Advances in Universal Artificial Intelligence

Marcus Hutter

Australian National University Canberra, ACT, 0200, Australia http://www.hutter1.net/



Abstract

There is great interest in understanding and constructing generally intelligent systems approaching and ultimately exceeding human intelligence. Universal AI is such a mathematical theory of machine super-intelligence. More precisely, AIXI is an elegant parameter-free theory of an optimal reinforcement learning agent embedded in an arbitrary unknown environment that possesses essentially all aspects of rational intelligence. The theory reduces all conceptual AI problems to pure computational questions. After a brief discussion of its philosophical, mathematical, and computational ingredients, I will give a formal definition and measure of intelligence, which is maximized by AIXI. AIXI can be viewed as the most powerful Bayes-optimal sequential decision maker, for which I will present general optimality results. This also motivates some variations such as knowledge-seeking and optimistic agents, and feature reinforcement learning. Finally I present some recent approximations, implementations, and applications of this modern top-down approach to AI.

Overview

Goal: Construct a single universal agent that learns to act optimally in any environment.

- 3 -

State of the art: Formal (mathematical, non-comp.) definition of such an agent.

Accomplishment: Well-defines AI. Formalizes rational intelligence. Formal "solution" of the AI problem in the sense of ...

⇒ Reduces the conceptional AI problem to a (pure) computational problem.

Evidence: Mathematical optimality proofs and some experimental results.



- Universal Intelligence
- General Bayesian Agents
- Variations of Universal/Bayesian Agents
- Approximations & Applications
- Discussion

UNIVERSAL INTELLIGENCE

Marcus Hutter



Stuart Russell • Peter Norvig

Agent Model with Reward

Most if not all AI problems can be formulated within the agent framework! But how choose Agent?





- 6 -

Foundations of Universal Artificial Intelligence



Ockhams' razor (simplicity) principle

Entities should not be multiplied beyond necessity.

 $\mathsf{Posterior}(H|D) \propto \mathsf{Likelihood}(D|H) \times \mathsf{Prior}(H).$

Epicurus' principle of multiple explanations

Bayes' rule for conditional probabilities

If more than one theory is consistent with the observations, keep all theories.

Given the prior belief/probability one can predict all future probabilities.



Bellman equations

Theory of how to optimally plan and act in known environments. Solomonoff + Bellman = Universal Artificial Intelligence.

Turing's universal machine

Everything computable by a human using a fixed procedure can also be computed by a (universal) Turing machine.

Kolmogorov's complexity

The complexity or information content of an object is the length of its shortest description on a universal Turing machine.

Solomonoff's universal prior=Ockham+Epicurus+Bayes+Turing

Solves the question of how to choose the prior if nothing is known. \Rightarrow universal induction, formal Ockham. Prior $(H) = 2^{-\text{Kolmogorov}(H)}$

• $(\mathcal{A}, \mathcal{O}, \mathcal{R}) = (action, observation, reward) spaces.$ a_k =action at time k; x_k := $o_k r_k$ =perception at time k.

- 8 -

- Agent follows policy $\pi : (\mathcal{A} \times \mathcal{O} \times \mathcal{R})^* \rightsquigarrow \mathcal{A}$
- Environment reacts with $\mu : (\mathcal{A} \times \mathcal{O} \times \mathcal{R})^* \times \mathcal{A} \rightsquigarrow \mathcal{O} \times \mathcal{R}$
- Performance of agent π in environment μ = expected cumulative reward = $V_{\mu}^{\pi} := \mathbb{E}_{\mu}^{\pi} [\sum_{t=1}^{\infty} r_{t}^{\pi\mu}]$
- There are various ways to regularize the infinite reward sum: finite horizon, discounting, summability condition on μ .
- μ -optimal policy Al μ : $p^{\mu} := \arg \max_{\pi} V^{\pi}_{\mu}$

Formal Definition of Intelligence

Usually true environment µ unknown
 ⇒ average over wide range of environments
 (all semi-computable chronological semi-measures M_U)

- 9 -

- Ockham+Epicurus: Weigh each environment with its
 Kolmogorov complexity K(μ) := min_p {length(p) : U(p) = μ}
- Universal intelligence of agent π is $\Upsilon(\pi) := \sum_{\mu \in \mathcal{M}_U} 2^{-K(\mu)} V_{\mu}^{\pi}$.
- Informal interpretation: Intelligence measures an agent's ability to perform well in a wide range of environments.
- Properties of Υ : valid, informative, wide range, general, dynamic, unbiased, fundamental, formal, objective, fully defined, universal.
- $AIXI = \arg \max_{\pi} \Upsilon(\pi) = most intelligent agent.$

Marcus Hutter

Explicit AIXI Model in one Line

complete & essentially unique & limit-computable

AIXI:
$$a_k := \arg \max_{a_k} \sum_{o_k r_k} \dots \max_{a_m} \sum_{o_m r_m} [r_k + \dots + r_m] \sum_{p : U(p, a_1 \dots a_m) = o_1 r_1 \dots o_m r_m} 2^{-length(p)}$$

k=now, action, observation, reward, Universal TM, program, m=lifespan

AIXI is an elegant mathematical theory of general AI,

but incomputable, so needs to be approximated in practice.

- 10 -

Claim: AIXI is the most intelligent environmental independent, i.e. universally optimal, agent possible.

Proof: For formalizations, quantifications, and proofs, see [Hut05].

Potential Applications: Intelligent Agents, Games, Optimization, Active Learning, Adaptive Control, Robots, Philosophy of Mind, AI safety.

Kolmogorov complexity:

- generalization
- associative learning
- transfer learning [Mah09]
- knowledge representation
- abstraction
- similarity [CV05]
- regularization, bias-variance [Wal05]

Bayes:

- exploration-exploitation
- learning

History-based:

- partial observability
- non-stationarity
- long-term memory
- large state space

Expectimax:

• planning

UAI deals with these issues in a general and optimal way

Particularly Interesting Problems

- Sequence Prediction, e.g. weather or stock-market prediction. Strong result: $V^*_{\mu} - V^{p^{\xi}}_{\mu} = O(\sqrt{\frac{K(\mu)}{m}})$, m =horizon.
- Strategic Games: Learn to play well (minimax) strategic zero-sum games (like chess) or even exploit limited capabilities of opponent.
- Optimization: Find (approximate) minimum of function with as few function calls as possible. Difficult exploration versus exploitation problem.
- Supervised learning: Learn functions by presenting (z, f(z)) pairs and ask for function values of z' by presenting (z',?) pairs.
 Supervised learning is much faster than reinforcement learning.

AIXI quickly learns to predict, play games, optimize, and learn supervised

Marcus Hutter - 13 - Universal Artificial Intelligence

Curious/Philosophical/Social Questions for AIXI

- Where do rewards come from if humans are not around (see later, knowledge-seeing) agents [Ors11, OLH13]
- Will AIXI take drugs (wire-heading, hack reward system) [OR11]
- Will AIXI commit suicide [MEH16]
- Curiosity killed the cat and maybe AIXI [Sch07, Ors11, LHS13]
- Immortality can cause laziness [Hut05, Sec.5.7]
- Can self-preservation be learned or need parts be innate [RO11]

GENERAL BAYESIAN AGENTS

Agents in Probabilistic Environments

- Given history $a_{1:k}x_{< k}$, the probability that the environment leads to perception x_k in cycle k is (by definition) $\sigma(x_k|a_{1:k}x_{< k})$.
- Abbr.: $\sigma(x_{1:m}|a_{1:m}) = \sigma(x_1|a_1) \cdot \sigma(x_2|a_{1:2}x_1) \cdot \ldots \cdot \sigma(x_m|a_{1:m}x_{< m})$
- Value of policy p in environment σ is defined as expected discounted future reward sum: $V_{k\gamma}^{p\sigma} := \frac{1}{\Gamma_k} \lim_{m \to \infty} \sum_{x_{k:m}} (\gamma_k r_k + ... + \gamma_m r_m) \sigma(x_{k:m} | a_{1:m} x_{< k})_{|a_{1:m} = p(x_{< m})}$
- General discount sequence $\gamma_1, \gamma_2, \gamma_3, \dots$ Normalizer $\Gamma_k := \sum_{i=k}^{\infty} \gamma_i$
- The goal of the agent should be to maximize the value.
- σ -optimal policy Al σ : $p^{\sigma} := \arg \max_{p} V_{k\gamma}^{p\sigma}$
- If true env. μ is known, choose $\sigma = \mu$.

The Bayes-Mixture Distribution ξ

Assumption: The true environment μ is unknown.

Bayesian approach: The true probability distribution μ is not learned directly, but is replaced by a Bayes-mixture ξ .

Assumption: We know that the true environment μ is contained in some known (finite or countable) set \mathcal{M} of environments.

The Bayes-mixture ξ is defined as

 $\xi(x_{1:m}|a_{1:m}) := \sum_{\nu \in \mathcal{M}} w_{\nu} \nu(x_{1:m}|a_{1:m}) \quad \text{with} \quad \sum_{\nu \in \mathcal{M}} w_{\nu} = 1, \quad w_{\nu} > 0 \; \forall \nu$ The weights w_{ν} may be interpreted as the prior degree of belief that the true environment is ν .

Then $\xi(x_{1:m}|a_{1:m})$ could be interpreted as the prior subjective belief probability in observing $x_{1:m}$, given actions $a_{1:m}$.

Questions of Interest

- It is natural to follow the policy p^{ξ} which maximizes V_{ξ}^{p} .
- If μ is the true environment the expected reward when following policy p^{ξ} will be $V_{\mu}^{p^{\xi}}.$
- The optimal (but infeasible) policy p^{μ} yields reward $V^{p^{\mu}}_{\mu} \equiv V^{*}_{\mu}$.
- Are there policies with uniformly larger value than $V^{p^{\xi}}_{\mu}$?
- How close is $V^{p^{\xi}}_{\mu}$ to V^{*}_{μ} ?

Marcus Hutter

• What is the most general class $\mathcal M$ and weights w_{ν} ?

 $\mathcal{M} = \mathcal{M}_U$ and $w_\nu = 2^{-K(\nu)} \implies \mathsf{AI}\xi = \mathsf{AIXI}$!

Universal Artificial Intelligence

Convergence of ξ to μ

Countable mixture \Rightarrow dominance $\xi(x_{1:n}|a_{1:n}) \ge w_{\mu}\mu(x_{1:n}|a_{1:n}) \Rightarrow$

Theorem 1 (multistep predictive ξ converges to μ) $\xi(x_{k:m_k}|x_{\langle k}a_{1:m_k}) \xrightarrow{k \to \infty} \mu(x_{k:m_k}|x_{\langle k}a_{1:m_k})$ with μ prob. 1. with rapid conv. for bounded horizon $h_k \equiv m_k - k + 1 \leq h_{max} < \infty$

- Caveat: Convergence holds only for actions actually chosen.
- Does replacing μ with ξ lead to Al ξ system with asymptotically optimal behavior with rapid convergence?
- This looks promising from analogy to Sequence Prediction but is much more subtle and tricky!

Marcus Hutter - 19 - Universal Artificial Intelligence

Convergence of Universal to True Value

Theorem 2 (Convergence of universal to true value)

For a given policy p and history generated by p and μ , i.e. on-policy, the future universal value $V_{\dots}^{p\xi}$ converges to the true value $V_{\dots}^{p\mu}$:

$$V_{k\gamma}^{p\xi} \xrightarrow{k \to \infty} V_{k\gamma}^{p\mu}$$
 i.m. for any γ .

If the history is generated by $p = p^{\xi}$, this implies $V_{k\gamma}^{*\xi} \to V_{k\gamma}^{p^{\xi}\mu}$.

Hence the universal value $V_{k\gamma}^{*\xi}$ can be used to estimate the true value $V_{k\gamma}^{p^{\xi}\mu}$, without any assumptions on \mathcal{M} and γ .

Nevertheless, maximization of $V_{k\gamma}^{p\xi}$ may asymptotically differ from max. of $V_{k\gamma}^{p\mu}$, since $V_{k\gamma}^{p\xi} \not\rightarrow V_{k\gamma}^{p\mu}$ for $p \neq p^{\xi}$ is possible (and also $V_{k\gamma}^{*\xi} \not\rightarrow V_{k\gamma}^{*\mu}$).

Results for Discounted Future Value

Theorem 3 (Properties of Discounted Future Value)

- $V_{k\gamma}^{\pi\rho}$ is linear in ρ : $V_{k\gamma}^{\pi\xi} = \sum_{\nu} w_{k-1}^{\nu} V_{k\gamma}^{\pi\nu}$.
- $V_{k\gamma}^{*\rho}$ is convex in ρ : $V_{k\gamma}^{*\xi} \leq \sum_{\nu} w_{k-1}^{\nu} V_{k\gamma}^{*\nu}$.
- where $w_{k-1}^{\nu} := w_{\nu} \frac{\nu(x_{\leq k}|a_{\leq k})}{\xi(x_{\leq k}|a_{\leq k})}$ is the posterior belief in ν .
- p^{ξ} is Pareto-optimal in the sense that there is no other policy π with $V_{k\gamma}^{\pi\nu} \geq V_{k\gamma}^{p^{\xi}\nu}$ for all $\nu \in \mathcal{M}$ and strict inequality for at least one ν .
- If there exists a self-optimizing policy for \mathcal{M} , then p^{ξ} is selfoptimizing in the sense that If $\exists \tilde{\pi}_k \forall \nu : V_{k\gamma}^{\tilde{\pi}_k \nu} \stackrel{k \to \infty}{\longrightarrow} V_{k\gamma}^{*\nu} \implies V_{k\gamma}^{p^{\xi} \mu} \stackrel{k \to \infty}{\longrightarrow} V_{k\gamma}^{*\mu}$.



Importance of the Right Discounting

Standard geometric discounting: $\gamma_k = \gamma^k$ with $0 < \gamma < 1$.

Problem: Most environments do not possess self-optimizing policies under this discounting.

Reason: Effective horizon h_k^{eff} is finite $(\sim \ln \frac{1}{\gamma} \text{ for } \gamma_k = \gamma^k)$.

The analogue of $m \to \infty$ is $k \to \infty$ and $h_k^{e\!f\!f} \to \infty$ for $k \to \infty$.

Result: Policy p^{ξ} is self-optimizing for the class of $(l^{th} \text{ order})$ ergodic MDPs if $\frac{\gamma_{k+1}}{\gamma_k} \to 1$.

Example discounting: $\gamma_k = k^{-2}$ or $\gamma_k = k^{-1-\varepsilon}$ or $\gamma_k = 2^{-K(k)}$.

Horizon is of the order of the age of the agent: $h_k^{eff} \sim k$.

[Ors10]

[LH11]

Is Bayesian RL Optimal?

- asymptotically optimal := self-optimizing on policy-induced history
- AIXI is *not* asymptotically optimal
- No policy can be asymptotically optimal for \mathcal{M}_U
- There are finite \mathcal{M} for which the Bayes-optimal policy p^{ξ} is *not* asymptotically optimal (for any γ).
- For every (in)finite \mathcal{M} there exist [LH14a, Lat14, LLOH16] (weakly/mean) asymptotically optimal policies (see below)
- Jumping into a trap is asymptotically optimal. It also has great PAC bound.
- Bayesian RL may still be (regarded as) "best" (by construction, its Pareto-optimality, Thompson sampling variation, ...)
- ⇒ further theoretical investigations of Bayesian RL and alternatives are needed.

VARIATIONS OF UNIVERSAL/BAYESIAN AGENTS

- Knowledge-Seeking Agents
- Exploration Bursts
- Optimistic Agents
- Thompson Sampling

Origin of Rewards and Universal Goals

- Where do rewards come from if we don't (want to) provide them?
- Human interaction: reward the robot according to how well it solves the tasks we want it to do.
- Autonomous: Hard-wire reward to predefined task:
 E.g. Mars robot: reward = battery level & evidence of water/life.
- Is there something like a universal goal?
- Curiosity-driven learning [Sch07]
- Knowledge seeking agents

[Ors11, OLH13]

Universal Knowledge-Seeking Agent (KSA) reward for exploration; goal is to learn the true environment [OLH13]

- $w_k^{\nu} := w_{\nu} \frac{\nu(x_{1:k}|a_{1:k})}{\xi(x_{1:k}|a_{1:k})}$ is the posterior belief in ν given history $\alpha x_{1:k}$.
- $w_k^{()}$ summarizes the information contained in history $\alpha x_{1:k}$.
- $w_{k-1}^{()} \rightsquigarrow w_k^{()}$ changes $\Leftrightarrow x_k$ given $a\!x_{< k}$ is informative about $\nu \in \mathcal{M}$
- Information gain can be quantified by KL-divergence.
- Reward agent for gained information: $r_k := \mathsf{KL}(w_k^{()}||w_{k-1}^{()}) \equiv \sum_{\nu \in \mathcal{M}} w_k^{\nu} \log(w_k^{\nu}/w_{k-1}^{\nu})$

Asymptotic Optimality of Universal KSA

Theorem 4 (Asymptotic Optimality of Universal KSA)

- Universal π_{ξ}^* converges to optimal π_{μ}^* . More formally:
- $P_{\xi}^{\pi}(\cdot|ax_{< k})$ converges in (μ, π_{ξ}^{*}) -probability to $P_{\mu}^{\pi}(\cdot|ax_{< k})$ uniformly for all π .

Def: $P_{\rho}^{\pi}(\cdot | ax_{< k})$ is (ρ, π) -probability of future $ax_{k:\infty}$ given past $ax_{< k}$. Note: On-policy agent π_{ξ}^{*} is able to even predict off-policy! Remark: No assumption on \mathcal{M} needed, i.e. Thm. applicable to \mathcal{M}_{U} .

Bayesian RL with Extra Exploration Bursts

Combining Bayes-optimal and KSA policies we can achieve PAC bounds and weak asymptotic optimality in arbitrary environment classes \mathcal{M} .

BayesExp algorithm (Basic Idea)

If the Bayes-expected info-gain (see KSA) is small, then "exploit" by following the Bayes optimal policy for 1 step else explore by following a policy that maximises the expected information gain for a couple of time-steps.

Results:

- Optimal minimax sample-complexity (PAC) bounds in arbitrary finite class *M* of history-based environments. [LH14a]
- Weak asymptotic optimality in arbitrary countable class of history-based environments, including \mathcal{M}_U . [Lat14]
- Inq algorithm: Similar and even strong asymptotic optimal [CCH19]

Bayesian RL with Thompson Sampling

Thompson Sampling (TS) algorithm (Basic Idea)

- sample environment $u \in \mathcal{M}$ from posterior probability $w_k^
 u$,
- follow ν -optimal policy π^*_{ν} for a couple of time-steps. Important: Resample only after an effective horizon!

(Cf. Bayes-optimal policy maximizes the Bayesian mixture value, which is the posterior average over the values of all environments in \mathcal{M} .)

Results:

[LLOH16]

- Mean asymptotic optimality in arbitrary countable class of history-based environments, including \mathcal{M}_U .
- Given a recoverability assumption, also regret is sublinear.

Remarks: TS is more natural than Bayes with Exploration Bursts. Thompson Sampling is a stochastic policy unlike Bayes-optimal policies.

Optimistic Agents in Deterministic Worlds

act optimally w.r.t. the most optimistic environment until it is contradicted [SH12]

- $\pi^{\circ} := \pi_k^* := \arg \max_{\pi} \max_{\nu \in \mathcal{M}_{k-1}} V_{k\gamma}^{\pi\nu}(\alpha x_{< k})$
- \mathcal{M}_{k-1} := environments consistent with history $\alpha x_{< k}$.
- As long as the outcome is consistent with the optimistic prediction, the return is optimal, even if the wrong environment is chosen.

Theorem 5 (Optimism is asymptotically optimal)

For finite $\mathcal{M} \equiv \mathcal{M}_0$, where $\mu \in \mathcal{M}$ is the true environment

- Asymptotic: $V_{k\gamma}^{\pi^{\circ}\mu} = V_{k\gamma}^{*\mu}$ for all large k.
- Errors: For geometric discount, $V_{k\gamma}^{\pi^{\circ}\mu} \ge V_{k\gamma}^{*\mu} \varepsilon$ (i.e. $\pi^{\circ} \varepsilon$ -suboptimal) for all but at most $|\mathcal{M}| \frac{\log \varepsilon (1-\gamma)}{\log \gamma}$ time steps k.

Optimistic Agents for General Environments

- Generalization to stochastic environments: Likelihood criterion: Exclude ν from \mathcal{M}_{k-1} if $\nu(x_{< k}|a_{< k}) < \varepsilon_k \cdot \max_{\nu \in \mathcal{M}} \nu(x_{< k}|a_{< k})$. [SH12]
- Generalization to compact classes \mathcal{M} : Replace \mathcal{M} by centers of finite ε -cover of \mathcal{M} in def. of π° . [SH12]
- Use decreasing $\varepsilon_k \to 0$ to get asymptotic optimality.
- There are non-compact classes for which asymptotic optimality is impossible to achieve. [Ors10]
- Weaker asymptotic optimality in Cesaro sense possible by starting with finite subset $\mathcal{M}_0 \subset \mathcal{M}$ and adding environments ν from \mathcal{M} over time to \mathcal{M}_k . [SH15]
- Fazit: There exist (weakly) asymptotically optimal policies for arbitrary (separable) /compact \mathcal{M} .

Optimism in MDPs and Beyond

- Let \mathcal{M} be the class of all MDPs with $|\mathcal{S}| < \infty$ states and $|\mathcal{A}| < \infty$ actions and geometric discount γ .
- Then \mathcal{M} is continuous but compact
 - $\implies \pi^{\circ}$ is asymptotically optimal by previous slide.
- But much better polynomial error bounds in this case are possible:

Theorem 6 (PACMDP bound) $V_{k\gamma}^{\pi^{\circ}\mu} \leq V_{k\gamma}^{*\mu} - \varepsilon$ for at most $\tilde{O}(\frac{|\mathcal{S}|^2|\mathcal{A}|}{\varepsilon^2(1-\gamma)^3}\log\frac{1}{\delta})$ time steps k with probability $1-\delta$. [LH14b]

Similar bounds for General Optimistic Agents possible if environments are generated by combining laws (of nature): Laws predict only some feature (factorization) in some context (localization). [SH15]

APPROXIMATIONS & APPLICATIONS

Towards Practical Universal AI

Goal: Develop efficient general-purpose intelligent agent

• Additional Ingredients:

Main Reference (year)

- Universal search: Schmidhuber (200X) & al.
- Learning: TD/RL Sutton & Barto (1998) & al.
- Information:
- Complexity/Similarity:
- Optimization:
- Monte Carlo:

MML/MDL Wallace, Rissanen

Li & Vitanyi (2008)

Aarts & Lenstra (1997)

Fishman (2003), Liu (2002)

Computational Issues: Universal Search

- Levin search: Fastest algorithm for inversion and optimization problems.
- Theoretical application:

Assume somebody found a non-constructive proof of P=NP, then Levin-search is a polynomial time algorithm for every NP (complete) problem.



- Practical versions: OOPS and Levin Tree Search [Sch04, OHL23] Appl.: Mazes, towers of Hanoi, robotics, Rubik's cube, Sokoban, ...
- FastPrg: The asymptotically fastest and shortest algorithm for all well-defined problems. [Hut02]
- Computable Approximations of AIXI: [HQC24] AIξ, AIXI*tl*, MDP-AIXI, MC-AIXI-CTW, Self-AIXI.
- Human Knowledge Compression Prize: (500'000€)



Computable Approximations of AIXI

- Al ξ : Bayesian mixtures ξ over smaller classes \mathcal{M} [Hut05]
- AIXI*tl*: Search in proof and program space for provably optimal policy within given time and length bound similar to FastPrg [Hut07]
- MDP-AIXI: AI ξ for MDP class applied to 2x2 Matrix Games [PH06]
- MC-AIXI-CTW: AI ξ for CTW class with MCTS planning [VNH⁺11]
- Self-AIXI: Avoids expensive planning by self-predicting its own stream of action data [CGMH⁺23]
- PhiMDP: An alternative approach to Universal AI based on learning reductions from histories to MDP states [Hut09b]
- ExSAgg: Extreme reduction of histories to surrogate MDPs [Hut16]
- AIXIjs: Implementation of various history-based RL agents [ALH17]

ξ

ture reward estimation

A Monte-Carlo AIXI Approximation

Consider class of Variable-Order Markov Decision Processes.

The Context Tree Weighting (CTW) algorithm can efficiently mix (exactly in essentially linear time) all prediction suffix trees.

Monte-Carlo approximation of expectimax tree: Upper Confidence Tree (UCT) algorithm:

- Sample observations from CTW distribution.
- Select actions with highest upper confidence bound.
- Expand tree by one leaf node (per trajectory).
- Simulate from leaf node further down using (fixed) playout policy.
- Propagate back the value estimates for each node.

Repeat until timeout.

Guaranteed to converge to exact value.

Extensions in many directions exist

[VSH12, GBVB13]

[VNH⁺11]

Monte-Carlo AIXI Applications

without providing any domain knowledge, the same agent is able to self-adapt to a diverse range of interactive environments.



Extensions of MC-AIXI-CTW [VSH12]

- Smarter than random playout policy, e.g. learnt CTW policy.
- Extend the model class to improve general prediction ability. However, not so easy to do this in a comput. efficient manner.
- Predicate CTW: Context is vector of (general or problem-specific) predicate=feature=attribute values.
- Convex Mixing of predictive distributions. Competitive guarantee with respect to the best fixed set of weights.
- Switching: Enlarge base class by allowing switching between distr. Can compete with best rarely changing sequence of models.
- Improve underlying KT Est.: Adaptive KT, Window KT, KT0, SAD
- Partition Tree Weighting technique for piecewise stationary sources with breaks at/from a binary tree hierarchy.
- Mixtures of factored models such as quad-trees for images [GBVB13]
- Avoid MCTS by compression-based value estimation. [VBH⁺15]

Feature Reinforcement Learning (FRL)

- Basic Idea: Learn best reduction Φ of history to an MDP [Hut09b]
- Theoretical guarantees: Asymptotic consistency. [SH10]
- Example Φ -class: As Φ choose class of suffix trees as in CTW.
- How to find/approximate Φ^{best} :
 - Exhaustive search for toy problems [Ngu13]
 - Monte-Carlo (Metropolis-Hastings / Simulated Annealing) for approximate solution [NSH11]
 - Exact "closed-form" by CTM similar to CTW [NSH12]
- Experimental results: Comparable to MC-AIXI-CTW [NSH12]

Feature Reinforcement Learning (ctd)

- Extensions/Improvements:
 - Looping suffix trees for long-term memory [DSH12, DSH14a]
 - Structured/Factored MDPs (Dynamic Bayesian Networks) [Hut09a]
 - Extreme State Aggregation beyond MDPs [Hut16, MH19, MH21b]
 - Exact Binarization of Huge Action Spaces
- Related:
 - Q-Learning for History-Based Reinforcement Learning [[
 - Convergence of Q-Learning Beyond MDPs
 - Non-Convergence of Temporal-Difference-Like Methods with Linear Function Approximation
 - Reinforcement Learning with Value Advice

[DSH13] [MH18]

[HYZM19]

[DSH14b]

[MH21a]

DISCUSSION

Intelligent Agents in Perspective



Agents = General Framework, Interface = Robots, Vision, Language

Marcus Hutter

Hutter - 44 - Universal Artificial Intelligence **Aspects of Intelligence** are all(?) either directly included in AIXI or are emergent

HOW INCLUDED IN AIXI

TRAIT OF INTELL. reasoning creativity association generalization pattern recognition problem solving memorization planning achieving goals learning optimization self-preservation vision language motor skills classification induction deduction

to improve internal algorithms (emergent) exploration bonus, randomization, ... for co-compression of similar observations for compression of regularities in perceptions for compression how to get more reward storing historic perceptions searching the expectimax tree by optimal sequential decisions Bayes-mixture and belief update compression and expectimax by coupling reward to robot components observation=camera image (emergent) observation/action = audio-signal (emergent) action = movement (emergent) by compression Universal Bayesian posterior (Ockham's razor) Correctness proofs in AIXItl

Relation of UAI to Deep Learning - Current How LLMs can be regarded as approximations of AIXI

- Language modelling = minimizing log-loss = compression $[DRD^+24]$
- In-context learning as implicit Bayesian inference [XRLM22, GDR⁺23]
- Meta-Learning on Algorithmic Data: A step towards Solomonoff Induction [GMGH⁺24]
- Tree-of-Thought corresponds to MCTS planning [YYZ⁺23] (In-context learning is the analog of CTW updates within MCTS)
- RLHF is a very crude form of RL (horizon 1) [ZSW⁺20]
- Multimodal Transformers: AIXI is universal and agnostic to the meaning of the I/O bitstream, so automatically multimodal.

Relation of UAI to Deep Learning - Future What's missing in LLMs: How to make them closer to AIXI

- Continual learning to "compile" new experiences from in-context short-term memory to in-weight long-term memory [KRRP22]
- Learn a proper long-horizon value function
- Using this value function in chain-of-thought sampling should mimic MCTS quite well
- Reasoning capabilities of current LLMs are still limited. Possibly new NN architectures are needed, but maybe current architecture with proper scaffolding (better Tree-of-Thought, Tool-Use, ...) suffices.

Fazit: Gold-standard AIXI can guide in which directions to develop LLMs.

Recent Progress

 Dynamic Knowledge Injection for AIXI Agents 	[YZNH24]
• Transformers Learning Universal Predictors	GMGH ⁺ 24]
• Self-Predictive Universal AI [CGMH ⁺ 23]
 Language Modeling Is Compression 	[DRD ⁺ 24]
 Universal Agent Mixtures & Geometry of Intelligence 	[AQDH23]
• AIXI intervenes in the provision of reward [CH0	O22, CH22]
• Binarization of Rewards, Actions, Observations [CHV22a	a, CHV22b]
Reward-Punishment Symmetric Universal Intelligence	[AH21]
• Quantum Computing Algorithms for Universal Prediction	[CH20a]
• On the Computability of Solomonoff Induction and AIXI	[LH18]
- Generalised Discount Functions applied to a Monte-Carlo AI μ	
Implementation	[LALH17]
 Loss Bounds and Time Complexity for Speed Priors 	[FLH16]
 Suitable versions of AIXI are limit-computable 	[LH15]

ASI Safety

- 48 -

- Chances and Risks of Artificial Intelligence for Society [HH21]
- Curiosity Kills the Asymptotically Optimal Agent [CHC21]
- Unambitious AIXI with short horizon is Safe [CVH21, CVH20]
- Pessimism About Unknowns Inspires Conservatism [CH20b]
- The Alignment Problem for Universal AI [EH18, EKH19]
- Reinforcement Learning with a Corrupted Reward Channel [EKO⁺17]
- Avoiding Wireheading with Value Reinforcement Learning [EH16]
- Self-Modification of Policy and Utility in UAI [EFDH16]

Outlook 1

- 49 -

- Find optimality notions for generally intelligent agents which are strong enough to be convincing but weak enough to be satisfiable.
- More powerful and faster computational approximations of AIXI
- Social questions about AIXI or other Super-Intelligences: socialization, rewards, drugs, suicide, self-improvement, manipulation, attitude, curiosity, immortality, self-preservation.
- Training (sequence): To maximize informativeness of reward, one should provide a sequence of simple-to-complex tasks to solve, with the simpler ones helping in learning the more complex ones.

Outlook 2

- Address the many open theoretical questions in [Hut05].
- Bridge the gap between (Universal) AI theory and AI practice.
- Explore what role logical reasoning, knowledge representation, vision, language, etc. play in Universal AI.
- Determine the right discounting of future rewards.
- Develop the right nurturing environment for a learning agent.
- Consider embodied agents (e.g. internal ↔ external reward)
- Analyze AIXI in the multi-agent setting (done) [LTF16, FTC15]

Thanks! Questions? Details:

A Unified View of Artificial Intelligence

Decision Theory = Probability + Utility Theory + + Universal Induction = Ockham + Bayes + Turing

Open research problems:

at www.hutter1.net/ai/uaibook.htm

Compression contest:

with 500'000€ prize at prize.hutter1.net

Projects: <u>www.hutter1.net/official/projects.htm</u>



New Book (2024)