very weak assumptions about its environment, what can an agent achieve, in principle?*
  - What is intelligent behavior?
  - What is a useful optimality criterion?
- Theoretically studied, but few to no experiments/reference implementations to date
- Contribution: experiments, along with open-source demo platform for several URL algorithms.

## Demo

Online demo: <http://aslanides.io/aixijs>

**AIXI.js**  
About Demos

**MC-AIXI**  
Monte Carlo AIXI on a known Gridworld.

**MC-AIXI-Dirichlet**  
AIXI with a Dirichlet model on an unknown Gridworld.

**Thompson Sampling**  
Thompson sampling on a known Gridworld.

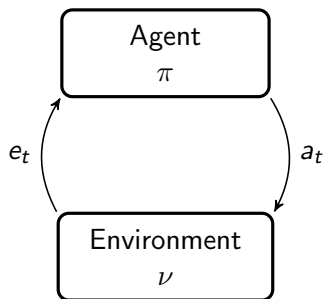
**Hooked on noise**  
Entropy-seeking agents get hooked on white noise and stop exploring, while the knowledge-seeking agent ignores it.

**Knowledge-seeking agents**  
Compare the behavior of the Square, Shannon, and Kullback-Leibler knowledge-seeking agents.

For use on GitHub

# Agent-environment model

- Environment class: **POMDPs** (possibly non-ergodic)
- Percepts ( $\neq$  states) are (**observation, reward**) pairs  $e_k = (o_k, r_k)$
- Interact to generate a **history**  $h_t := a_1 e_1 a_2 e_2 \dots a_t e_t$ .



AI $\xi$  (Hutter, 2005)

- Non-parametric **Bayesian mixture** over some countable **model class**  $\mathcal{M}$ :

$$\xi(e_t|h_t) \doteq \sum_{\nu \in \mathcal{M}} w(\nu|h_t) \nu(e_t|h_t)$$

# AI $\xi$ (Hutter, 2005)

- Non-parametric **Bayesian mixture** over some countable **model class**  $\mathcal{M}$ :

$$\xi(e_t|h_t) \doteq \sum_{\nu \in \mathcal{M}} w(\nu|h_t) \nu(e_t|h_t)$$

- The **Bayes-optimal policy** (AI $\xi$ ; Hutter, 2005) is

$$a_t = \arg \max_{a_t} \underbrace{\sum_{e_t} \cdots \max_{a_m} \sum_{e_m}}_{\text{expectimax search}} \underbrace{\sum_{k=t}^m \gamma_k u(h_k)}_{\text{return}} \underbrace{\prod_{j=t}^k \xi(e_j|h_t)}_{\text{environment model}} .$$



# AI $\xi$ (Hutter, 2005)

- Non-parametric **Bayesian mixture** over some countable **model class**  $\mathcal{M}$ :

$$\xi(e_t|h_t) \doteq \sum_{\nu \in \mathcal{M}} w(\nu|h_t) \nu(e_t|h_t)$$

- The **Bayes-optimal policy** (AI $\xi$ ; Hutter, 2005) is

$$a_t = \arg \max_{a_t} \underbrace{\sum_{e_t} \cdots \max_{a_m} \sum_{e_m}}_{\text{expectimax search}} \underbrace{\sum_{k=t}^m \gamma_k u(h_k)}_{\text{return}} \underbrace{\prod_{j=t}^k \xi(e_j|h_t)}_{\text{environment model}} .$$

- In practice:
  - Forward planning by MCTS
  - Use manageable model class  $\mathcal{M}$

# AI $\xi$ (Hutter, 2005)

- Non-parametric **Bayesian mixture** over some countable **model class**  $\mathcal{M}$ :

$$\xi(e_t|h_t) \doteq \sum_{\nu \in \mathcal{M}} w(\nu|h_t) \nu(e_t|h_t)$$

- The **Bayes-optimal policy** (AI $\xi$ ; Hutter, 2005) is

$$a_t = \arg \max_{a_t} \underbrace{\sum_{e_t} \cdots \max_{a_m} \sum_{e_m}}_{\text{expectimax search}} \underbrace{\sum_{k=t}^m \gamma_k u(h_k)}_{\text{return}} \underbrace{\prod_{j=t}^k \xi(e_j|h_t)}_{\text{environment model}} .$$

- In practice:
  - Forward planning by MCTS
  - Use manageable model class  $\mathcal{M}$

- [Video]

# Issues

Problems:

- Not asymptotically optimal (Orseau, 2010)

# Issues

## Problems:

- Not asymptotically optimal (Orseau, 2010)
- Won't overcome the bias of bad priors, *c.f.* supervised learning (Leike & Hutter, 2015)

## Knowledge-seeking agents (Orseau, 2013)

- **Utility agents** are intrinsically motivated, and don't need an extrinsic reward signal.
- Knowledge-seeking agent (KSA) – motivated to reduce uncertainty
- No exploration/exploitation tradeoff

Agent	Utility function	Description
$AI\xi$	$r$	Reward
Square-KSA	$-\xi$	Entropy
Shannon-KSA	$-\log[\xi]$	Entropy
Kullback-Leibler-KSA	$-\Delta\text{Ent}[w(\cdot)]$	Information gain

## Knowledge-seeking agents (Orseau, 2013)

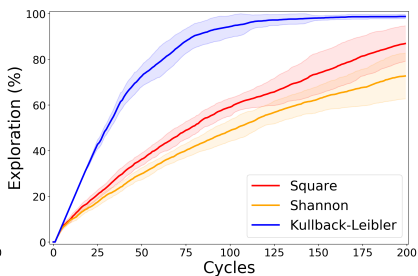
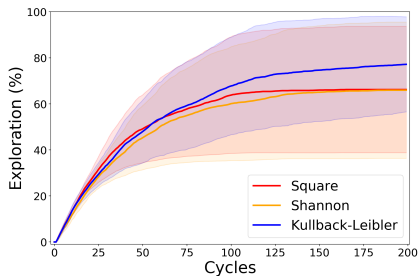
- **Utility agents** are intrinsically motivated, and don't need an extrinsic reward signal.
- Knowledge-seeking agent (KSA) – motivated to reduce uncertainty
- No exploration/exploitation tradeoff

Agent	Utility function	Description
$AI\xi$	$r$	Reward
Square-KSA	$-\xi$	Entropy
Shannon-KSA	$-\log[\xi]$	Entropy
Kullback-Leibler-KSA	$-\Delta\text{Ent}[w(\cdot)]$	Information gain

- [Video]

# Experiments

- Qualitative behavior is highly model-sensitive
- KL-KSA outperforms entropy-seeking in stochastic environments



# Outlook

- Open-source online JavaScript demo: <https://aslanides.io/aixijs>
- Used to run experiments for another IJCAI paper (Reinforcement Learning with a Corrupted Reward Channel, Everitt et al. 2017)
- Come and talk to me at the ANU booth downstairs :)