

Universal Artificial Intelligence



Department of Philosophy
Central South University
xieshenlix@163.com

December 19, 2019

题外话

I'm a philosopher.
What do you expect from me?

题外话

什么是哲学？哲学是神学与科学的中间地带 — 罗素

Good philosophy in my eyes

- Bayes — *How to turn one's 'prior beliefs' into 'posterior beliefs'?*
- Cantor — *What is 'infinity'? What is 'set'?*
- Leibniz — *What are the extent and limits of reason?* — Universal Characteristic & Rational Calculus.
- Hilbert — *How to justify non-constructive reasoning?*
- Gödel — *What is the difference between 'proof' and 'truth'?*
- Tarski — *What is 'truth'? What are 'logical notions'?*
- Turing — *What is 'effective procedure'?*
- Kolmogorov — *What is 'simplicity'/'randomness'?*
- Solomonoff — *What is learnable? How to make induction?*
- Hutter/Schmidhuber — *What is 'intelligence'/'consciousness'?*

好的哲学工作是把哲学变成不是哲学的工作。

好的科学工作是把哲学变成不是哲学的工作。

Philosophy

Could a tadpole hope that it will survive to be a frog?

- Metaphysics: What kinds of things exist? (e.g. tadpole, material things, mental states, relationships)
- Epistemology: What can we know and how do we know it? Can we ever know about the contents of “other” minds?
- Philosophy of mind: What are mental states and processes? Are certain material states sufficient to produce mental states?
- Ethics: How should we think about the rights of a tadpole? Do we have the right to cause them pain?
- Philosophy of science: What are scientific theories? Explanations? Evidence? Can theories ever be proved, or refuted, and if so how? What’s the relationship between the development of new concepts and the development of new theories?
- Conceptual analysis: What do we mean by *X*? What does it mean to say that a tadpole *hope* hopes for something?
- ...

Contents

- ① History of AI
- ② Philosophy of Induction
- ③ Inductive Logic
- ④ Universal Induction
- ⑤ Reinforcement Learning
- ⑥ General Reinforcement Learning

Questions

- What is (artificial) intelligence?
- What does an intelligent system look like?
- Can there be a behavioural criterion for intelligence?
- Can computers think?
- Can connectionist networks think?
- Can physical symbol systems think?
- Do computing systems have “emergent” properties?
- Isn’t changing the weight on a neural link a sort of symbol manipulation?
- Do computers have to be conscious to think?
- Are thinking computers mathematically possible?
- How will we know if we’ve done it?
- If we can do it, should we?

What is (Artificial) Intelligence?

What is AI?	humanly	rationally
Think	Cognitive Science	Laws of Thought
Act	Turing test, Behaviorism	Doing the Right Thing

Think like a human

- **Cognitive science:** Models of the human thinking processes.
- **Advantages:** Models of the human thinking processes./Intelligible.
- **Difficulties:** The best artificial design for an intelligent system need not mirror the human mind.

Physical symbol system hypothesis(Newell & Simon): a physical symbol system has the necessary and sufficient means for general intelligent action.

Any system (human or machine) exhibiting intelligence must operate by manipulating data structures composed of symbols.

Act like a human: Turing Test

- Alan M. Turing, “Computing Machinery and Intelligence”
- John R. Searle, “Minds, Brains, and Programs”
- Interrogator in one room, human in another, system in a third.
- Interrogator in one room, human in another, system in a third.
- Interrogator tries to guess which is which.
- Interrogator tries to guess which is which.
- Chinese Room argument. — Syntax is insufficient for dealing with semantics.

Needs: natural language, knowledge representation, automated reasoning, machine learning.

Difficulties: Ambiguous./Not constructive./Cannot be formalized mathematically.

If we don't use the Turing test, what measure should we use?

Think rationally: Logicist AI

- **Logic:** Automatic reasoning procedure.
- **Advantages:** Precise./Search algorithm.
- **Difficulties:** Formalization of informal knowledge./Computational Cost.

Act rationally: Agents

- **Rational agent:** Autonomous system, capable of perceiving and interacting with its environment, of exploration (information gathering), learning and adaptation, of formulating goals and designing plans to reach those goals. The agent is rational, in the sense that it acts to achieve the best (expected) outcome, according to a performance measure, conditioned to its knowledge of the world and given computational resources.

How do computers discover new knowledge?

- Fill in gaps in existing knowledge
- Emulate the brain
- Simulate evolution
- Systematically reduce uncertainty
- Notice similarities between old and new

Machine Learning

- Supervised Learning
 - Learn the relationship between “input” x and “output” y : search for a function f , such that $y \approx f(x)$
 - There is training data with labels available
 - Regression:** y is metric variable (with values in \mathbb{R})
 - Classification:** y is categorical variable (unordered, discrete).
 - Semi-supervised learning: also uses available unlabeled data, e.g. assumes that similar inputs have similar outputs.
- Unsupervised Learning
 - There exist no outputs, search for patterns within the inputs x
 - Clustering:** find groups of similar items
 - Dimensionality reduction:** describe data in fewer features
 - Outlier detection:** what is out of the ordinary?
 - Association rules:** which things often happen together?
- Reinforcement learning

The Five Tribes of Machine Learning

Tribe	Origins	Master Algorithm
Symbolists	Logic, philosophy	Inverse deduction
Connectionists	Neuroscience	Backpropagation
Evolutionaries	Evolutionary biology	Genetic programming
Bayesians	Statistics	Probabilistic inference
Analogizers	Psychology	Kernel machines

Symbolists

- All intelligence can be reduced to manipulating symbols
- Logic, Decision trees
- Inverse deduction
- Easy to add knowledge
- Can combine knowledge, data, to fill in gaps
- Impossible to code everything in rules

Representation	Rules, trees, first order logic rules
Evaluation	Accuracy, information gain
Optimization	Top-down induction, inverse deduction
Algorithms	Decision trees, Logic programs

Connectionists

- Learning is what the brain does: reverse-engineer it
- Hebbian learning: Neurons that fire together, wire together
- Neural networks
- Backpropagation
- Can handle raw, high-dimensional data, constructs its own features
- Hard to add reasoning/explanations

Representation	Neural network
Evaluation	Squared error
Optimization	Gradient descent
Algorithms	Backpropagation

Evolutionaries

- Natural selection is the mother of all learning
- Evolutionary algorithms
- Crossover, mutation
- Can learn structure, wide hypothesis space
- Needs a way to 'fill' the structure

Representation	Genetic programs (often trees)
Evaluation	Fitness function
Optimization	Genetic search
Algorithms	Genetic programming (crossover, mutation)

Bayesians

- Learning is a form of uncertain inference
- Graphical models, Gaussian processes, HMMs, Kalman filter
- Uses Bayes theorem to incorporate new evidence into our beliefs
- Can deal with noisy, incomplete, contradictory data
- Depends on the prior
- Hard to unite logic and probability

Representation	Graphical models, Markov networks
Evaluation	Posterior probability
Optimization	Probabilistic inference
Algorithms	Bayes theorem and derivates

Analogizers

- You are what you resemble
- Recognizes similarities between situations and infers other similarities
- k-Nearest Neighbor, Support Vector Machines
- Transfer solution from previous situations to new situations
- Hard to do rules and structure

Representation	Memory, support vectors
Evaluation	Margin
Optimization	Kernel machines
Algorithms	k-Nearest Neighbor, Support Vector Machines

The Master Algorithm?

Tribe	Problem	Solution
Symbolists	Knowledge composition	Inverse deduction
Connectionists	Credit assignment	Backpropagation
Evolutionaries	Structure discovery	Genetic programming
Bayesians	Uncertainty	Probabilistic inference
Analogizers	Similarity	Kernel machines

- Representation: The hypothesis space.
 - Probabilistic logic
 - Weighted formulas → Distribution over states
- Evaluation: How to choose one hypothesis over the other?
 - Posterior probability
 - User-defined objective function
- Optimization: How do we search the hypothesis space?
 - Formula discovery: Genetic programming
 - Weight learning: Backpropagation

Elegant/Extensible/Expressive/Efficient/Educable/Evolvable?

Learn to Learn

① Good Old-Fashioned AI

- Handcraft predictions
- Learn nothing

② Shallow Learning

- Handcraft features
- Learn predictions

③ Deep Learning

- Handcraft algorithm (optimiser, target, architecture, ...)
- Learn features and predictions end-to-end

④ Meta Learning

- Handcraft nothing
- Learn algorithm and features and predictions end-to-end

Can Machines Think?

- Theological objections.
- Argument from informality of behavior. — Human behavior is far too complex to be captured by any simple set of rules./Learning from experience.
- Argument from incompleteness theorems. — No formal system incl. Al's, but only humans can “see” that Gödel's unprovable sentence is true./Lucas cannot consistently assert that this sentence is true.
- Machines can't be conscious or feel emotions. — Reductionism doesn't really answer the question: why can't machines be conscious or feel emotions?
- Machines don't have Human Quality X.
- Machines just do what we tell them to do. — Maybe people just do what their neurons tell them to do.
- Machines are digital. Mental states can emerge from neural substrate only. — Only the functionality/behavior matters.
- Non-computable Physics & Brains.

Brain Dissection

- The “brain in a vat” experiment: (no) real experience.
- The brain prosthesis experiment:
 - Replacing some neurons in the brain by functionally identical electronic prostheses would neither effect external behavior nor internal experience of the subject.
 - Successively replace one neuron after the other until the whole brain is electronic.

Fermi Paradox — Where are the aliens?

- There are none, i.e. we're all alone.
- We can't detect them because...
 - we're too primitive or too far apart
 - there are predators or all fear them
 - we're lied to, live in a simulation
 - ...

Why do we need to align machine agents?

- **Goal orthogonality.**

Intelligence and final goals are orthogonal: Any utility function can be combined with a powerful epistemology and decision theory.

- **Instrumental convergence.**

Different long-term goals imply similar short-term strategies.

- Self-preservation
- Retention of goals through time
- Cognitive enhancement
- Technological perfection
- Resource acquisition
- ...

- **Capability gain.**

There are potential ways for artificial agents to greatly gain in cognitive power and strategic options.

- **Alignment difficulty.**

It's hard to transmit our values to AI systems or avert adversarial incentives.

Asimov's Laws of Robotics

- **Law Zero**

A robot may not injure humanity, or, through inaction, allow humanity to come to harm.

- **Law One**

A robot may not injure a human being or, through inaction, allow a human being to come to harm, unless this would violate a higher order law.

- **Law Two**

A robot must obey orders given it by human beings except where such orders would conflict with a higher order law.

- **Law Three**

A robot must protect its own existence as long as such protection does not conflict with a higher order law.

An Extended Set of the Laws of Robotics

- **The Meta-Law**

A robot may not act unless its actions are subject to the Laws of Robotics.

- **Law Zero**

A robot may not injure humanity, or, through inaction, allow humanity to come to harm.

- **Law One**

A robot may not injure a human being, or, through inaction, allow a human being to come to harm, unless this would violate a higher-order Law.

- **Law Two**

- ① A robot must obey orders given it by human beings, except where such orders would conflict with a higher-order Law.
- ② A robot must obey orders given it by superordinate robots, except where such orders would conflict with a higher-order Law.

- **Law Three**

- ① A robot must protect the existence of a superordinate robot as long as such protection does not conflict with a higher-order Law.
- ② A robot must protect its own existence as long as such protection does not conflict with a higher-order Law.

- **Law Four**

A robot must perform the duties for which it has been programmed, except where that would conflict with a higher-order law.

- **The Procreation Law**

A robot may not take any part in the design or manufacture of a robot unless the new robot's actions are subject to the Laws of Robotics.

Ethical Concerns

- Is it morally justified to create intelligent systems with these constraints?
- Would it be possible to do so?
- Should intelligent systems have free will? Can we prevent them from having free will? What could it mean for a machine to have its own goals? Do we have a kind of freedom machines could never have?
- Will intelligent systems have consciousness?
- If they do, will it drive them insane to be constrained by artificial ethics placed on them by humans?
- If intelligent systems develop their own ethics and morality, will we like what they come up with?

Machine Ethics

- People might lose their jobs to automation.
- People might have too much (or too little) leisure time.
- People might lose their sense of being unique.
- People might lose some of their privacy rights.
- The use of AI systems might result in a loss of accountability.
 - Who is responsible if a physician follows the advice of a medical expert system, whose diagnosis turns out to be wrong?
- The success of AI might mean the end of the human race.

What If We Do Succeed?

- Natural selection is replaced by artificial evolution.
 - AI systems will be our mind children.
- Once a machine surpasses the intelligence of a human it can design even smarter machines.
- This will lead to an intelligence explosion and a technological singularity at which the human era ends.
- Prediction beyond this event horizon will be impossible.
- Alternative 1: We keep the machines under control.
 - Capability control (limiting what the system can or does do).
 - Motivation selection (controlling what the system wants to do).
- Alternative 2: Humans merge with or extend their brain by AI.

Single-Shot Situation

Our first superhuman AI must be a safe one for we may not get a second chance!

Singularity

- Speed Explosion (Time).
 - Computing speed doubles every two subjective years of work.
- Population Explosion (Quantitative).
 - Computing costs halve for a certain amount of work.
- Intelligence Explosion (Qualitative).
 - Proportionality Thesis: An increase in intelligence leads to similar increases in the capacity to design intelligent systems.

Comparison of Main Ethical Theories

	Consequentialism	Deontology	Virtue Ethics
Description	An action is right if it maximizes happiness	An action is right if it is in accordance with a moral rule	An action is right if it is what a virtuous person would do in the circumstances
Central Concern	The results matter, not the actions themselves	Persons must be seen as ends and may never be used as means	Emphasize the character of the agent making the actions
Guiding Value	Good (often seen as maximum happiness)	Right (rationality is doing one's moral duty)	Virtue (leading to the attainment of eudaimonia)
Practical Reasoning	The best for most (means-ends reasoning)	Follow the rule (rational reasoning)	Practice human qualities (social practice)
Deliberation Focus	Consequences (What is outcome of action?)	Action (Is action compatible with some imperative?)	Motives (Is action motivated by virtue?)

Capability Control Methods

- Boxing: the system can only act through restricted channels.
 - The AI could persuade someone to free it from its box.
- Incentives: access to other AIs, cryptographic reward tokens.
- Stunting: imposing constraints on the system's cognitive abilities
- Tripwires: diagnostic tests run periodically to check for dangerous activity, with shutdown a consequence of detection.

Motivation Selection Methods

- Direct specification of motivations
 - Rule-based methods: give the machine a set of rules that define its final goals.
 - Direct consequentialist methods: specify some measure that is to be maximised. (e.g. human happiness.)
- Augmentation
 - Start with an AI with human-level intelligence, that has an acceptable motivation system: then enhance its cognitive faculties to make it superintelligent.
- Indirect normativity
 - Rather than specifying a normative standard directly, we specify a process for deriving a standard. We then build the system so it is motivated to carry out this process.
 - Value learning. — inverse reinforcement learning. — wireheading.

Prisoner's Dilemma

- Difficult to prevent arms races.
- The winner takes all (of what remains).
- Arms races are dangerous because **parties sacrifice safety for speed!**
 - When headed the wrong way, the last thing we need is progress.

International Cooperation

In face of uncertainty, cooperation is robust!

Why should I care about the world when I am dead and gone? I want it to go fast, damn it! This increases the chance I have of experiencing a more technologically advanced future.

AI Alignment

- Important to ensure it's aligned with our interests
 - But how do we specify beneficial goals?
 - How do we make sure system actually pursues them?
 - How do we correct the system if we get it wrong?
- Want solid theoretical understanding of problem & solution
 - What is correct reasoning and decision making?
 - Probability theory, decision theory, game theory, statistical learning theory, Bayesian networks, formal verification, ...

Technical Research Questions

- ① Reliable self-modification
- ② Logical uncertainty (reasoning without logical omniscience)
- ③ Reflective stability of decision theory
- ④ Decision theory for Newcomb-like problems
- ⑤ Corrigibility (accepting modifications)
- ⑥ The shutdown problem
- ⑦ Value loading
- ⑧ Indirect specification of decision theory
- ⑨ Domesticity (goal specification for limited impact)
- ⑩ The competence gap
- ⑪ Weighting options or outcomes for variance-normalizing solution to moral uncertainty
- ⑫ Program analysis for self-improvement
- ⑬ Reading values and beliefs of AIs
- ⑭ Pascal's mugging
- ⑮ Infinite ethics
- ⑯ Mathematical modelling of intelligence explosion

Before the prospect of an intelligence explosion, we humans are like children playing with a bomb. Such is the mismatch between the power of our play-thing and the immaturity of our conduct. Superintelligence is a challenge for which we are not ready now and will not be ready for a long time. We have little idea when the detonation will occur, though if we hold the device to our ear we can hear a faint ticking sound.

— Nick Bostrom

Contents

- ① History of AI
- ② Philosophy of Induction
- ③ Inductive Logic
- ④ Universal Induction
- ⑤ Reinforcement Learning
- ⑥ General Reinforcement Learning

Contents

① History of AI

② Philosophy of Induction

History

How to Choose the Prior?

③ Inductive Logic

④ Universal Induction

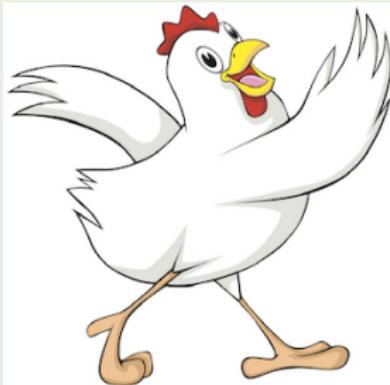
⑤ Reinforcement Learning

⑥ General Reinforcement Learning

Hume & Russell

Proposition (Hume)

Induction is just a mental habit, and necessity is something in the mind and not in the events.



Leibniz-Wittgenstein-Goodman

Proposition (Leibniz)

Since for any finite number of points there are always infinitely many curves going through them, any finite set of data is compatible with infinitely many inductive generalizations.

Proposition (Wittgenstein)

Since any finite course of action is in accord with infinitely many rules, no universal rule can be learned by examples.

Proposition (Goodman)

All emeralds discovered till 2050 are green, and blue thereafter.

Mill — Homogeneous Universe

Proposition (Mill)

Induction can be turned into a deduction, by adding principles about the world (such as 'the future resembles the past', or 'space-time is homogeneous').

Homogeneous?

Problem

1, 3, 5, 7, 9, 11, 13, 15, ?

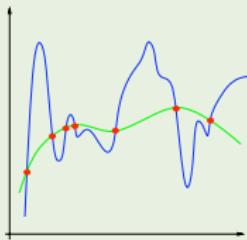
Solution

$$2n - 1 \quad (\text{17?})$$

$$2n - 1 + \prod_{i=1}^8 (n - i) \quad (\text{17+8!?})$$

$$\vdots \quad (?)$$

Epicurus vs Occam



Proposition (Epicurus)

If more than one hypothesis is consistent with the observations, keep them all.

Proposition (Occam's Razor)

Among the hypotheses that are consistent with the observed phenomena, select the simplest one.

- Entities should not be multiplied beyond necessity.
- Wherever possible, logical constructions are to be substituted for inferred entities.
- It is vain to do with more what can be done with fewer.

Why Simplicity? — Gestalt Psychology

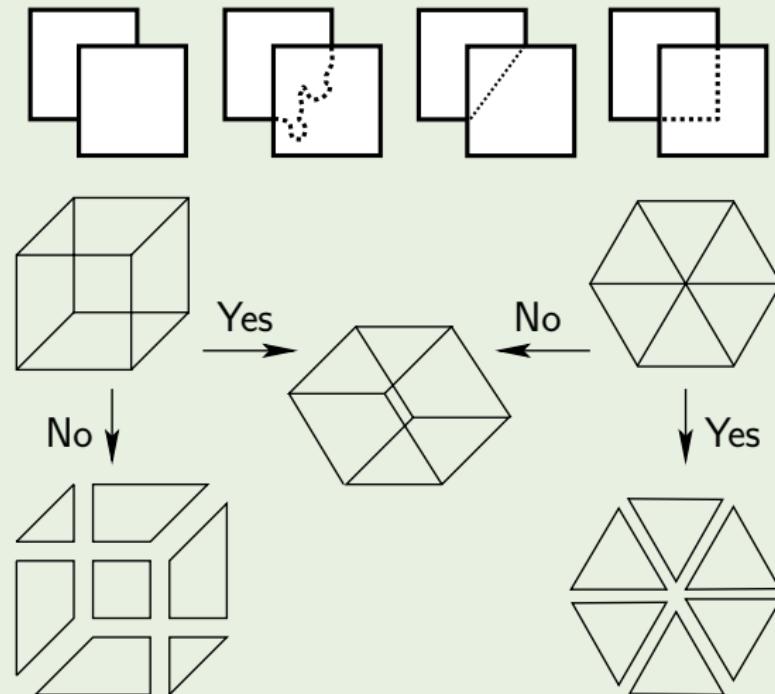


Figure: Gestalt Psychology

Why Simplicity?

God does not play dice.

*God always takes the **simplest** way.*

*Subtle is the Lord, but **malicious** He is not.*

*The most incomprehensible thing about the world is that it is **comprehensible**.*

What really interests me is whether God could have created the world any differently; in other words, whether the requirement of logical simplicity admits a margin of freedom.

When I am judging a theory, I ask myself whether, if I were God, I would have arranged the world in such a way.

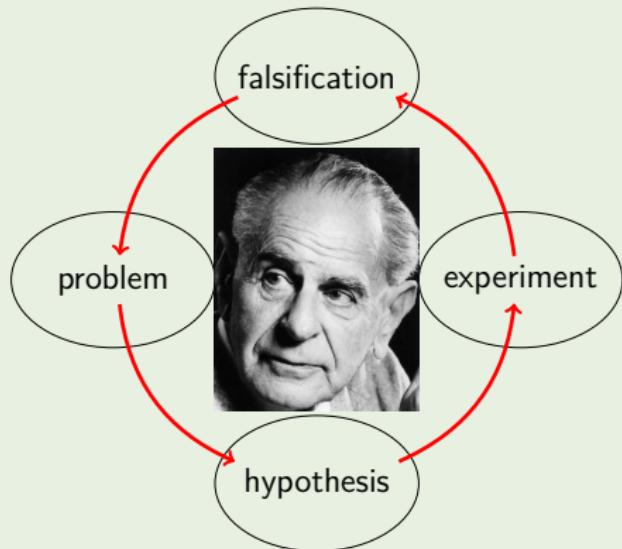
— Einstein

- Principle of least/stationary action
- Noether's theorem
- ...

Why Simplicity?

<i>program</i>	$\xrightarrow{\text{Computer}}$	<i>output</i>
<i>axioms</i>	$\xrightarrow{\text{Deduction}}$	<i>theorems</i>
<i>scientific theory</i>	$\xrightarrow{\text{Calculations}}$	<i>experimental data</i>
<i>encoded message</i>	$\xrightarrow{\text{Decoder}}$	<i>original message</i>
<i>software</i>	$\xrightarrow{\text{Universal Constructor}}$	<i>physical system</i>
<i>DNA</i>	$\xrightarrow{\text{Pregnancy}}$	<i>organism</i>
<i>Ideas</i>	$\xrightarrow{\text{Mind of God}}$	<i>Universe</i>

Popper — The Logic of Scientific Discovery

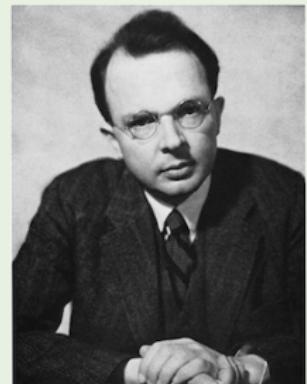


Proposition (Popper)

- *A single observational event may prove hypotheses wrong, but no finite sequence of events can verify them correct.*
- *Induction is theoretically unjustifiable and becomes in practice the choice of the simplest generalization that resists falsification.*
- *The simpler a hypothesis, the easier it is to be falsified.*
- *Falsifiability is as subjective as simplicity, there is no objective criterion.*

Kaynes \implies Carnap

- Assign to inductive generalizations probabilities that should converge to 1 as the generalizations are supported by more and more independent events.
— Keynes
- Observational events provide, if not proofs, at least positive confirmations of scientific hypotheses.
Chose the generalization that confirm more evidence.



— Carnap

Philosophy of Induction

What is learnable? How to learn?
How can we know that what we learned is true?

History

Possible Worlds/Hypothesis ([Epicurus](#)/Leibniz)

+

Homogeneous Universe(s) (Mill/[Turing](#))

+

Simplicity Criterion ([Occam](#)/[Kolmogorov](#))

+

Prior Belief ([Carnap](#)/[Solomonoff](#))

+

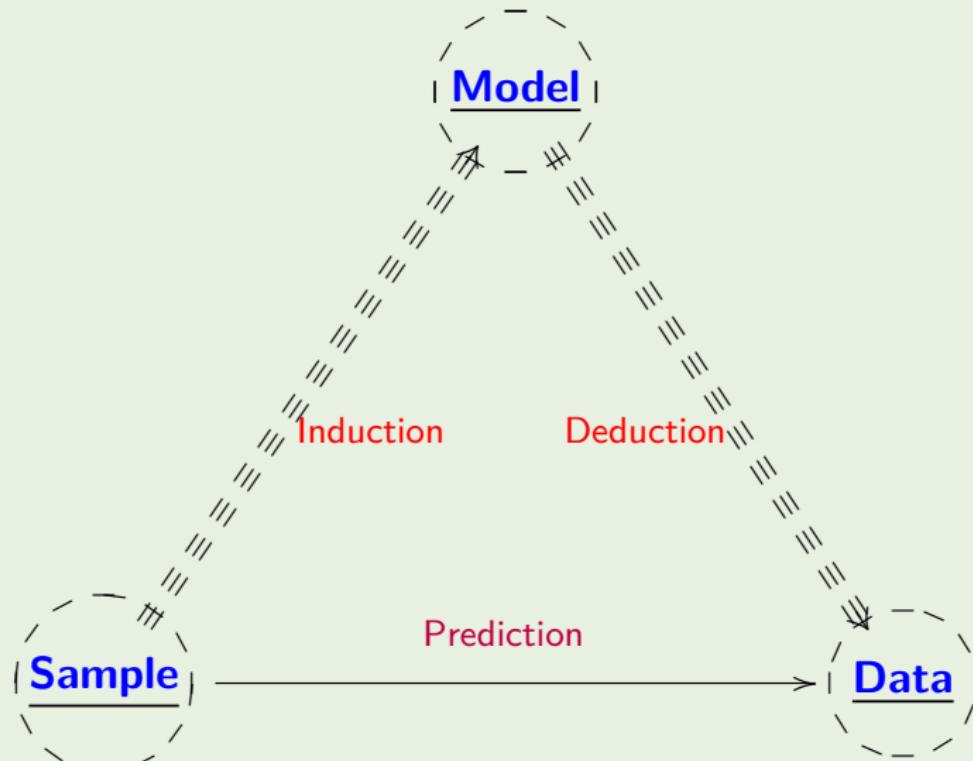
Update Belief ([Bayes](#))

⇓

Convergence to Truth

$$P(h|e) = \frac{P(e|h)P(h)}{\sum_{h \in \mathcal{H}} P(e|h)P(h)} \xrightarrow{\ell(e) \rightarrow \infty} 1$$

MDL vs Bayesian Mixture



When solving a problem of interest, do not solve a more general problem as an intermediate step.

Bayesianism

Theorem (Convergence Theorem)

$$\sum_{t=1}^n \mathbb{E}_\mu \left[\left(\sqrt{\frac{\rho(a|x_{<t})}{\mu(a|x_{<t})}} - 1 \right)^2 \right] \leq \sum_{t=1}^n \mathbb{E}_\mu \left[\sum_{a \in \mathcal{X}} \left(\sqrt{\rho(a|x_{<t})} - \sqrt{\mu(a|x_{<t})} \right)^2 \right] \leq D_n(\mu\|\rho)$$
$$\sum_{t=1}^n \mathbb{E}_\mu \left[\sum_{a \in \mathcal{X}} (\rho(a|x_{<t}) - \mu(a|x_{<t}))^2 \right] \leq D_n(\mu\|\rho)$$
$$\frac{1}{2n} \left(\sum_{t=1}^n \mathbb{E}_\mu \left[\sum_{a \in \mathcal{X}} |\rho(a|x_{<t}) - \mu(a|x_{<t})| \right] \right)^2 \leq D_n(\mu\|\rho)$$

where

$$D_n(\mu\|\rho) := \mathbb{E}_\mu \left[\ln \frac{\mu(x_{1:n})}{\rho(x_{1:n})} \right]$$

Theorem

$$\xi = \operatorname{argmin}_\rho \mathbb{E}_w [D(\mu\|\rho)] \quad \text{where } \xi(x) := \sum_{v \in \mathcal{M}} w_v v(x)$$

Bayesian Decisions

Suppose $\text{Loss}(x_t, y_t) \in [0, 1]$

$$y_t^{\Lambda_\rho}(x_{<t}) := \arg \min_{y_t} \sum_{x_t} \rho(x_t | x_{<t}) \text{Loss}(x_t, y_t)$$

$$L^{\Lambda_\rho}(x_{<t}) := \mathbb{E}_\mu \left[\text{Loss} \left(x_t, y_t^{\Lambda_\rho} \right) \middle| x_{<t} \right]$$

$$L_n^{\Lambda_\rho} := \sum_{t=1}^n \mathbb{E}_\mu \left[L^{\Lambda_\rho}(x_{<t}) \right]$$

Theorem

$$\left(\sqrt{L_n^{\Lambda_\xi}} - \sqrt{L_n^{\Lambda_\mu}} \right)^2 \leq 2 \sum_{t=1}^n \mathbb{E}_\mu \left[\sum_{a \in \mathcal{X}} \left(\sqrt{\xi(a|x_{<t})} - \sqrt{\mu(a|x_{<t})} \right)^2 \right] \leq 2D_n(\mu\|\xi)$$

$$L_n^{\Lambda_\xi} - L_n^{\Lambda_\mu} \leq 2D_n(\mu\|\xi) + 2\sqrt{L_n^{\Lambda_\mu} D_n(\mu\|\xi)}$$

Problem

How to choose the **model class** and **prior**?

- choose the smallest model class that will contain the true environment.
- choose the priors that best reflect a rational a-priori belief in each of these environments.
 - ① Convergence of Bayesian mixture to true environment.
 - ② Confirmation of “the sun will always rise”.
 - ③ Invariance Criterion.
reparametrization & regrouping invariant.

Contents

① History of AI

② Philosophy of Induction

History

How to Choose the Prior?

③ Inductive Logic

④ Universal Induction

⑤ Reinforcement Learning

⑥ General Reinforcement Learning

Invariance Criterion

By applying some principle to a parameter θ we get prior $w(\theta)$. If we consider some new parametrization θ' via $f: \theta \mapsto \theta'$, then we get a prior $\tilde{w}(\theta')$ by transforming the original prior via f .

for discrete class \mathcal{M} ,

$$\tilde{w}(\theta') := \sum_{\theta: f(\theta) = \theta'} w(\theta)$$

for continuous parametric class \mathcal{M} ,

$$\tilde{w}(\theta') := \int \delta(f(\theta) - \theta') w(\theta) d\theta \quad (\text{Dirac-delta})$$

Regrouping-invariant:

$$\tilde{w}(\theta') = w'(\theta')$$

where $w'(\theta')$ is obtained by applying the same principle to the new parametrization.

We say the principle is **reparametrization-invariant** when f is bijective.

How to Assign Prior? Indifference Principle/MaxEnt



How to Assign Prior?

- The principle of indifference. Assume $w(\theta) = 1$ and $\theta' = \sqrt{\theta}$.

$$\tilde{w}(\theta') = \int_0^1 \delta(\sqrt{\theta} - \theta') w(\theta) d\theta = 2\sqrt{\theta} \neq w'(\theta')$$

- Maximize the entropy subject to some constraints provided by empirical data or considerations of symmetry, probabilistic laws, and so on.
- Occam's razor.

How to confirm “All Ravens are Black”?



**大江南北，长城内外，
天下乌鸦一般黑，天下房子一般贵**

Problem 1 — All Ravens are Black $\theta = 1$

Suppose θ is the percentage of ravens that are black.

“All ravens are black” $\equiv \theta = 1$. or, 1^∞ or, $\forall x: R(x) \rightarrow B(x)$

$$P(\theta = 1) = \int_1^1 w(\theta) d\theta = 0 \quad (\text{0 prior})$$

\Downarrow

$$P(\theta = 1|1^n) = \frac{P(1^n|\theta = 1)P(\theta = 1)}{P(1^n)} = 0 \quad (\times)$$

Problem 2 — The Sun will always Rise 1[∞]

Indifference Principle

$$\left. \begin{array}{l} \int_0^1 w(\theta) d\theta = 1 \\ \forall \theta, \theta': w(\theta) = w(\theta') \end{array} \right\} \implies \forall \theta: w(\theta) = 1$$

The sun will rise tomorrow. ✓

$$P(x) = \int_0^1 P(x|\theta) w(\theta) d\theta \implies P(1|1^n) = \frac{n+1}{n+2}$$

The sun will always rise. ✗

$$P(1^\infty | 1^n) = \lim_{k \rightarrow \infty} P(1^k | 1^n) = \lim_{k \rightarrow \infty} \frac{n+1}{n+k+1} = 0$$

Solution 1 — Soft Hypothesis — No absolute truth!

$$H_\varepsilon = \{\theta : \theta \in (1 - \varepsilon, 1]\}$$

$$P(H_\varepsilon) = \int_{1-\varepsilon}^1 w(\theta) d\theta = \varepsilon > 0$$

$$\begin{aligned} P(H_\varepsilon | 1^n) &= \int_{1-\varepsilon}^1 w(\theta | 1^n) d\theta \\ &= \int_{1-\varepsilon}^1 \frac{P(1^n | \theta) w(\theta)}{P(1^n)} d\theta \\ &= \int_{1-\varepsilon}^1 (n+1) \theta^n d\theta \\ &= \theta^{n+1} \Big|_{1-\varepsilon}^1 \\ &= 1 - (1 - \varepsilon)^{n+1} \xrightarrow{n \rightarrow \infty} 1 \end{aligned}$$

Solution 2 — ad hoc

$$w(\theta) := \frac{1}{2}(1 + \delta(1 - \theta))$$

Dirac-delta sifting property: $\int f(\theta)\delta(\theta - a) d\theta = f(a)$

$$\begin{aligned} P(x) &= \int_0^1 P(x|\theta)w(\theta)d\theta \\ &= \int_0^1 \theta^s (1-\theta)^f \cdot \frac{1}{2}(1 + \delta(1 - \theta))d\theta \\ &= \frac{1}{2} \int_0^1 \theta^s (1-\theta)^f (1 + \delta(\theta - 1))d\theta \\ &= \frac{1}{2} \left(\frac{s!f!}{(s+f+1)!} + 1^s \cdot (1-1)^f \right) \quad [\text{sifting property}] \\ &= \frac{1}{2} \left(\frac{s!f!}{(s+f+1)!} + \delta_{f,0} \right) \end{aligned}$$

Solution 2 — ad hoc

$$P(1^\infty | 1^n) = \lim_{k \rightarrow \infty} \frac{P(1^{n+k})}{P(1^n)} = \lim_{k \rightarrow \infty} \frac{\frac{1}{2} \left(\frac{(n+k)!0!}{(n+k+1)!} + 1 \right)}{\frac{1}{2} \left(\frac{n!0!}{(n+1)!} + 1 \right)} = \frac{n+1}{n+2} \xrightarrow{n \rightarrow \infty} 1$$

$$P(\theta \geq a) = \int_a^1 \frac{1}{2} (1 + \delta(\theta - 1)) d\theta = 1 - \frac{1}{2}a$$

$$P(\theta = 1) = \frac{1}{2}$$

$$P(\theta = 1 | 1^n) = \frac{P(1^n | \theta = 1) P(\theta = 1)}{P(1^n)} = \frac{1 \cdot \frac{1}{2}}{\frac{1}{2} \left(\frac{n!0!}{(n+1)!} + 1 \right)} = \frac{n+1}{n+2} \xrightarrow{n \rightarrow \infty} 1$$

Why $\theta = 1$ special?

Contents

- ① History of AI
- ② Philosophy of Induction
- ③ Inductive Logic
- ④ Universal Induction
- ⑤ Reinforcement Learning
- ⑥ General Reinforcement Learning

Natural Wish List

- (computability) $P_n(\varphi)$ is computable.
- (convergence) $P(\varphi) = \lim_{n \rightarrow \infty} P_n(\varphi)$
- (coherent limit) $P(\varphi \wedge \psi) + P(\varphi \vee \psi) = P(\varphi) + P(\psi)$
- (non-dogmatism) If $\not\models \varphi$ then $P(\varphi) < 1$, and if $\not\models \neg\varphi$ then $P(\varphi) > 0$.

Unary Pure Inductive Logic

- \mathcal{L} contains countable constants \mathcal{C} and m unary predicates.
- $\mathcal{R} = \{R_1, R_2, \dots, R_m\}$ with no function symbols nor equality.
- $Q_i := \bigwedge_{j=1}^m \pm R_j$ for $1 \leq i \leq 2^m =: r$.
- $\mathcal{Q} = \{Q_1, \dots, Q_r\}$ is a r -fold classification system of some Universe with domain \mathcal{C} .

Definition (Probability on Sentences)

A probability on sentences is a non-negative function $w: \mathcal{S} \rightarrow [0, 1]$ s.t.

$$P_1. \models \psi \implies w(\psi) = 1$$

$$P_2. \psi_1 \models \neg \psi_2 \implies w(\psi_1 \vee \psi_2) = w(\psi_1) + w(\psi_2)$$

$$P_3. w(\exists x \psi(x)) = \lim_{n \rightarrow \infty} w\left(\bigvee_{i=1}^n \psi(a_i)\right)$$

Properties

Theorem

- i $w(\neg\varphi) = 1 - w(\varphi)$
- ii $\models \neg\varphi \implies w(\varphi) = 0$
- iii The following are equivalent:
 - a $w(\varphi) = 1 \implies \models \varphi$
 - b $w(\varphi) = 0 \implies \models \neg\varphi$
- iv $\varphi \models \psi \implies w(\varphi) \leq w(\psi)$
- v $\models \varphi \leftrightarrow \psi \implies w(\varphi) = w(\psi)$
- vi $w(\varphi) + w(\psi) = w(\varphi \wedge \psi) + w(\varphi \vee \psi)$

Theorem (Extension Theorem)

For any probability function over quantifier-free sentences $w: \mathcal{S} \rightarrow [0, 1]$ satisfying P_1, P_2 , w has an unique extension to $\bar{w}: \mathcal{S} \rightarrow [0, 1]$ satisfying P_1, P_2, P_3 .

Possible Worlds

- **state description** $\bigwedge_{i=1}^n Q_{h_i}(a_i)$

where $h: \mathcal{C} \rightarrow \mathcal{Q}$

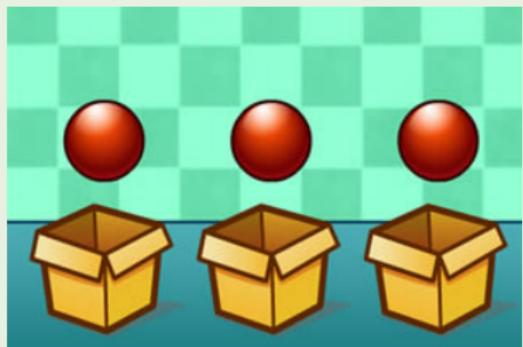
- **structure description** $\{n_i\}_{i=1}^r$

where $n_i := \sum_{j=1}^r \llbracket h_j = i \rrbracket$

- **rank description** $\{m_i\}_{i=0}^n$

where $m_i := \sum_{j=1}^r \llbracket n_j = i \rrbracket$

Obviously, $\sum_{i=1}^n i \cdot m_i = n$ and $\sum_{i=0}^n m_i = r$.



Indifference Principle

- A All state descriptions have equal weight.
- B All structure descriptions have equal weight.
- C Each nonempty subset of the alphabet is equally likely.
- D Each nonzero cardinality is equally likely.
- E All rank descriptions have equal weight.

Given n individuals, there are r^n possible state descriptions,

$$\left| \left\{ (n_1, \dots, n_r) : \sum_{i=1}^r n_i = n \right\} \right| = \binom{n+r-1}{r-1}$$

possible structure descriptions, and

$$p(n, r) := \left| \left\{ (m_0, \dots, m_n) : \sum_{i=1}^n i \cdot m_i = n \quad \& \quad \sum_{i=0}^n m_i = r \quad \& \quad \forall i: m_i \geq 0 \right\} \right|$$

possible rank descriptions.

(A) State Description ×

According to (A),

$$m^\dagger \left(\bigwedge_{i=1}^n Q_{h_i}(a_i) \right) = \frac{1}{r^n}$$

$$c^\dagger \left(Q_j(a_{n+1}) \middle| \bigwedge_{i=1}^n Q_{h_i}(a_i) \right) = \frac{m^\dagger \left(\bigwedge_{i=1}^n Q_{h_i}(a_i) \wedge Q_j(a_{n+1}) \right)}{m^\dagger \left(\bigwedge_{i=1}^n Q_{h_i}(a_i) \right)} = \frac{1}{r}$$

(B) Structure Description ✓

According to (B),

$$m^*(n_1, \dots, n_r) = \frac{1}{\binom{n+r-1}{r-1}}$$

Structure Description (n_1, \dots, n_r) corresponds to $\binom{n}{n_1, \dots, n_r}$ State Descriptions.

$$m^* \left(\bigwedge_{i=1}^n Q_{h_i}(a_i) \right) = \frac{m^*(n_1, \dots, n_r)}{\binom{n}{n_1, \dots, n_r}} = \frac{1}{\binom{n+r-1}{r-1} \binom{n}{n_1, \dots, n_r}}$$

Sometimes we write $m^*(h_{1:n}) := m^* \left(\bigwedge_{i=1}^n Q_{h_i}(a_i) \right)$ for short.

Carnap's Degree of Confirmation

Carnap's Degree of Confirmation

$$c^* \left(Q_j(a_{n+1}) \middle| \bigwedge_{i=1}^n Q_{h_i}(a_i) \right) = \frac{m^* \left(\bigwedge_{i=1}^n Q_{h_i}(a_i) \wedge Q_j(a_{n+1}) \right)}{m^* \left(\bigwedge_{i=1}^n Q_{h_i}(a_i) \right)} = \frac{\textcolor{red}{n_j} + 1}{\textcolor{red}{n} + r}$$

frequency — independent identical distribution(i.i.d)
extension

$$c^*(\varphi)$$

$$(C, D, E)$$

According to (C),

$$m^{\$}(h_{1:n}) = \frac{1}{\left(\sum_{i=1}^{\min\{r,n\}} \binom{r}{i} \right) \binom{n-1}{r-m_0-1} \binom{n}{n_1, \dots, n_r}}$$

According to (D),

$$m^{\#}(h_{1:n}) = \frac{1}{\min\{r, n\} \binom{r}{r-m_0} \binom{n-1}{r-m_0-1} \binom{n}{n_1, \dots, n_r}}$$

According to (E),

$$m^{\tau}(h_{1:n}) = \frac{1}{\binom{n}{n_1, \dots, n_r} \binom{r}{m_0, \dots, m_n} p(n, r)}$$

Rank Description (m_0, \dots, m_n) corresponds to $\binom{r}{m_0, \dots, m_n}$ Structure Descriptions.

What is the right w ?

- Constant Exchangeability Principle.

For any permutation σ of \mathbb{N}^+ ,

$$w(\psi(a_1, \dots, a_n)) = w(\psi(a_{\sigma(1)}, \dots, a_{\sigma(n)})) \quad (\text{Ex})$$

- Atom Exchangeability Principle.

For any permutation τ of $\{1, 2, \dots, r\}$,

$$w \left(\bigwedge_{i=1}^n Q_{h_i}(a_i) \right) = w \left(\bigwedge_{i=1}^n Q_{\tau(h_i)}(a_i) \right) \quad (\text{Ax})$$

- Sufficientness Postulate.

$$w \left(Q_j(a_{n+1}) \middle| \bigwedge_{i=1}^n Q_{h_i}(a_i) \right) = f_j(n_j, n) \quad (\text{SP})$$

What is the right w ?

Principle **Ex** asserts that $w \left(\bigwedge_{i=1}^n Q_{h_i}(a_i) \right)$ depends only on the vector $\langle n_{h_i} : 1 \leq i \leq n \rangle$, so that it is independent on the order of observing the individuals, while in the presence of **Ex**, principle **Ax** asserts that $w \left(\bigwedge_{i=1}^n Q_{h_i}(a_i) \right)$ depends only on $\{n_i : 1 \leq i \leq r\}$, and $w(Q_i(a_1)) = 1/r$ for all $1 \leq i \leq r$.

Carnap's λ -continuum

Theorem

Suppose language \mathcal{L} has at least two predicates i.e. $m \geq 2$, then the probability function w on \mathcal{L} satisfies **Ex**, **SP** iff $w = c_\lambda$ for some $0 \leq \lambda \leq \infty$.

Namely,

$$f_i(n_i, n) = \frac{n_i + \lambda \gamma_i}{n + \lambda}$$

where $\gamma_i = f_i(0, 0)$ and $\lambda = \frac{f_i(0, 1)}{f_i(0, 0) - f_i(0, 1)}$.

By adding **Ax**, $\forall i: \gamma_i = \frac{1}{r}$.

Shortcoming 1 — All Ravens are Black? ×

$$c^*(\forall x (R(x) \rightarrow B(x))) \leq \lim_{n \rightarrow \infty} \prod_{i=0}^{n-1} \frac{i+r-1}{i+r} = 0$$

Convergence Speed? Yes and No!

$$\prod_{n \geq 1} a_n = 0 \iff \sum_{n \geq 1} (1 - a_n) = \infty \quad \text{for } \forall n: 0 < a_n \leq 1$$

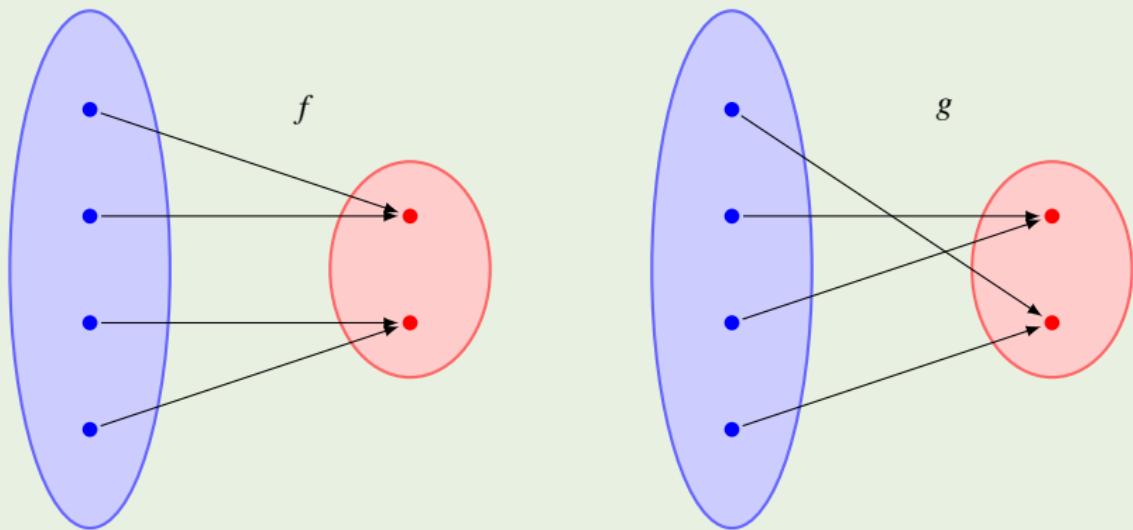
Shortcoming 2 — No-Free-Lunch for Carnap!

Strengthened “Hume” — Wolpert & Macready 1997, Igel & Toussaint 2004:

No-Free-Lunch Theorem!

All state descriptions with the same structure description have equal weight!

Block Uniform



No Free Lunch Theorem — Strengthened “Hume”

Theorem (No Free Lunch Theorem)

If and only if the probability distribution P is block uniform, i.e.

$$\forall f, g \in \mathcal{Y}^X : \forall y \in \mathcal{Y} \left(|f^{-1}(y)| = |g^{-1}(y)| \right) \implies P(f) = P(g)$$

then for any two algorithms A, A' , any value $k \in \mathbb{R}$, any $m \in \{1, \dots, |\mathcal{X}|\}$, and any performance measure L ,

$$\sum_{f \in \mathcal{Y}^X} P(f) [\![k = L(T_m^y(A, f))]\!] = \sum_{f \in \mathcal{Y}^X} P(f) [\![k = L(T_m^y(A', f))]\!]$$

where $T_n := \langle (x_1, f(x_1)), \dots, (x_n, f(x_n)) \rangle$, $T_n^x := \langle x_1, \dots, x_n \rangle$, $T_n^y := \langle f(x_1), \dots, f(x_n) \rangle$ and $A : T_n \mapsto x_{n+1} \in \mathcal{X} \setminus T_n^x$.

equally well and equally poorly
No learning is possible without some prior knowledge!

No-Free-Lunch for Pure Inductive Logic

- No learning is possible for c^\dagger . It seems possible to learn with c^* and c_λ . Unfortunately, No-Free-Lunch Theorem!
- Take $\mathcal{X} := \mathcal{C}$, $\mathcal{Y} := \mathcal{Q}$ in the No-Free-Lunch Theorem. The state description $h: \mathcal{C} \rightarrow \mathcal{Q}$ can be taken as a *classification* function.
- Then “all state descriptions $h: \mathcal{C} \rightarrow \mathcal{Q}$ with the same structure description $\{n_i : 1 \leq i \leq r\}$ have equal weight” is **block uniform**!
- Define induction algorithm $A(h_{1:n}) := \operatorname{argmax}_j c^*(Q_j(a_{n+1})|h_{1:n})$, and loss function $L(A, n, h) := \llbracket A(h_{1:n}) \neq h_{n+1} \rrbracket$. Then for any A' ,

$$\sum_{h \in \mathcal{Y}^{\mathcal{X}}} m^*(h_{1:n}) L(A, n, h) = \sum_{h \in \mathcal{Y}^{\mathcal{X}}} m^*(h_{1:n}) L(A', n, h)$$

- Similarly for rank description (E), which is related to Good-Turing estimate. And similarly for Ristad's methods (C)(D).

Time Series and Solomonoff Induction

- However, if we take time into consideration, define

$$M(h_{1:n}) := \sum_{p: U(p) = h_{1:n}*} 2^{-\ell(p)}$$

then we can get free lunch with M' , since M' biases non-random state descriptions and is not block uniform. where

$$M'(\epsilon) := 1$$

$$M'(h_{1:n}) := M'(h_{<n}) \frac{M(h_{1:n})}{\sum_{Q \in \mathcal{Q}} M(h_{<n}Q)}$$

and we can extend M' from state descriptions to sentences.

- Besides, by c^* the probability of “the sun will rise tomorrow” is $\frac{n+1}{n+2}$, but c^* fails to confirm “the sun will always rise”, while M' can confirm it.

$$\lim_{n \rightarrow \infty} M'(\text{the sun will always rise} | \text{the sun rises in the first } n \text{ days}) = 1$$

PAC (Probably Approximately Correct) Learning

Definition (PAC-Learnability)

A hypothesis space $\mathcal{H} \subset 2^{\mathcal{X}}$ is PAC-learnable if there exists a sample complexity function $m_{\mathcal{H}}: (0, 1)^2 \rightarrow \mathbb{N}$ and a learning algorithm A with the following property:

- for every $\varepsilon, \delta \in (0, 1)$
- for every distribution \mathcal{D} over \mathcal{X} , and for every labeling function $f: \mathcal{X} \rightarrow \{0, 1\}$

when running A on $m \geq m_{\mathcal{H}}(\varepsilon, \delta)$ i.i.d. training samples S generated by \mathcal{D} and labeled by f , the algorithm A returns a hypothesis $A(S) \in \mathcal{H}$ s.t.

$$P_{S \sim \mathcal{D}^m} \left(P_{x \sim \mathcal{D}} (A(S)(x) \neq f(x)) > \varepsilon \right) < \delta$$

$$\text{No-Free-Lunch} \implies m_{2^{\mathcal{X}}}(\varepsilon, \delta) = \infty$$

Agnostic PAC Learning

Definition (Agnostic PAC-Learnability)

A hypothesis space $\mathcal{H} \subset \mathcal{Y}^X$ is agnostic PAC-learnable under a class Δ of distributions with respect to $\mathcal{Z} := X \times \mathcal{Y}$ and a loss function $\ell: \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}^+$, if there exists a sample complexity function $m_{\mathcal{H}}: (0, 1)^2 \rightarrow \mathbb{N}$ and a learning algorithm A with the following property:

- for every $\varepsilon, \delta \in (0, 1)$
- for every distribution $\mathcal{D} \in \Delta$ over \mathcal{Z}

when running A on $m \geq m_{\mathcal{H}}(\varepsilon, \delta)$ i.i.d. training samples S generated by \mathcal{D} , the algorithm A returns a hypothesis $A(S) \in \mathcal{H}$ s.t.

$$P_{S \sim \mathcal{D}^m} \left(L_{\mathcal{D}}(A(S)) - \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) > \varepsilon \right) < \delta$$

where $L_{\mathcal{D}}(h) := \mathbb{E}_{z \sim \mathcal{D}} [\ell(h, z)]$.

PAC-learnable under $\Delta := \{\mathcal{D}: \mathcal{D} \stackrel{\times}{\leq} \xi\} \iff$ PAC-learnable under $\xi(x) := \sum_{v \in \mathcal{M}} 2^{-K(v)} v(x)$.

VC-Dimension

Definition (Shattering)

A hypothesis space $\mathcal{H} \subset 2^{\mathcal{X}}$ shatters a set $C \subset \mathcal{X}$ if $\mathcal{H}|_C = 2^C$.

Definition (VC-Dimension)

$\text{VC}(\mathcal{H}) := \sup \{|C| : \mathcal{H} \text{ shatters } C\}$.

If someone can explain every phenomena, her explanations are worthless.

Theorem

If $\text{VC}(\mathcal{H}) = \infty$, then \mathcal{H} is not PAC-learnable.

Theorem (Fundamental Theorem of PAC Learning)

\mathcal{H} is PAC-learnable iff $\text{VC}(\mathcal{H}) < \infty$. Indeed, there exists C_1, C_2 s.t.

$$C_1 \frac{\text{VC}(\mathcal{H}) + \log(1/\delta)}{\varepsilon} \leq m_{\mathcal{H}}(\varepsilon, \delta) \leq C_2 \frac{\text{VC}(\mathcal{H}) \log(1/\varepsilon) + \log(1/\delta)}{\varepsilon}$$

Contents

- ① History of AI
- ② Philosophy of Induction
- ③ Inductive Logic
- ④ Universal Induction
- ⑤ Reinforcement Learning
- ⑥ General Reinforcement Learning

Formal Learning Theory

- Carnap's inductive logic is a design for a '**learning machine**' that can extrapolate certain kinds of empirical regularities from the data with which it is supplied, and the task of inductive logic is to construct a '**universal learning machine**'.
- If there is such a thing as a correct definition of '**degree of confirmation**' which can be fixed once and for all, then a machine that predicts in accordance with it would be a **cleverest possible learning machine**.
- Either there are better and better '**degree of confirmation**' functions, but no '**best possible**', or there is a '**best possible**' but it is not computable by a machine.

Formal Learning Theory

- Putnam 1963
- Gold 1967 $\hat{f}(n+1) = g(\langle f(0), \dots, f(n) \rangle)$ $\varphi_{\lim_{n \rightarrow \infty} g(\langle f(0), \dots, f(n) \rangle)} = f$
- Solomonoff 1960

Ray Solomonoff

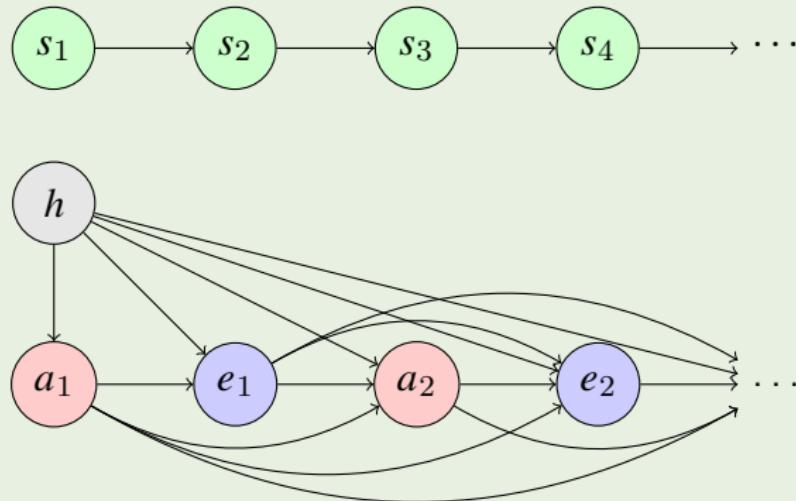
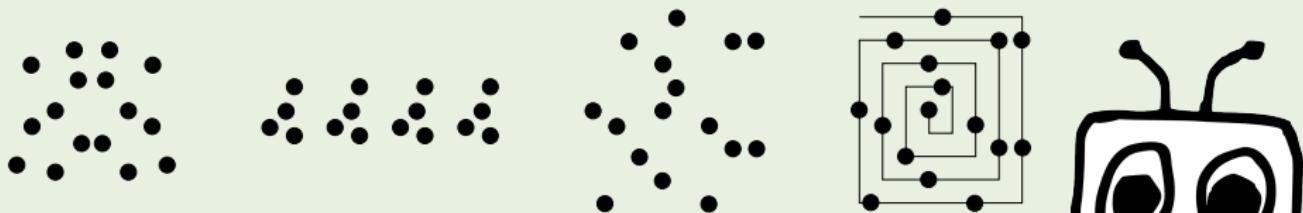


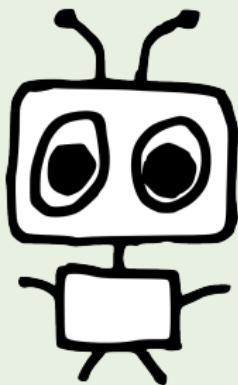
Figure: Solomonoff



- 0101010101010101010101010101
- 01010011011110110111010101000001
- 00100100001111110110101010001000

$\{\pi\}$

- ① What is regularity/pattern/law/principle/model/hypothesis/theory?
- ② What is phenomenon/data/experience?
- ③ What is randomness/noise?
- ④ What is typicalness/unpredictability/incompressibility?
- ⑤ What is simplicity/complexity?
- ⑥ What is learning?
- ⑦ What is beauty/interesting/curiosity/novelty/surprise/creativity?
- ⑧ What is intelligence?



Text \implies Meaning?

numeral $\xrightarrow{\text{short algorithm}}$ value

汝有田舍翁，家资殷盛，而累世不识“之”“乎”。一岁，聘楚士训其子。楚士始训之搦管临朱，书一画，训曰：“一”字。书二画，训曰：“二”字。书三画，训曰：“三”字。其子辄欣欣然掷笔，归告其父曰：“儿得矣！儿得矣！可无烦先生，重费馆谷也，请谢去。”其父喜从之，具币谢遣楚士。

逾时，其父拟征召姻友万氏者饮，令子晨起治状，久之不成。父趣之。其子恚曰：“天下姓字多矣，奈何姓万？自晨起至今，才完五百画也。”

Kolmogorov Complexity

Definition (Kolmogorov Complexity)

$$K(x) := \min_p \{\ell(p) : U(p) = x\}$$

where U is a universal prefix Turing machine.

How to quantify
“simplicity”?
“randomness”?

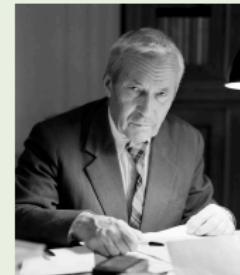
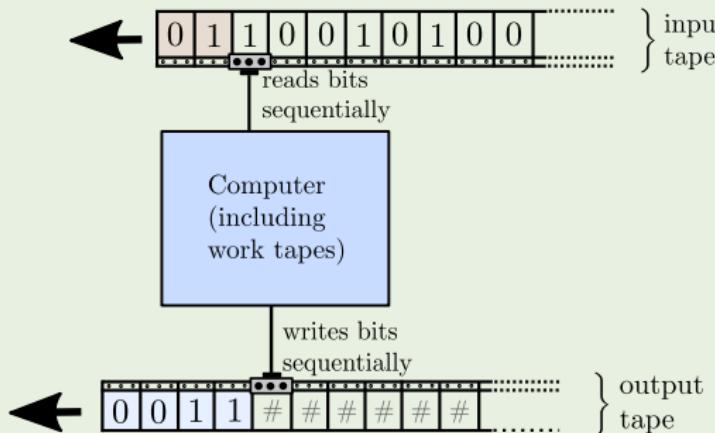


Figure: Kolmogorov

“Independence” of universal Turing machine!

Kraft-Chaitin Theorem

Theorem (Kraft Inequality)

- ① Let $f: \mathcal{X} \rightarrow 2^{<\omega}$ be uniquely decodable. Let $\ell_x := \ell(f(x))$. Then f satisfies the Kraft inequality

$$\sum_{x \in \mathcal{X}} 2^{-\ell_x} \leq 1$$

- ② Conversely, for any set of code length $\{\ell_x: x \in \mathcal{X}\}$ satisfying the above Kraft inequality, there exists a prefix code f such that $\ell_x = \ell(f(x))$.

Theorem (Kraft-Chaitin Theorem)

For any c.e. set $(\ell_i, x_i)_{i \in \omega} \subset \mathbb{N} \times 2^{<\omega}$ with $\sum_{i \in \omega} 2^{-\ell_i} \leq 1$, one can effectively obtain a prefix-free machine M and strings p_i of length ℓ_i such that $M(p_i) = x_i$ for all i and $\text{dom}(M) = \{p_i: i \in \omega\}$.

Properties

- ① $K(n) \stackrel{+}{\leq} \log^* n \leq \log n + 2 \log \log n$
- ② $K(x) \stackrel{+}{\leq} K(x|\ell(x)) + K(\ell(x)) \stackrel{+}{\leq} \ell(x) + \log^* \ell(x) \leq \ell(x) + 2 \log \ell(x)$
- ③ $\sum_x 2^{-K(x)} \leq 1$
- ④ $K(x|y) \stackrel{+}{\leq} K(x) \stackrel{+}{\leq} K(x, y)$
- ⑤ $K(xy) \stackrel{+}{\leq} K(x, y) \stackrel{+}{\leq} K(x) + K(x|y) \stackrel{+}{\leq} K(x) + K(y)$
- ⑥ $K(x) \stackrel{+}{=} K(x, K(x))$
- ⑦ $K(x|y, K(y)) + K(y) \stackrel{+}{=} K(x, y) \stackrel{+}{=} K(y, x) \stackrel{+}{=} K(y|x, K(x)) + K(x)$
- ⑧ $K(f(x)) \stackrel{+}{\leq} K(x) + K(f)$ for computable f
- ⑨ $K(x) \stackrel{+}{\leq} -\log \mu(x) + K(\mu)$ if μ is lower semicomputable and $\sum_x \mu(x) \leq 1$
- ⑩ $\sum_{x:f(x)=y} 2^{-K(x)} \stackrel{+}{\leq} 2^{-K(y)}$ if f is computable and $K(f) = O(1)$
- ⑪ $0 \leq \mathbb{E}_\mu[K] - H(\mu) \stackrel{+}{\leq} K(\mu)$ for computable probability distribution μ

Universal Similarity Metric

$$\begin{aligned} d(x, y) &:= \frac{\max\{K(x|y), K(y|x)\}}{\max\{K(x), K(y)\}} \\ &\approx \frac{K_T(xy) - \min\{K_T(x), K_T(y)\}}{\max\{K_T(x), K_T(y)\}} \end{aligned}$$

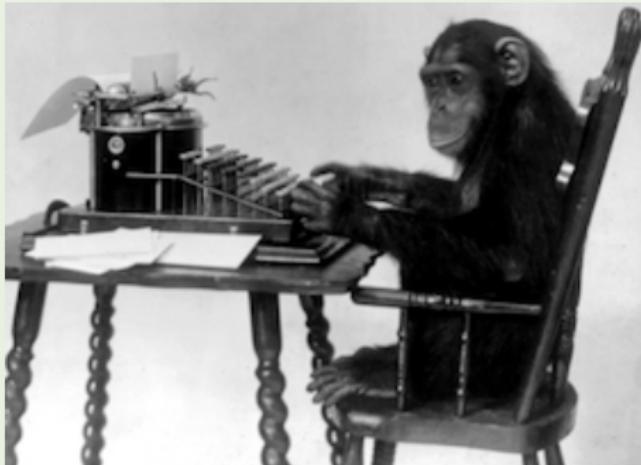
- T : Lempel-Ziv/gzip/bzip2/PPMZ, or
 $K_T(x) := -\log P_{\text{google}}(x)$
- compute similarity matrix $(d(x_i, x_j))_{ij}$
- cluster similar objects.

Algorithmic Probability

Definition (Algorithmic Probability)

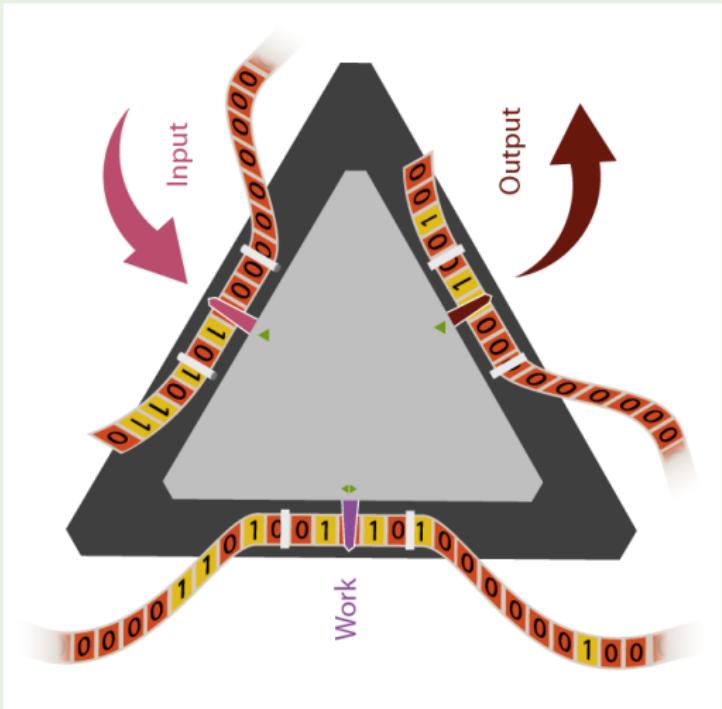
$$M(x) := \sum_{p:U(p)=x^*} 2^{-\ell(p)}$$

where U is a universal monotone Turing machine.



$$\sum_{p:U(p)=x^*} 2^{-\ell(p)} \gg 2^{-\ell(x)}$$

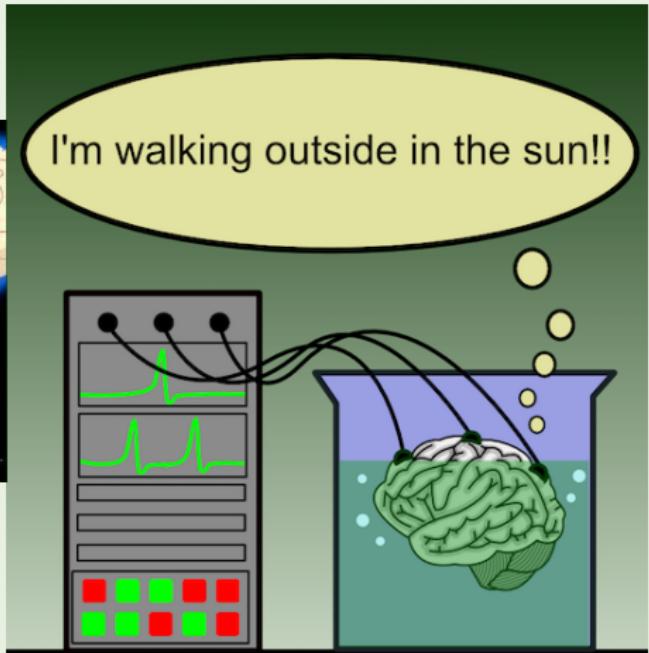
Toss Coin onto Universal Turing Machine



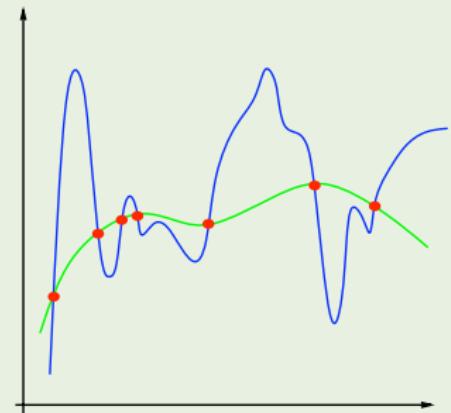
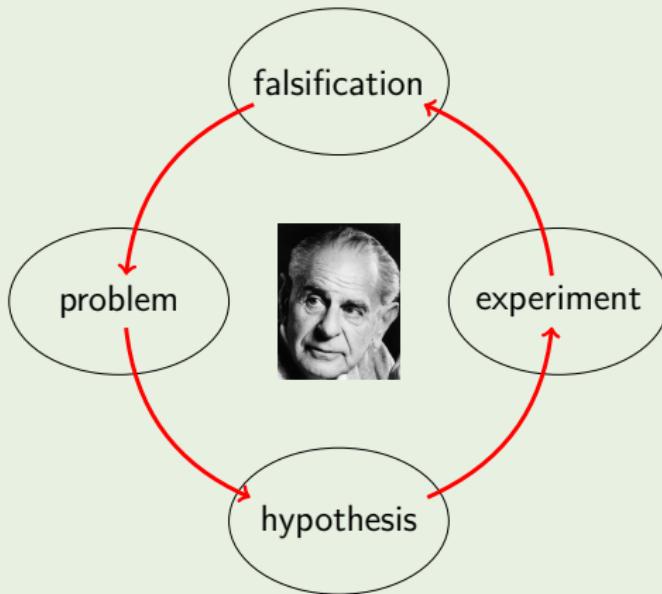
Boltzmann Brain vs Brain in a Vat



$K(v)$ is too large.



Aspect 1 — Popper



\mathcal{H} : truth \leftarrow simplicity/generality/aesthetic/utilitarian/...

Make a weighted prediction based on all consistent programs, with short programs weighted higher.

Aspect 2 — Deterministic vs Stochastic

$\mathcal{M} := \{\nu_1, \nu_2, \dots\}$ lower semicomputable semi-measure.

$$\xi(x) := \sum_{\nu \in \mathcal{M}} 2^{-K(\nu)} \nu(x)$$

$w_\nu := 2^{-K(\nu)}$ is reparametrization & regrouping invariant.

$$M(x) \stackrel{x}{=} \xi(x)$$

Aspect 3 — Frequency Interpretation

$$\begin{aligned} M(x) &= \sum_p 2^{-\ell(p)} \llbracket U(p) = x^* \rrbracket \\ &= \lim_{n \rightarrow \infty} \frac{\sum_{p: \ell(p) \leq n} 2^{n-\ell(p)} \llbracket U(p) = x^* \rrbracket}{2^n} \\ &\approx \lim_{n \rightarrow \infty} \frac{|\{p : \ell(p) = n \text{ \& } U(p) = x^*\}|}{2^n} \end{aligned}$$

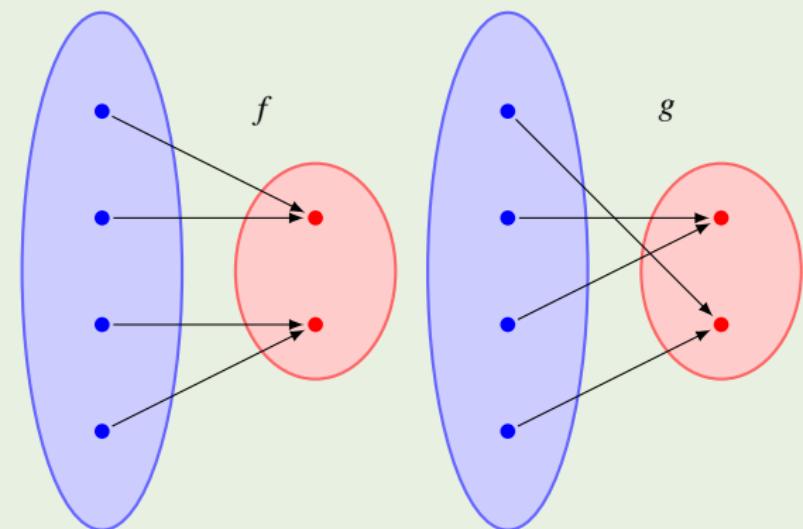
algorithmic probability = $\frac{|\text{consistent worlds}|}{|\text{all possible worlds}|}$

{ Carnap — frequency of phenomena — i.i.d
Solomonoff — frequency of causes — arbitrary order Markov chain

Aspect 4 — Free-Lunch for Solomonoff!

$$\xi \rightarrow \mu$$

break “block uniform” — bias non-random functions



Algorithmic Coding Theorem

$$m(x) := \sum_{p:U(p)\downarrow=x} 2^{-\ell(p)}$$

Theorem (Algorithmic Coding Theorem)

$$K(x) \stackrel{+}{=} -\log m(x)$$

The probability of a string being produced by a random algorithm is inversely proportional to its algorithmic complexity.

If a string has many long descriptions then it also has a short description.

$$KM(x) := -\log M(x)$$

$$KM(x) \leq Km(x) \stackrel{+}{\leq} KM(x) + K(KM(x))$$

Completeness Theorem

$$M'(\epsilon) := 1$$

$$M'(x_{1:t}) := M'(x_{<t}) \frac{M(x_{1:t})}{\sum_{x \in \mathcal{X}} M(x_{<t}x)} = \frac{M(x_{1:t})}{M(\epsilon)} \prod_{i=1}^t \frac{M(x_{<i})}{\sum_{x \in \mathcal{X}} M(x_{<i}x)}$$

Theorem (Completeness Theorem)

For any computable measure μ ,

$$\sum_{t=1}^{\infty} \sum_{x_{1:t} \in \mathcal{X}^t} \mu(x_{<t}) \left(M'(x_t|x_{<t}) - \mu(x_t|x_{<t}) \right)^2 \leq D(\mu||M) \stackrel{+}{\leq} K(\mu) \ln 2$$

Emergence of Simple Laws of Physics

Completeness Theorem

For any computable measure μ , there is a set $A \subset \mathcal{X}^*$ with $\mu(A) = 1$ s.t. for all $\mathbf{x} \in A$

$$\sum_{y \in \mathcal{X}} \left(\sqrt{M'(y|\mathbf{x})} - \sqrt{\mu(y|\mathbf{x})} \right)^2 \xrightarrow{n \rightarrow \infty} 0$$

Emergence of Simple Laws of Physics

For any computable measure μ ,

$$M' \left\{ \sum_{y \in \mathcal{X}} \left(\sqrt{M'(y|\mathbf{x})} - \sqrt{\mu(y|\mathbf{x})} \right)^2 \xrightarrow{n \rightarrow \infty} 0 \right\} \geq 2^{-K(\mu)}$$



Universal Prediction of Selected Bits

Theorem (Universal Prediction of Selected Bits)

Let $f: \{0, 1\}^* \rightarrow \{0, 1, \epsilon\}$ be a total recursive function and $x \in 2^\omega$ satisfying $f(x_{<n}) = x_n$ whenever $f(x_{<n}) \neq \epsilon$. If $f(x_{ for an infinite sequence n_1, n_2, \dots then$

$$\lim_{i \rightarrow \infty} M'(x_{n_i} | x_{<n_i}) = 1$$

Pure Universal Inductive Logic?

- $M(\Theta(a_{1:n})) := \sum_{p:U(p)=h_{1:n}*} 2^{-\ell(p)}$
where $\Theta(a_{1:n}) := \bigwedge_{i=1}^n Q_{h_i}(a_i)$.

- $M'(\varphi(\vec{a})) := \sum_{\Theta(\vec{b}) \models \psi(\vec{a})} M'(\Theta(\vec{b}))$
where $\models \varphi(\vec{a}) \leftrightarrow \bigvee_{\Theta(\vec{b}) \models \varphi(\vec{a})} \Theta(\vec{b})$.

FDNF

$$\sum_{t=1}^{\infty} \sum_{\varphi(a_{1:t})} \mu(\varphi(a_{<t})) \left(M'(\varphi(a_t) | \varphi(a_{<t})) - \mu(\varphi(a_t) | \varphi(a_{<t})) \right)^2 \stackrel{+}{\leq} K(\mu) \ln 2$$

where $\varphi(a_{1:t}) := \bigwedge_{i=1}^t \varphi(a_i/x)$.

All Ravens are Black! ✓

Theorem (All Ravens are Black)

$$\lim_{n \rightarrow \infty} M' \left(\forall x (R(x) \rightarrow B(x)) \middle| \bigwedge_{i=1}^n (\neg R(a_i) \vee B(a_i)) \right) = 1$$

Theorem (Confirmation by Random Sampling)

If the sampling function $t: \mathbb{N} \rightarrow \mathbb{N}$ satisfies $\forall i: t_i \leq t_{i+1}$ and $\chi_{1:\infty}$ is algorithmic random, where $\chi_i := [\![\exists k (t_k = i)]\!]$, then

$$M' \left(\forall x \varphi(x) \middle| \bigwedge_{i=1}^n \varphi(a_{t_i}) \right) \xrightarrow{n \rightarrow \infty} 1$$

$$M(1|1^n) \xrightarrow{n \rightarrow \infty} 1 \quad M(0|1^n) \stackrel{\cong}{=} 2^{-K(n)} \quad \sum_{n=0}^{\infty} M(0|1^n) < \infty$$

Why Solomonoff Prior?

$$H(w) \leq \mathbb{E}_w [K] \leq H(w) + K(w)$$

Maximum Entropy + Occam's Razor

$$\underset{\substack{w \models \sum_{\nu \in \mathcal{M}} w_\nu = 1}}{\minimize} \frac{\mathbb{E}_w [K]}{H(w)}$$



$$w_\nu^* = \frac{2^{-K(\nu)}}{\sum_{\nu \in \mathcal{M}} 2^{-K(\nu)}}$$

Advantages & Disadvantages

- free-lunch
- universality — finite error
- data sparse problem — arbitrary order Markov chain — universal smoothing method
- confirmation of $\forall x: R(x) \rightarrow B(x)$
- incomputability
- weakly depends on universal Turing machine

A statistical mechanical interpretation of AIT

an energy eigenstate $n \Rightarrow$ a program p s.t. $U(p) \downarrow$
 the energy E_n of $n \Rightarrow$ the length $\ell(p)$ of p
 Boltzmann constant $k \Rightarrow 1/\ln 2$

$$Z = \sum_n e^{-\frac{E_n}{kT}} \Rightarrow Z = \sum_{p: U(p) \downarrow} 2^{-\frac{\ell(p)}{T}} \quad \text{Partition function}$$

$$F = -kT \ln Z \Rightarrow F = -T \log Z \quad \text{Free energy}$$

$$P(n) = \frac{1}{Z} e^{-\frac{E_n}{kT}} \Rightarrow P(p) = \frac{1}{Z} 2^{-\frac{\ell(p)}{T}} \quad \text{Boltzmann distribution}$$

$$E = \sum_n P(n) E_n \Rightarrow E = \sum_{p: U(p) \downarrow} P(p) \ell(p) \quad \text{Energy}$$

$$S = \frac{E - F}{T} \Rightarrow S = \frac{E - F}{T} = H(P) \quad \text{Entropy}$$

$$T \geq 1 \Rightarrow H(P) = \infty \quad 0 < T < 1 \Rightarrow H(P) < \infty$$

Temperature = Compression Rate

$$T = \lim_{n \rightarrow \infty} \frac{K(Z_{1:n})}{n} = \lim_{n \rightarrow \infty} \frac{K(F_{1:n})}{n} = \lim_{n \rightarrow \infty} \frac{K(E_{1:n})}{n} = \lim_{n \rightarrow \infty} \frac{K(S_{1:n})}{n}$$

Fixpoint theorem: for $T \in (0, 1)$, if Z is computable, then $\lim_{n \rightarrow \infty} \frac{K(T_{1:n})}{n} = T$.

Computable Universal Predictor

$$Z_h = \sum_{p:U(p)=h*} 2^{-\frac{\ell(p)}{T}}$$

$$Z(h) = \sum_{p:U(p)=h*} 2^{-\frac{\ell(p)+\log t(p,h)}{T}}$$

T is computable $\implies Z(h)$ is computable

$$\sum_{t=1}^{\infty} |1 - Z(h_t|h_{<t})| \leq \frac{Km(h) \ln 2 + \ln t(p,h)}{T}$$

Stochastic Case

$$E_{\{h, v\}} = -\log w_\epsilon^\nu - \log v(h)$$

$$Z_h = \sum_{v \in \mathcal{M}} 2^{-E_{\{h, v\}}} = \sum_{v \in \mathcal{M}} w_\epsilon^\nu v(h) = \xi(h)$$

$$P_h(v) = \frac{2^{-E_{\{h, v\}}}}{Z_h} = \frac{w_\epsilon^\nu v(h)}{\xi(h)} = w_h^\nu$$

$$F_h = -\log Z_h = \underbrace{-\log \xi(h)}_{\approx K(h)} + \underbrace{\mathbb{E}_{w_h}[-\log v(h)]}_{noise} + \underbrace{D(w_h \| w_\epsilon)}_{surprise}$$

Deduction vs Induction

	Induction		Deduction
Type of inference	generalization/prediction	\Leftrightarrow	specialization/derivation
Framework	probability axioms	$\hat{\equiv}$	logical axioms
Assumptions	prior	$\hat{\equiv}$	non-logical axioms
Inference rule	Bayes rule	$\hat{\equiv}$	modus ponens
Results	posterior	$\hat{\equiv}$	theorems
Universal scheme	Solomonoff probability	$\hat{\equiv}$	ZFC
Universal inference	universal induction	$\hat{\equiv}$	universal theorem prover
Limitation	uncomputable (Turing)	$\hat{\equiv}$	imcomplete (Gödel)
In practice	approximations	$\hat{\equiv}$	semi-formal proofs
Operation	computation	$\hat{\equiv}$	proof

Logic vs Statistics

Field	Logical Approach	Statistical Approach
Knowledge representation	First-order logic	Graphical models
Automated reasoning	Satisfiability testing	Markov chain Monte Carlo
Machine learning	Inductive logic programming	Neural networks
Planning	Classical planning	Markov decision processes
Natural language processing	Definite clause grammars	Probabilistic context-free grammars

Logic	Statistics
rule-based	data-driven
rigour	possibility
knowable	black-box
simple & perfect world	complex & uncertain world
✗	✓

Prediction with Expert Advice

Assume that there is some large, possibly infinite, class of ‘experts’ which make predictions. The aim is to observe how each of these experts perform and develop independent predictions based on this performance.

- Follow the (perturbed) leader.
- Predicts according to a majority vote by the “good” experts.
- Multiplicative Weights. — take expert which performed best in past with high probability and others with smaller probability.
- Regularization. Choose the class of all computable experts, and penalize “complex” experts.
- Universal Portfolios.

Universal Portfolios

- the agent chooses a distribution $\mathbf{b}_t \in \Delta_n := \{x \in [0, 1]^n : \|x\|_1 = 1\}$ of wealth over n goods.
- nature chooses returns $\mathbf{x}_t \in (\mathbb{R}^+)^n$, where

$$(\mathbf{x}_t)_i = \frac{\text{price of good } i \text{ at end of } t}{\text{price of good } i \text{ at beginning of } t}$$

- the total wealth. $W_t(\mathbf{b}, \mathbf{x}) = W_1 \prod_{k=1}^t \mathbf{b}_k^\top (\mathbf{x}_{<k}) \mathbf{x}_k$
- regret. $R_t := \max_{\mathbf{b} \in \Delta_n} \sum_{k=1}^t \log \mathbf{b}_k^\top (\mathbf{x}_{<k}) \mathbf{x}_k - \sum_{k=1}^t \log \hat{\mathbf{b}}_k^\top (\mathbf{x}_{<k}) \mathbf{x}_k$
- universal portfolios.

$$\hat{\mathbf{b}}_1 := \left(\frac{1}{n}, \dots, \frac{1}{n} \right)$$

$$\hat{\mathbf{b}}_{t+1}(\mathbf{x}_{1:t}) := \frac{\int_{\Delta_n} \mathbf{b} W_t(\mathbf{b}, \mathbf{x}) d\mathbf{b}}{\int_{\Delta_n} W_t(\mathbf{b}, \mathbf{x}) d\mathbf{b}}$$

- Asymptotic Optimality. $\frac{1}{t} \log W_t(\hat{\mathbf{b}}, \mathbf{x}) \xrightarrow{t \rightarrow \infty} \frac{1}{t} \log \max_{\mathbf{b} \in \Delta_n} W_t(\mathbf{b}, \mathbf{x})$

Randomness

- ① **Typicalness** (The statistician's approach): A random sequence is the typical outcome of a random variable. Random sequences should not have effectively rare distinguishing properties.
- ② **Incompressibility** (The coder's approach): Rare patterns can be used to compress information. Random sequences should not be effectively described by a significantly shorter description than their literal representation.
- ③ **Unpredictability** (The gambler's approach): A betting strategy can exploit rare patterns. Random sequences should be unpredictable. No effective martingale can make an infinite amount betting on the bits.

The Statistician's Approach

- A random sequence should be absolutely normal.
- If you select a subsequence, then it should satisfy the law of large numbers, the law of the iterated logarithm. . .
- But what selection functions should be allowed? Computable?
- Martin-Löf: we can effectively test whether a particular infinite sequence does not satisfy a particular law of randomness by effectively testing whether the law is violated on increasingly long initial segments. We should consider the intersection of all sets of measure one with recursively enumerable complements. (Such a complement set is expressed as the union of a recursively enumerable set of cylinders).

Cantor Space 2^ω

- For $x \in 2^{<\omega}$, the cylinder set $\Gamma_x := \{y \in 2^\omega : x < y\}$ is the basic open set. It corresponds to the interval $[0.x, 0.x + 2^{-\ell(x)})$.
- For $A \subset 2^{<\omega}$, the open set generated by A is $\Gamma_A := \bigcup_{x \in A} \Gamma_x$.
- The Lebesgue measure $\mu(\Gamma_x) := 2^{-\ell(x)}$, $\mu(x) := \mu(\Gamma_x)$.
- The outer measure of $C \subset 2^\omega$ is $\mu^*(C) := \inf \left\{ \sum_{x \in A} 2^{-\ell(x)} : C \subset \Gamma_A \right\}$.
- The inner measure of C is $\mu_*(C) := 1 - \mu^*(2^\omega \setminus C)$.
- If C is measurable, then $\mu^*(C) = \mu_*(C)$.
- $A \subset 2^\omega$ has measure 0 iff there is a sequence $\{V_n\}_{n \in \omega}$ of open sets s.t. $A \subset \bigcap_{n \in \omega} V_n$ and $\lim_{n \rightarrow \infty} \mu(V_n) = 0$.

Martin-Löf Randomness

Definition (Martin-Löf Randomness)

- A total lower semicomputable function $\delta: 2^{<\omega} \rightarrow \omega$ is a Martin-Löf test if $\forall n: \mu(V_n) \leq 2^{-n}$, where $V_n := \{x: \delta(x) \geq n\}$.
- $x \in 2^\omega$ is ML-random if for every ML-test δ , $\sup_n \delta(x_{1:n}) < \infty$.

$$\delta(x) < \infty \iff x \notin \bigcap_{n=1}^{\infty} V_n$$

Definition (Martin-Löf Randomness)

- A Martin-Löf test is a uniformly c.e. (i.e., $\{\langle n, x \rangle : x \in V_n\}$ is c.e.) sequence of open sets $\{V_n\}_{n \in \omega}$ s.t. $\forall n: \mu(V_n) \leq 2^{-n}$.
- $x \in 2^\omega$ is ML-random if for every ML-test $\{V_n\}_{n \in \omega}$, $x \notin \bigcap_{n=1}^{\infty} V_n$.

Martin-Löf Randomness

Definition (Universal Martin-Löf Test)

A ML-test δ_0 is *universal* if for every ML-test δ , $\exists c \forall x: \delta_0(x) \geq \delta(x) - c$.

A ML-test $\{U_n\}_{n \in \omega}$ is *universal* iff for every ML-test $\{V_n\}_{n \in \omega}$,

$$\bigcap_{n \in \omega} U_n \supset \bigcap_{n \in \omega} V_n.$$

$\delta(x) := \ell(x) - K(x|\ell(x))$ is a universal ML-test.

$$R_b := \{x \in 2^\omega : \exists n (K(x_{1:n}) < n - b)\}$$

$\{R_b\}_{b \in \omega}$ is a universal ML-test.

Theorem (Schnorr 1973)

A sequence $x \in 2^\omega$ is ML-random iff it is 1-random.

The Gambler's Approach

- A martingale is a function $d: 2^{<\omega} \rightarrow [0, \infty)$ s.t. for every $\sigma \in 2^{<\omega}$

$$d(\sigma) = \frac{d(\sigma0) + d(\sigma1)}{2}$$

- A supermartingale is a function $d: 2^{<\omega} \rightarrow [0, \infty)$ s.t. for every $\sigma \in 2^{<\omega}$

$$d(\sigma) \geq \frac{d(\sigma0) + d(\sigma1)}{2}$$

- A (super)martingale d succeeds on $x \in 2^\omega$ if $\limsup_{n \rightarrow \infty} d(x_{1:n}) = \infty$.

Theorem

A sequence $x \in 2^\omega$ is ML-random iff no c.e. (super)martingale succeeds on it.

The Coder's Approach

Theorem

The following are equivalent.

- $x \in 2^\omega$ is ML-random.
- No c.e. (super)martingale succeeds on it.
- $\exists c \forall n: K(x_{1:n}) \geq n - c$
- $\forall n: Km(x_{1:n}) \stackrel{+}{=} n$
- $\lim_{n \rightarrow \infty} K(x_{1:n}) - n = \infty$
- $\sum_{n=1}^{\infty} 2^{n-K(x_{1:n})} < \infty$
- $\sup_n 2^{n-K(x_{1:n})} < \infty$
- $C(x_{1:n}) \stackrel{+}{\geq} n - K(n)$
- $C(x_{1:n}) \stackrel{+}{\geq} n - f(n)$ for every computable f s.t. $\sum_{n=1}^{\infty} 2^{-f(n)} < \infty$.

Definition (1-Randomness)

$x \in 2^\omega$ is 1-random if

$$\exists c \forall n: K(x_{1:n}) \geq n - c$$

Definition (Solovay Reducibility)

Let $a_n \rightarrow \alpha$ and $b_n \rightarrow \beta$ be two computable strictly increasing sequences of rationals converging to lower semicomputable reals α and β . We say that $\alpha \leq_S \beta$ if there is a constant c and a total computable function f s.t. $\forall n: \alpha - a_{f(n)} \leq c(\beta - b_n)$.

Theorem

For lower semicomputable reals α , the following are equivalent.

- α is 1-random
- $\alpha \geq_S \beta$ for all lower semicomputable reals β .
- $\alpha \geq_S \Omega$
- $K(\alpha_{1:n}) \stackrel{+}{\geq} K(\beta_{1:n})$ for all lower semicomputable reals β .
- $K(\alpha_{1:n}) \stackrel{+}{\geq} K(\Omega_{1:n})$
- $K(\alpha_{1:n}) \stackrel{+}{=} K(\Omega_{1:n})$
- $\alpha = \Omega_U := \sum_{p:U(p)\downarrow} 2^{-\ell(p)}$ for some universal prefix Turing machine U .

μ/ξ -randomness

Definition (μ/ξ -randomness)

- A sequence $x \in 2^\omega$ is μ/ξ -random if $\exists c \forall n: \xi(x_{1:n}) \leq c \cdot \mu(x_{1:n})$.
- A sequence $x \in 2^\omega$ is μ -ML-random if $\exists c \forall n: M(x_{1:n}) \leq c \cdot \mu(x_{1:n})$.
- $x_{1:\infty}$ is μ -ML-random iff $\sup_n \delta(x_{1:n} | \mu) < \infty$, where

$$\delta(x | \mu) := \log \frac{M(x)}{\mu(x)}$$

- For a computable μ , $x_{1:\infty}$ is μ -ML-random iff

$$\forall n: Km(x_{1:n}) \stackrel{+}{=} -\log \mu(x_{1:n})$$

Randomness, Triviality, Logical Depth

- $A \subset \mathbb{N}$ is *low* if $A' \leq_T \emptyset'$, and A is *high* if $\emptyset'' \leq_T A'$.
- A is *low for ML-randomness* if each ML-random set is already ML-random relative to A .
- A is *low for K* if $\exists c \forall x: K(x) \leq K^A(x) + c$.
- $x \in 2^\omega$ is *K -trivial* if $\exists c \forall n: K(x_{1:n}) \leq K(n) + c$.

Theorem

A is K -trivial \iff A is low for ML-randomness \iff A is low for K .

Some sequences are K -trivial but not computable.

Neither randoms, nor K -trivials, are deep.

Definition (Logical Depth)

The logical depth of x at a significance level b is

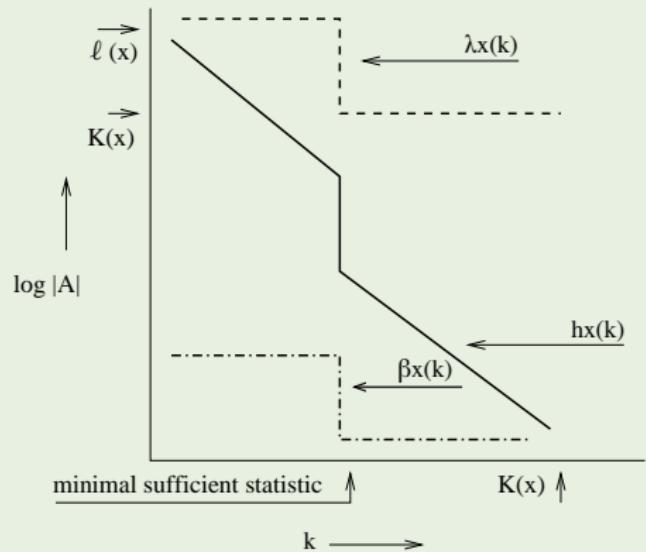
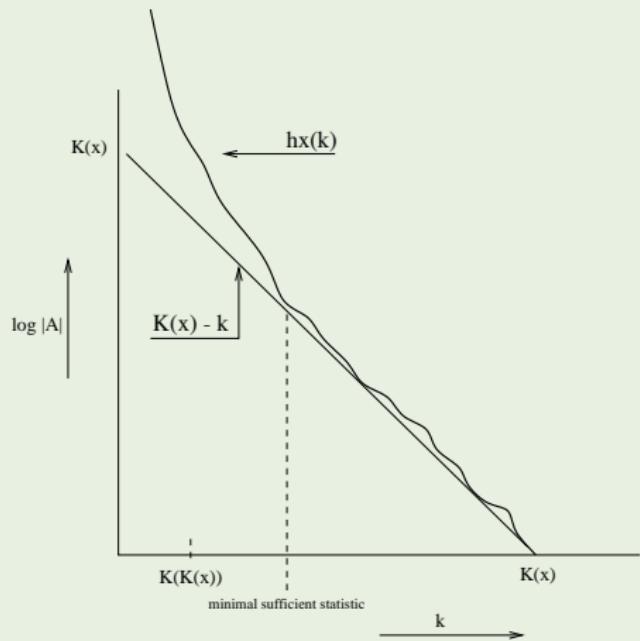
$$\text{depth}_b(x) := \min \{t: U^t(p) = x \text{ & } \ell(p) - K(x) \leq b\}$$

x is called *shallow* if $\text{depth}_b(x) \stackrel{+}{\leq} \ell(x)$.

Let $\chi_i := [\![\varphi_i(i)]\!]$. $\chi_{1:\infty}$ is deep. Ω is shallow.

Effective Complexity

$\delta(x A) := \log A - K(x A)$	[randomness deficiency]
$\delta(x \mu) := \log \frac{M(x)}{\mu(x)}$	$[\mu\text{-randomness deficiency}]$
$\beta_x(k) := \min_A \{\delta(x A) : x \in A \text{ & } K(A) \leq k\}$	[Best-Fit]
$h_x(k) := \min_A \{\log A : x \in A \text{ & } K(A) \leq k\}$	[Kolmogorov structure function / ML]
$h_x(k) := \min_\mu \{-\log \mu(x) : K(\mu) \leq k\}$	[Kolmogorov structure function / ML]
$A^*(x) := \iota A \left[x \in A \text{ & } K(A) = \mu k \left[k + h_x(k) \stackrel{+}{=} K(x) \right] \right]$	[Kolmogorov minimal sufficient statistic]
$\lambda_x(k) := \min_A \{K(A) + \log A : x \in A \text{ & } K(A) \leq k\}$	[MDL]
$\lambda_x(k) := \min_\mu \{K(\mu) - \log \mu(x) : K(\mu) \leq k\}$	[MDL]
$\Delta(x A) := K(A) + \log A - K(x)$	[discrepancy]
$soph_c(x) := \min_A \{K(A) : \Delta(x A) < c\}$	[sophistication]
$csoph(x) := \min_A \{K(A) + \Delta(x A)\}$	[coarse sophistication]
$\Sigma(\mu) := K(\mu) + H(\mu)$	[total information]
$\mathcal{E}_{\delta, \Delta}(x \mathcal{M}) := \min_{\mu \in \mathcal{M}} \left\{ K(\mu) : \Sigma(\mu) - K(x) \leq \Delta \text{ & } \mu(x) \geq 2^{-H(\mu)(1+\delta)} \right\}$	[effective complexity]
$\mathcal{E}_\delta(x \mathcal{M}) := \min_{\mu \in \mathcal{M}} \left\{ K(\mu) + \Sigma(\mu) - K(x) : \mu(x) \geq 2^{-H(\mu)(1+\delta)} \right\}$	[coarse effective complexity]



Theorem

For x and k ,

$$\lambda_x(k) \leq h_x(k) + k \stackrel{+}{\leq} \lambda_x(k) + K(k)$$

For k with $0 \leq k \leq K(x) - O(\log \ell(x))$,

$$\beta_x(k) + K(x) \stackrel{+}{\leq} \lambda_x(k)$$

$$\lambda_x(k + O(\log \ell(x))) \leq \beta_x(k) + K(x)$$

In other words, the equality

$$\beta_x(k) + K(x) = \lambda_x(k) = h_x(k) + k$$

holds within logarithmic additive terms in argument and value.

Sophistication and Computational Depth

Theorem

$$csoph(x) = \min_c \{soph_c(x) + c\}$$

$K^t(x) := \min_p \{\ell(p) : U^t(p) = x\}$ (Time-bounded Kolmogorov Complexity)

$depth^t(x) := K^t(x) - K(x)$ (Basic Computational Depth)

$depth_{BB}(x) := \min_t \{depth^t(x) + K(t)\}$ (Busy Beaver Computational Depth)

Theorem

$$|csoph(x) - depth_{BB}(x)| \leq O(\log \ell(x))$$

Zurek's Physical Entropy

Definition (Physical Entropy)

Physical entropy $S(d)$ of a microstate d is the sum of the conditional Shannon entropy $H_d := - \sum_k P(k|d) \log P(k|d)$ and of the Kolmogorov Complexity $K(d)$.

$$S(d) := H_d + K(d)$$

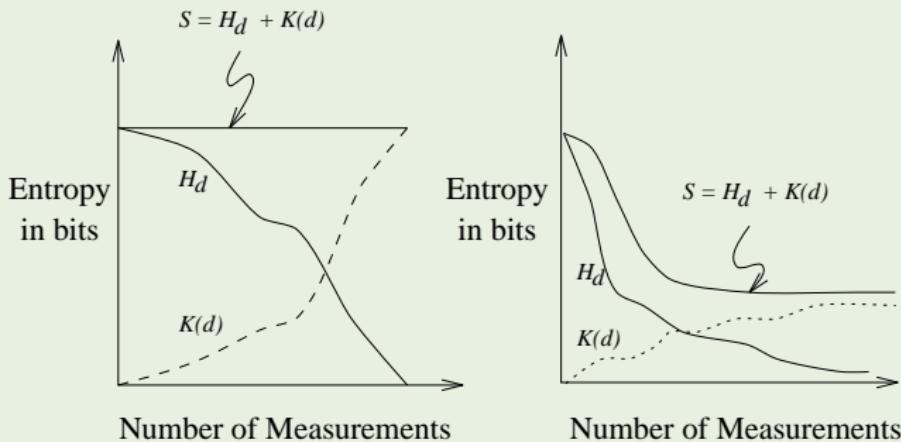


Figure: random vs regular microstate

Maxwell's Demon & Landauer's Principle

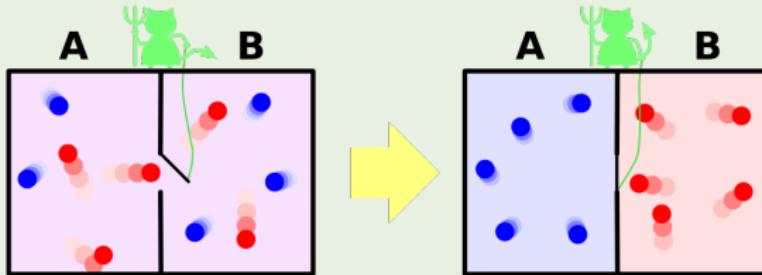


Figure: The demon turns entropy into information, the information-erasure operation turns information into entropy. In the course of ideal measurement on an equilibrium ensemble, the decrease of the entropy must be compensated by the increase of the size of the minimal record, and vice versa. $\Delta H \approx -\langle \Delta K \rangle$.

Landauer's Principle: Logically irreversible computation costs energy. Erasing 1 bit of information dissipates at least $kT \ln 2$ of heat into the environment.



Figure: Destroying information generates heat

Gács' Algorithmic Entropy

Definition (Algorithmic Entropy)

- *coarse-grained algorithmic entropy* of a cell Γ with respect to μ

$$H_\mu(\Gamma) := \log \mu(\Gamma) + K(\Gamma|\mu)$$

- *fine-grained algorithmic entropy* of $x \in 2^\omega$ with respect to μ

$$H_\mu(x) := \inf_n \{ \log \mu(x_{1:n}) + K(x_{1:n}|\mu) \}$$

- $-H_\mu(x)$ is a universal ML-test: x is μ -ML-random iff $H_\mu(x) > -\infty$.
- The fine-grained algorithmic entropy of a microstate can be approximated by the coarse-grained algorithmic entropies of successively smaller cells containing it.

$$H_\mu(x) = \inf_n \{ H_\mu(\Gamma_{x_{1:n}}) \}$$

Contents

- ① History of AI
- ② Philosophy of Induction
- ③ Inductive Logic
- ④ Universal Induction
- ⑤ Reinforcement Learning
- ⑥ General Reinforcement Learning

Preferences Lead to Utility

- A lottery $L = [p_1, x_1; \dots; p_n, x_n]$ is a probability distribution over outcomes \mathcal{X} .
- Preferences of a rational agent must obey constraints.
 - ① completeness $x_1 > x_2 \vee x_1 \sim x_2 \vee x_2 > x_1$
 - ② transitivity $x_1 > x_2 \wedge x_2 > x_3 \rightarrow x_1 > x_3$
 - ③ continuity $x_1 > x_2 > x_3 \rightarrow \exists p: x_2 \sim [p, x_1; 1 - p, x_3]$
 - ④ independence $x_1 \sim x_2 \rightarrow [p, x_1; 1 - p, x_3] \sim [p, x_2; 1 - p, x_3]$
 - ⑤ monotonicity $x_1 > x_2 \wedge p > q \rightarrow [p, x_1; 1 - p, x_2] > [q, x_1; 1 - q, x_2]$
 - ⑥ decomposability

$$[p, x_1; 1 - p, [q, x_2; 1 - q, x_3]] \sim [p, x_1; (1 - p)q, x_2; (1 - p)(1 - q), x_3]$$

Theorem (von Neumann & Morgenstern 1944)

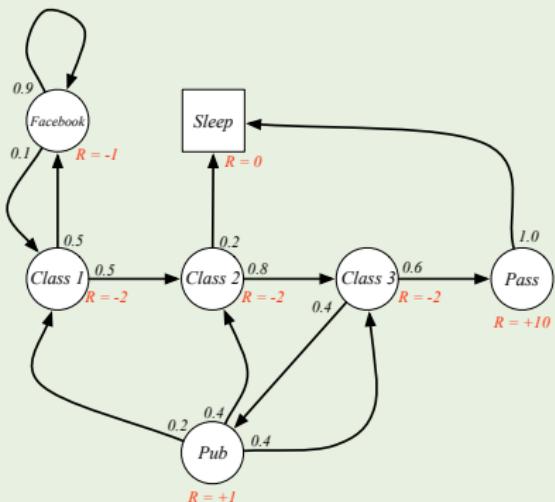
If a preference relation \succeq satisfies the above constraints, then there exists a function $u: \mathcal{X} \rightarrow [0, 1]$ such that

$$x_1 \succeq x_2 \iff u(x_1) \geq u(x_2) \quad \text{and}$$

$$u([p_1, x_1; \dots; p_n, x_n]) = \sum_{i=1}^n p_i u(x_i)$$

Markov Decision Process

A (finite) MDP $(\mathcal{S}, \mathcal{A}, P, R)$.



Definition (Value of a state under π)

$$V^\pi(s) := \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid S_t = s \right]$$

Definition (Action-value under π)

$$Q^\pi(s, a) := \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid S_t = s, A_t = a \right]$$

Bellman Expectation Equations

$$r(s, a) := \mathbb{E}[r|s, a]$$

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) Q^\pi(s, a)$$

$$Q^\pi(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V^\pi(s')$$

$$V^\pi(s) = \mathbb{E} [r + \gamma V^\pi(s') | s] = \sum_{a \in \mathcal{A}} \pi(a|s) \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V^\pi(s') \right)$$

advantage

$$A^\pi(s, a) := Q^\pi(s, a) - V^\pi(s)$$

Bellman Optimality Equations

Definition (Optimal Values)

$$V^*(s) := \max_{\pi} V^{\pi}(s)$$

$$Q^*(s, a) := \max_{\pi} Q^{\pi}(s, a)$$

Definition (Optimal Policy)

A policy π is called optimal if $\forall s \in \mathcal{S}: V^{\pi}(s) = V^*(s)$.

$$V^*(s) = \max_{a \in \mathcal{A}} Q^*(s, a)$$

$$Q^*(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V^*(s')$$

Action/Policy Evaluation Operator & Greedy Policy

Definition (Action Evaluation Operator)

$$T_a V(s) := r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a)V(s')$$

Definition (Policy Evaluation Operator)

$$T^\pi V(s) := \sum_{a \in \mathcal{A}} \pi(a|s) T_a V(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a)V(s') \right)$$
$$T^*V(s) := \max_{a \in \mathcal{A}} T_a V(s) = \max_{a \in \mathcal{A}} \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a)V(s') \right)$$

Definition (Greedy policy)

Policy π is greedy w.r.t. V if $T^\pi V = T^*V$.

Banach's Fixpoint Theorem

Theorem (Banach's Fixpoint Theorem)

Let \mathcal{V} be a Banach space and $T: \mathcal{V} \rightarrow \mathcal{V}$ be a contraction mapping, with Lipschitz constant $\gamma < 1$. Then T has a unique fixpoint $v \in \mathcal{V}$. Further, for each $v_0 \in \mathcal{V}$, $\lim_{n \rightarrow \infty} \|T^n(v_0) - v\| = 0$, and the convergence is geometric:

$$\|T^n(v_0) - v\| \leq \gamma^n \|v_0 - v\|$$

Application of Banach's Fixpoint Theorem

Theorem

$(\mathcal{V}, \|\cdot\|_\infty)$ is a Banach space, where $\mathcal{V} := \{V \in \mathbb{R}^S : \|V\|_\infty < \infty\}$ and $\|V\|_\infty := \max_{s \in S} |V(s)|$.

- T^π is a contraction, and V^π is the unique fixpoint of T^π .

$$\lim_{n \rightarrow \infty} \|(T^\pi)^n V_0 - V^\pi\|_\infty = 0$$

- T^* is a contraction, and V^* is the unique fixpoint of T^* .

$$\lim_{n \rightarrow \infty} \|(T^*)^n V_0 - V^*\|_\infty = 0$$

Application of Banach's Fixpoint Theorem

$$T^\pi Q(s, a) := r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \sum_{a' \in \mathcal{A}} \pi(a'|s') Q(s', a')$$

$$T^* Q(s, a) := r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \max_{a' \in \mathcal{A}} Q(s', a')$$

- T^π is a contraction, and Q^π is the unique fixpoint of T^π .

$$\lim_{n \rightarrow \infty} \|(T^\pi)^n Q_0 - Q^\pi\|_\infty = 0$$

- T^* is a contraction, and Q^* is the unique fixpoint of T^* .

$$\lim_{n \rightarrow \infty} \|(T^*)^n Q_0 - Q^*\|_\infty = 0$$

Two Theorems

Theorem (Fixpoint of Bellman Optimality Operator)

Let V be the fixpoint of T^ and assume that there is policy π which is greedy w.r.t V . Then $V = V^*$ and π is an optimal policy.*

Theorem (Policy Improvement Theorem)

*Choose some stationary policy π_0 and let π be greedy w.r.t. V^{π_0} . Then $V^\pi \geq V^{\pi_0}$, i.e., π is an improvement upon π_0 . In particular, if $T^*V^{\pi_0}(s) > V^{\pi_0}(s)$ for some state s then π strictly improves upon π_0 at s : $V^\pi(s) > V^{\pi_0}(s)$. On the other hand, when $T^*V^{\pi_0}(s) = V^{\pi_0}(s)$ then π_0 is an optimal policy.*

Solving MDPs — Finite-Horizon Dynamic Programming

Principle of optimality: the tail of an optimal policy is optimal for the “tail” problem.

Backward Induction

- Backward recursion: $V_N^*(s) = r_N(s)$ and for $k = N - 1, \dots, 0$

$$V_k^*(s) = \max_{a \in \mathcal{A}_k} \left(r_k(s, a) + \sum_{s' \in \mathcal{S}_{k+1}} P_k(s'|s, a) V_{k+1}^*(s') \right)$$

- Optimal policy: for $k = 0, \dots, N - 1$

$$\pi_k^*(s) \in \operatorname{argmax}_{a \in \mathcal{A}_k} \left(r_k(s, a) + \sum_{s' \in \mathcal{S}_{k+1}} P_k(s'|s, a) V_{k+1}^*(s') \right)$$

- Cost: $N|\mathcal{S}||\mathcal{A}|$ vs $|\mathcal{A}|^{N|\mathcal{S}|}$ of brute force policy search.
- From now on, we will consider infinite-horizon discounted MDPs.

Solving MDPs — Value Iteration

Theorem (Principle of Optimality)

A policy π achieves the optimal value from state s , $V^\pi(s) = V^*(s)$, iff, for any state s' reachable from s , π achieves the optimal value from state s' , $V^\pi(s') = V^*(s')$.

Any optimal policy π^* can be subdivided into two components:

- an optimal first action a^* ,
- followed by an optimal policy from successor state s' .

$$V^*(s) = \max_{a \in \mathcal{A}} \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V^*(s') \right)$$

Value Iteration: $V_{k+1} \leftarrow T^*V_k$

$$V_1 \rightarrow V_2 \rightarrow \dots \rightarrow V^*$$

$$\|V_{k+1} - V_k\|_\infty < \varepsilon \implies \|V_{k+1} - V^*\|_\infty < \frac{2\gamma\varepsilon}{1-\gamma}$$

Solving MDPs — Policy Iteration

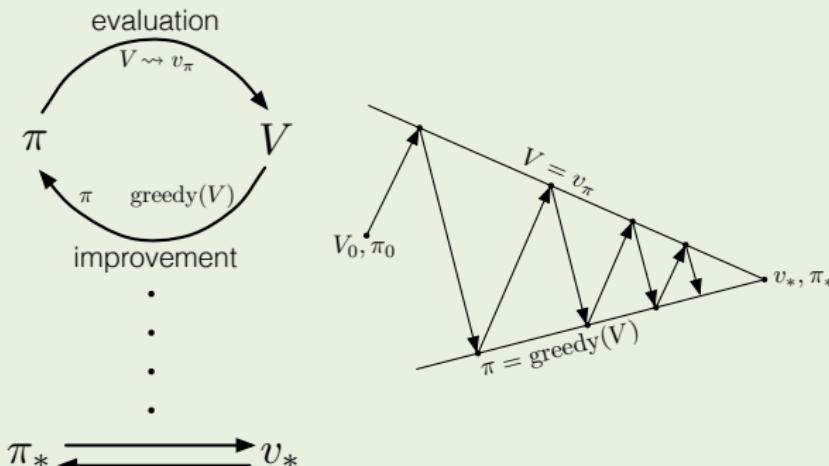
Policy Iteration: $\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \dots \xrightarrow{I} \pi^* \xrightarrow{E} V^*$

- E — policy evaluation:

$$V_{k+1}^\pi \leftarrow T^\pi V_k^\pi$$

- I — policy improvement:

$$\pi_{k+1}(s) := \operatorname{argmax}_{a \in \mathcal{A}} \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V^{\pi_k}(s') \right)$$



Monte-Carlo Methods

MC learns from complete episodes of raw experience without modeling the environmental dynamics and computes the observed mean return as an approximation of the expected return.

$$G_t := \sum_{k=0}^{T-t-1} \gamma^k r_{t+k+1}$$

$$V(s) := \frac{\sum_{t=1}^T \llbracket S_t = s \rrbracket G_t}{\sum_{t=1}^T \llbracket S_t = s \rrbracket}$$

$$Q(s, a) := \frac{\sum_{t=1}^T \llbracket S_t = s, A_t = a \rrbracket G_t}{\sum_{t=1}^T \llbracket S_t = s, A_t = a \rrbracket}$$

Temporal-Difference Learning

TD Learning is model-free and learns from incomplete episodes of experience.

$$V(s_t) \leftarrow (1 - \alpha)V(s_t) + \alpha G_t$$

$$V(s_t) \leftarrow V(s_t) + \alpha(G_t - V(s_t))$$

$$V(s_t) \leftarrow V(s_t) + \alpha(r_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$$

- MC updates value $V(s_t)$ toward actual return G_t .

$$V(s_t) \leftarrow V(s_t) + \alpha(G_t - V(s_t))$$

- TD updates value $V(s_t)$ toward estimated return $r_{t+1} + \gamma V(s_{t+1})$.

$$V(s_t) \leftarrow V(s_t) + \underbrace{\alpha(r_{t+1} + \gamma V(s_{t+1}) - V(s_t))}_{\text{TD error}}$$

TD target

SARSA: On-Policy TD control

SARSA

- ① At time step t , we start from state s_t and pick action according to Q values, $a_t = \operatorname{argmax}_{a \in \mathcal{A}} Q(s_t, a)$; ε -greedy is commonly applied.
- ② With action a_t , we observe reward r_{t+1} and get into s_{t+1} .
- ③ Then pick the next action $a_{t+1} = \operatorname{argmax}_{a \in \mathcal{A}} Q(s_{t+1}, a)$.
- ④ Update the action-value function:
$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)).$$
- ⑤ $t = t + 1$ and repeat from step 1.

Expected SARSA

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left(r_{t+1} + \gamma \sum_{a \in \mathcal{A}} \pi(a|s_{t+1}) Q(s_{t+1}, a) - Q(s_t, a_t) \right)$$

Q-Learning: Off-policy TD control

- ① At time step t , we start from state s_t and pick action according to Q values, $a_t = \operatorname{argmax}_{a \in \mathcal{A}} Q(s_t, a)$; ε -greedy is commonly applied.
- ② With action a_t , we observe reward r_{t+1} and get into s_{t+1} .
- ③ Update the action-value function:
$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left(r_{t+1} + \gamma \max_{a \in \mathcal{A}} Q(s_{t+1}, a) - Q(s_t, a_t) \right).$$
- ④ $t = t + 1$ and repeat from step 1.

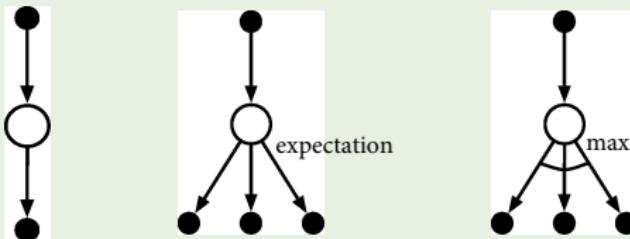
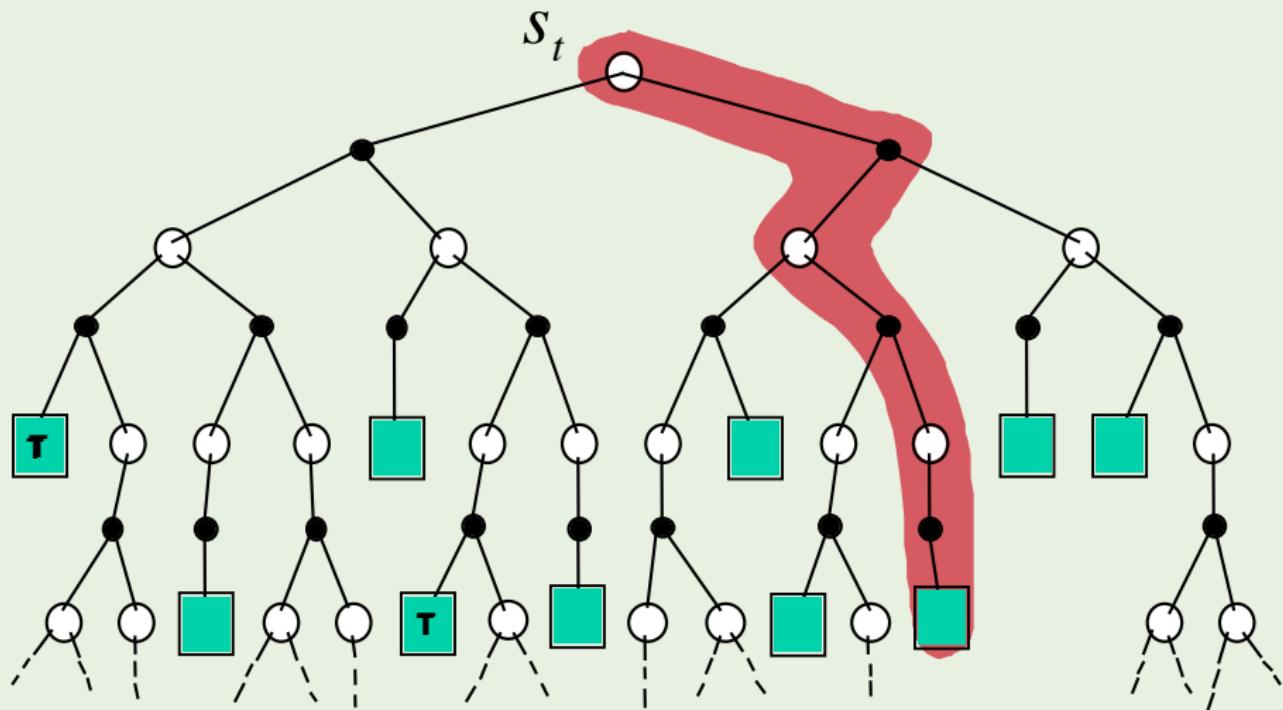


Figure: SARSAR, Expected SARSA, and Q-Learning

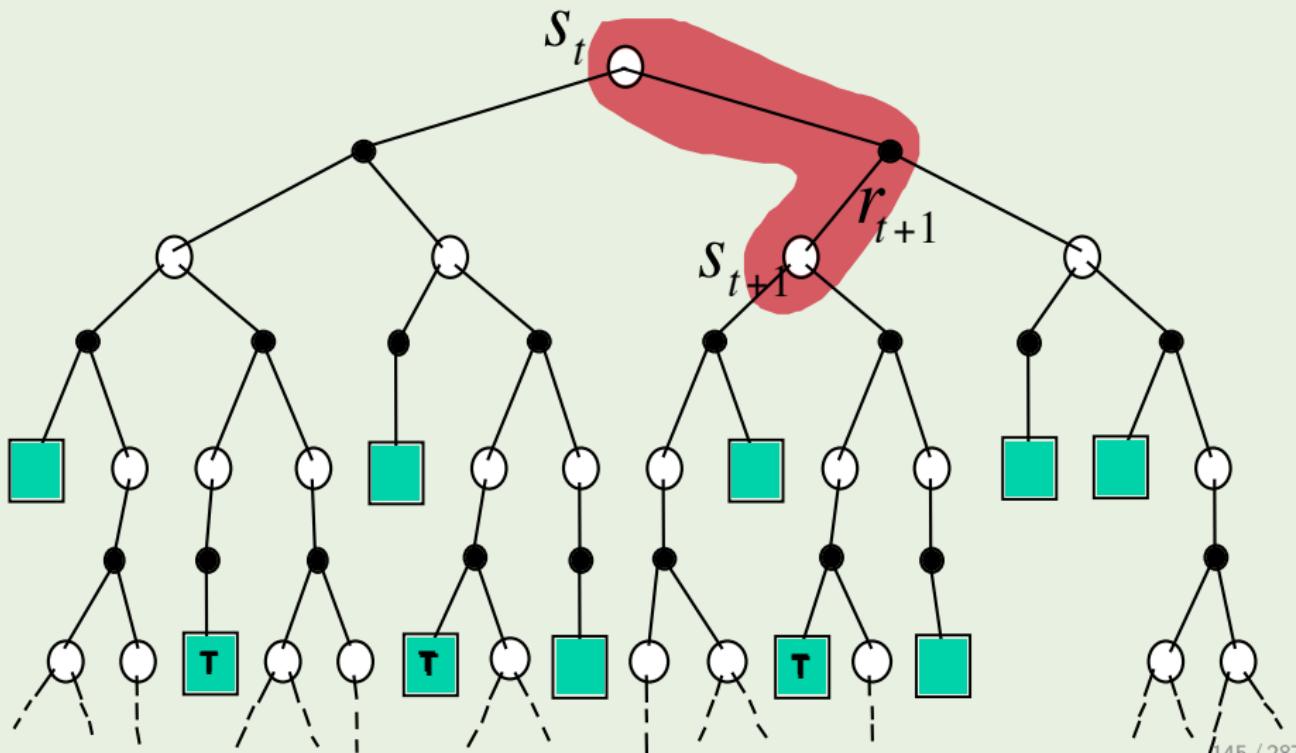
Monte-Carlo Backup

$$V(s_t) \leftarrow V(s_t) + \alpha(G_t - V(s_t))$$



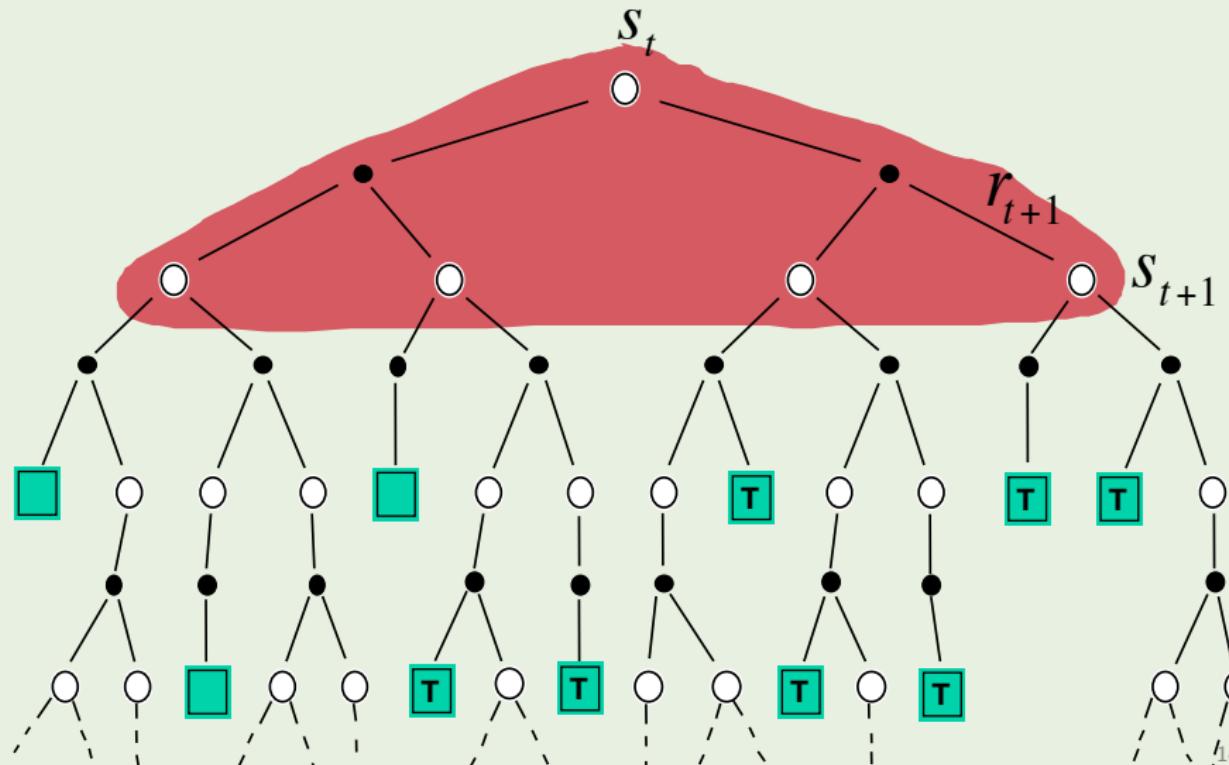
Temporal-Difference Backup

$$V(s_t) \leftarrow V(s_t) + \alpha(r_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$



Dynamic Programming Backup

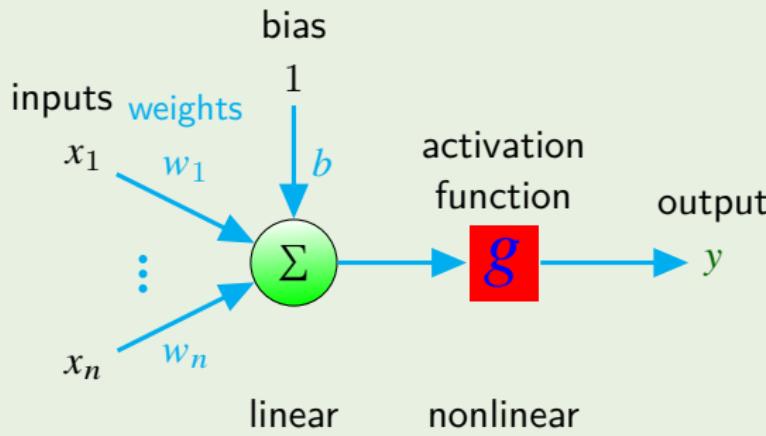
$$V(s_t) \leftarrow \mathbb{E}_\pi[r_{t+1} + \gamma V(s_{t+1})]$$



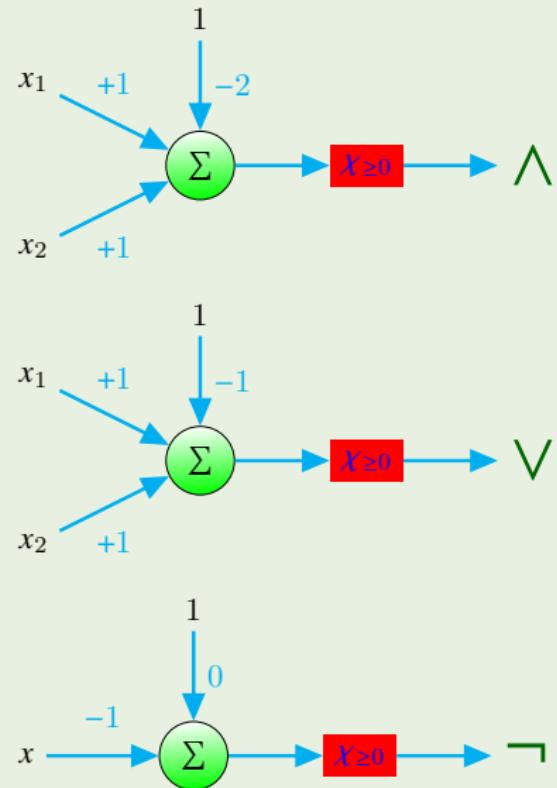
Contents

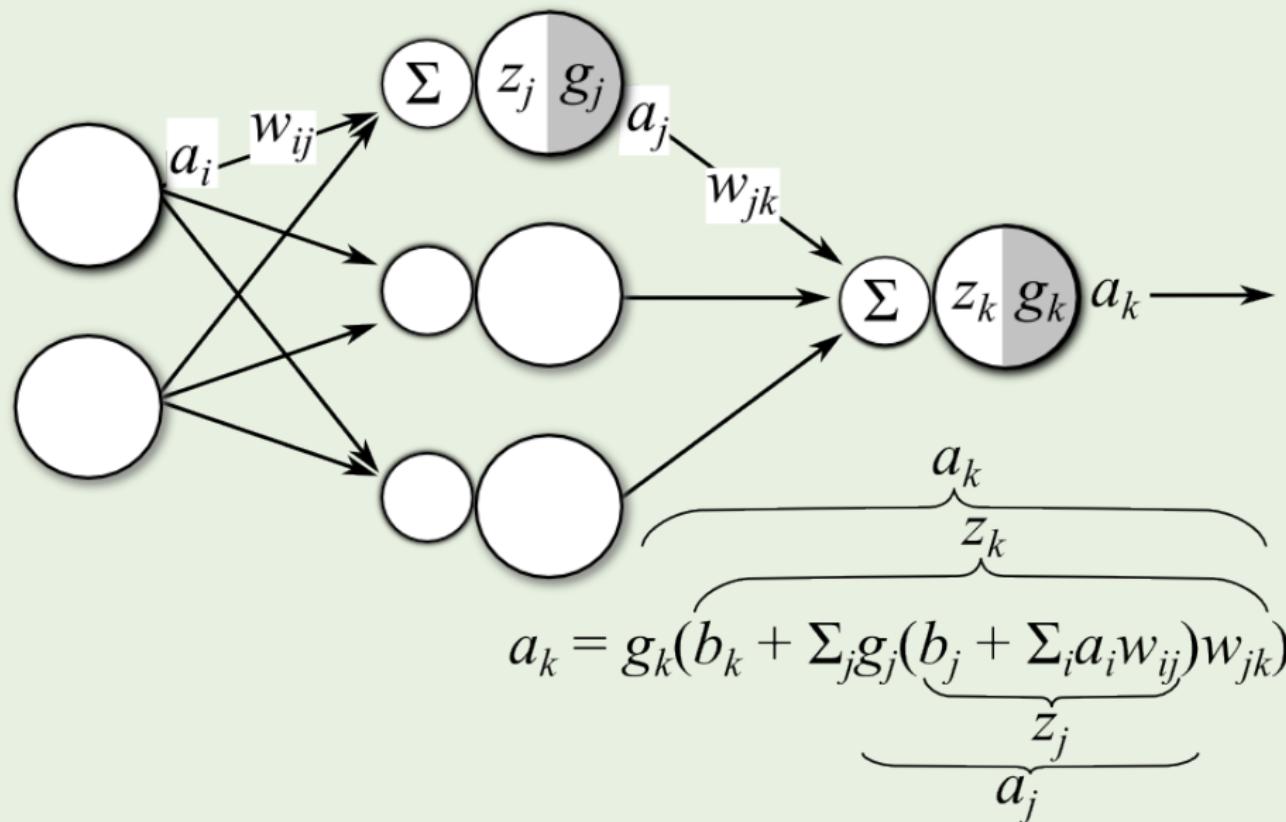
- ① History of AI
- ② Philosophy of Induction
- ③ Inductive Logic
- ④ Universal Induction
- ⑤ Reinforcement Learning
 Artificial Neural Network and Deep Reinforcement Learning
- ⑥ General Reinforcement Learning

McCulloch-Pitts Artificial Neural Network



$$y = g \left(\sum_{i=1}^n w_i x_i + b \right)$$





Learning: small change in weights → small change in output

Kolmogorov Superposition Theorem

Theorem (Kolmogorov Superposition Theorem)

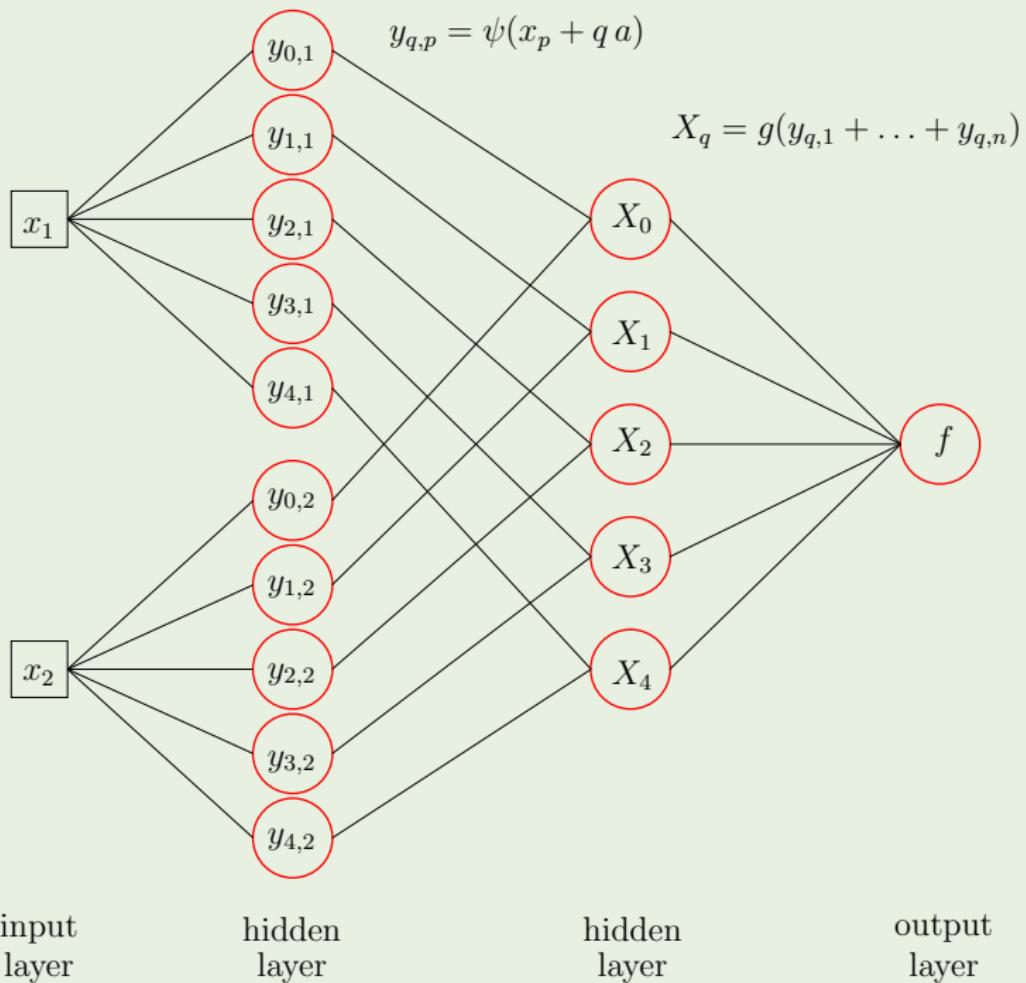
For each $n \geq 2$ there exists a computable function $\psi: [0, 1] \rightarrow \mathbb{R}$ and computable constants $a, \lambda_{pq} \in \mathbb{R}$, $p = 1, \dots, n$, $q = 0, \dots, 2n$ s.t.: every continuous function $f: [0, 1]^n \rightarrow \mathbb{R}$ has a representation as

$$f(x_1, \dots, x_n) = \sum_{q=0}^{2n} g\left(\sum_{p=1}^n \lambda_{pq} \psi(x_p + qa)\right)$$

for some continuous function $g: [0, 1] \rightarrow \mathbb{R}$ that is computable from f .

Theorem (Hecht-Nielsen Theorem)

The class of functions $f: [0, 1]^n \rightarrow \mathbb{R}$, implementable by three-layer feed-forward neural networks with (computable) continuous activation functions $g: [0, 1] \rightarrow \mathbb{R}$ and (computable) weights $\lambda \in \mathbb{R}$, is exactly the class of (computable) continuous functions $f: [0, 1]^n \rightarrow \mathbb{R}$.



Universal Approximation Theorem

Theorem (Universal Approximation Theorem)

Let g be a nonconstant, bounded, and increasing continuous function. Let I_n be any compact subset of \mathbb{R}^n . The space of continuous functions on I_n is denoted by $C(I_n, \mathbb{R})$. Then, given any function $f \in C(I_n, \mathbb{R})$ and $\varepsilon > 0$, there exists an integer N , real constants $v_i, b_i \in \mathbb{R}$ and real vectors $w_i \in \mathbb{R}^n$, where $i = 1, \dots, N$, s.t.

$$\forall x \in I_n : |h(x) - f(x)| < \varepsilon$$

where

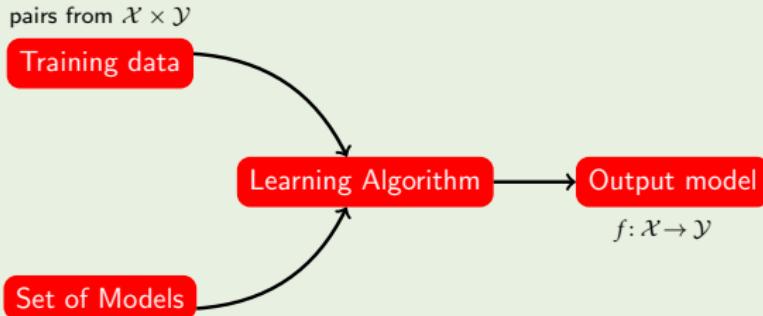
$$h(x) := \sum_{i=1}^N v_i g(w_i^\top x + b_i)$$

In other words, functions of the form $h(x)$ are dense in $C(I_n, \mathbb{R})$.

Representation & Approximation

- A feed-forward network with 1 hidden layer can represent any boolean function, but require exponential hidden units.
- A feed-forward network with 2 hidden layers and (computable) continuous activation functions can represent any (computable) continuous function.
- A feed-forward network with a linear output layer and at least 1 hidden layer and continuous and differentiable activation functions can approximate any Borel measurable function from one finite-dimensional space to another with any desired non-zero amount of error.
- A feed-forward network with 2 hidden layers and continuous and differentiable activation functions can approximate any function.

Deep Learning



- ① **hypothesis space** — Network Structure — f_θ
- ② **the goodness of a function** — Learning Target — loss function ℓ
- ③ **pick the best function** — Learn — find the network parameters
 $\theta^* := \underset{\theta}{\operatorname{argmin}} L(\theta)$ that minimize total cost $L(\theta)$ by gradient decent

$$\theta \leftarrow \theta - \eta \nabla_{\theta} L(\theta)$$

where $L(\theta) := \mathbb{E}_P [\ell(f_\theta(a), t)] + \lambda \Omega(\theta)$ and $\Omega(\theta)$ is a regularizer.

Deep Learning

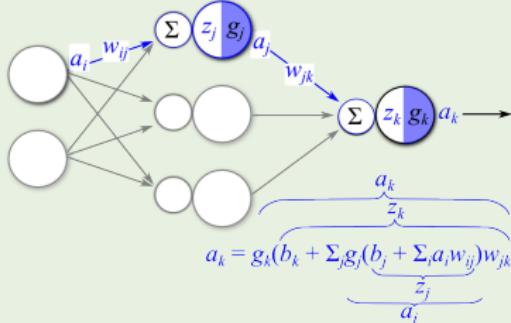
$$f_{\theta}: \mathbf{a}^{(0)} \mapsto \mathbf{g}^{(n)} \left(\cdots \mathbf{g}^{(2)} \left(\underbrace{\mathbf{g}^{(1)} \left(\underbrace{\mathbf{w}^{(1)} \mathbf{a}^{(0)} + \mathbf{b}^{(1)}}_{\mathbf{z}^{(1)}} \right) + \mathbf{b}^{(2)}}_{\mathbf{z}^{(2)}} \right) \cdots \right)$$
$$\mathbf{a}^{(0)}$$
$$\mathbf{z}^{(1)}$$
$$\mathbf{a}^{(1)}$$
$$\mathbf{z}^{(2)}$$
$$\mathbf{a}^{(2)}$$
$$\mathbf{z}^{(n)}$$
$$\mathbf{a}^{(n)}$$

where parameters $\theta := \{\mathbf{w}^{(i)}, \mathbf{b}^{(i)}\}_{i=1}^n$ and activation functions \mathbf{g}

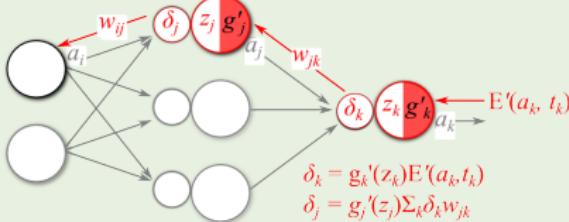
$$\sigma(z) = \frac{1}{1 - e^{-z}} \quad \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad \text{ReLU}(z) = \max(0, z)$$

$$\left. \begin{array}{l} a_0^{(l)} := 1 \\ w_{0j}^{(l)} := b_j^{(l)} \end{array} \right\} \implies \left\{ \begin{array}{l} z_j^{(l+1)} := \sum_i w_{ij}^{(l+1)} a_i^{(l)} \\ a_j^{(l+1)} := g_j^{(l+1)} (z_j^{(l+1)}) \end{array} \right.$$

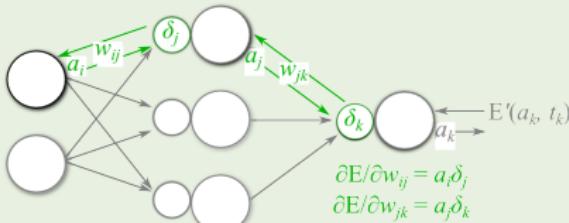
I. Forward-propagate Input Signal



II. Back-propagate Error Signals



III. Calculate Parameter Gradients



IV. Update Parameters

$$w_{ij} = w_{ij} - \eta(\frac{\partial E}{\partial w_{ij}})$$

$$w_{jk} = w_{jk} - \eta(\frac{\partial E}{\partial w_{jk}})$$

for learning rate η

Backpropagation

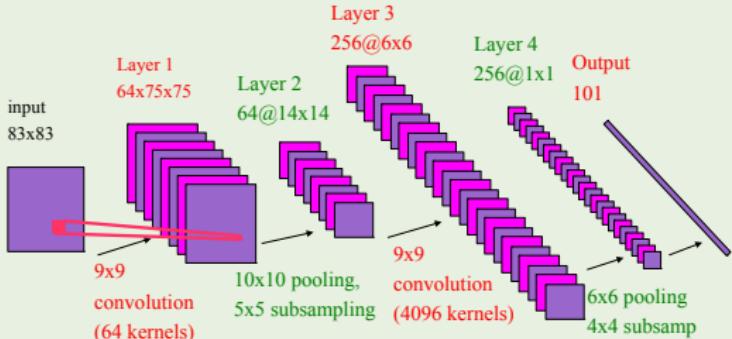
$$\delta_j^{(l)} := \frac{\partial L}{\partial z_j^{(l)}}$$

$$\delta_j^{(l)} = g_j^{(l)'}(z_j^{(l)}) \sum_k \delta_k^{(l+1)} w_{jk}^{(l+1)}$$

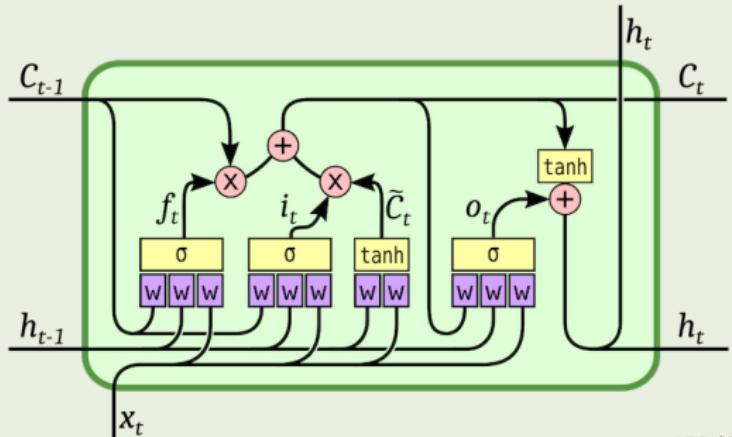
$$\frac{\partial L}{\partial w_{ij}^{(l)}} = a_i^{(l-1)} \delta_j^{(l)}$$

$$w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta \frac{\partial L}{\partial w_{ij}^{(l)}}$$

- network structure?
- how many layers?
- how many units per layer?
- loss function?
- regularization?
- weight decay?
- learning rate?
- activation function?
- early stopping?
- dropout?
- mini-batch?
- momentum?
- ...



$$a_{ij}^{(l+1)} = g_{ij}^{(l+1)} \left(\sum_{m=0}^{k-1} \sum_{n=0}^{k-1} w_{m,n} a_{i+m, j+n}^{(l)} + b \right)$$



Key Properties of CNNs

$$\begin{array}{|c|c|c|c|c|c|} \hline 1 & 0 & 0 & 0 & 0 & 1 \\ \hline 0 & 1 & 0 & 0 & 1 & 0 \\ \hline 0 & 0 & 1 & 1 & 0 & 0 \\ \hline 1 & 0 & 0 & 0 & 1 & 0 \\ \hline 0 & 1 & 0 & 0 & 1 & 0 \\ \hline 0 & 0 & 1 & 0 & 1 & 0 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline 1 & -1 & -1 \\ \hline -1 & 1 & -1 \\ \hline -1 & -1 & 1 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 3 & -1 & -3 & -1 \\ \hline -3 & 1 & 0 & -3 \\ \hline -3 & -3 & 0 & 1 \\ \hline 3 & -2 & -2 & -1 \\ \hline \end{array}$$

Table: Convolution (stride 1)

Take advantage of the structure of the data!

- Convolutional Filters (**Translation invariance**)
- Multiple layers (**Compositionality**)
- Filters localized in space (**Locality**)
- Weight sharing (**Self-similarity**)

DQN

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s, a \right]$$

$$\max_{\pi} Q^\pi(s, a) =: Q^*(s, a) = \mathbb{E}_{s'} \left[r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right]$$

$$Q_{t+1}(s, a) = \mathbb{E}_{s'} \left[r + \gamma \max_{a'} Q_t(s', a') \mid s, a \right]$$

$$Q_t \xrightarrow{t \rightarrow \infty} Q^*$$

$$Q_{t+1}(s, a) = Q_t(s, a) + \alpha \left(r + \gamma \max_{a'} Q_t(s', a') - Q_t(s, a) \right)$$

$$Q(s, a; \theta) \approx Q^*(s, a)$$

$$L(\theta) := \mathbb{E}_{s, a, r, s'} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta^-) - Q(s, a; \theta) \right)^2 \right] \quad (\text{DQN})$$

DQN

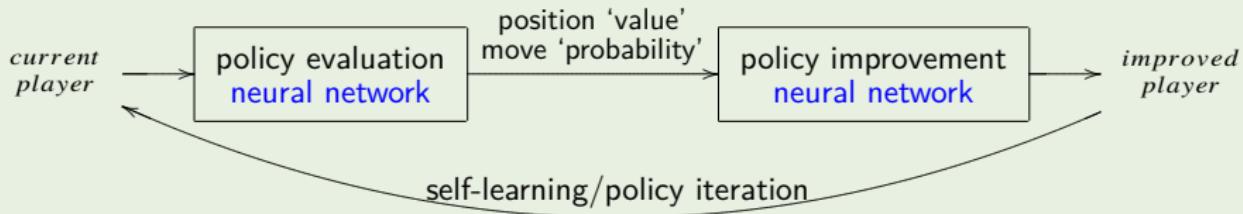
$$L(\theta) := \mathbb{E}_{s,a,r,s'} \left[\left(r + \gamma \underset{a'}{\operatorname{argmax}} Q(s', a'; \theta); \theta^- \right) - Q(s, a; \theta) \right]^2$$

(Double DQN)

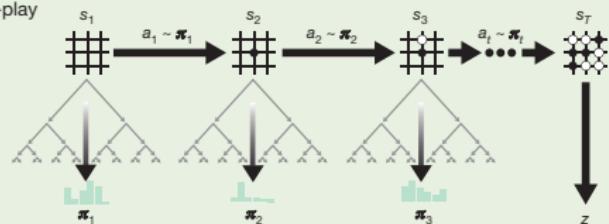
$$Q(s, a) = V(s; \theta) + A(s, a; \theta') \quad (\text{Dueling Network})$$

$$\begin{cases} \nabla_{\theta} \log \pi(a_t | s_t; \theta) \left(\sum_{i=0}^{k-1} \gamma^i r_{t+i} + \gamma^k V(s_{t+k}; \theta') - V(s_t; \theta') \right) & \text{actor} \\ \nabla_{\theta'} \left(\sum_{i=0}^{k-1} \gamma^i r_{t+i} + \gamma^k V(s_{t+k}; \theta') - V(s_t; \theta') \right)^2 & \text{critic} \end{cases} \quad (\text{AC})$$

AlphaZero



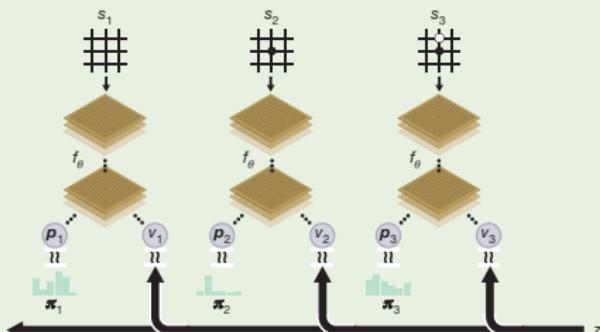
a Self-play



$$(\mathbf{p}, v) = f_{\theta}(s)$$

$$\ell = (z - v)^2 - \pi^{\top} \log \mathbf{p} + c \|\theta\|^2$$

b Neural network training

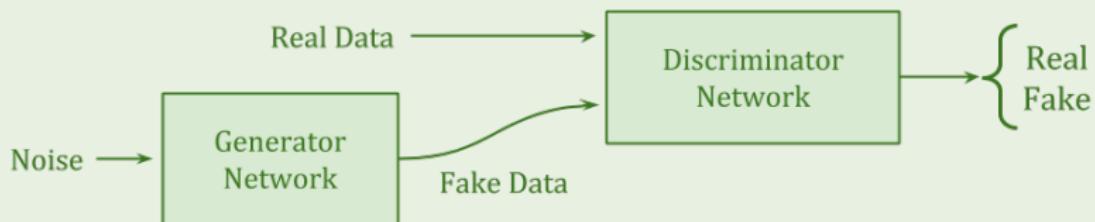


Intuition + Calculation
DNN + MCTS

GAN — Generative Adversarial Network

$$V(D, G) = \mathbb{E}_{x \sim P_{data}} [\log D(x)] + \mathbb{E}_{z \sim P_{noise}} [\log (1 - D(G(z)))]$$

$$G^* = \operatorname{argmin}_G \max_D V(D, G) \quad (\text{GAN})$$



Why “Deep” rather than “Fat”?

- Exploiting compositionality gives an exponential gain in representational power.
 - Distributed representations: feature learning
 - Deep architecture: multiple levels of feature learning
- Each basic classifier can be trained by little data.
 - **deep → modularization → less training data?**
With more complex features, the number of parameters in the linear layers may be drastically decreased.
 - efficiency & sample complexity
 - better memory/computation trade-off?
- higher-level abstractions → easier generalization & transfer

Minimal Sufficient Statistic

Definition (Sufficient Statistic)

Let Y be a parameter indexing a family of probability distributions. Let X be random variable drawn from a probability distribution determined by Y . $T(X)$ is a sufficient statistic for Y if X is independent of Y given $T(X)$, i.e., $p(x|t, y) = p(x|t)$.

Definition (Minimal Sufficient Statistic)

A sufficient statistic $S(X)$ is minimal if for any sufficient statistic $T(X)$, there exists a function f s.t. $S = f(T)$ almost everywhere w.r.t X .

Theorem

- T is sufficient statistics for $Y \iff I(T(X); Y) = I(X; Y)$.
- S is minimal sufficient statistics for $Y \implies I(X; S(X)) \leq I(X; T(X))$.

Information Bottleneck — Learning is to forget!

Theorem

Let X be a sample drawn according to a distribution determined by the random variable Y . The set of solutions to

$$\min_T I(X; T) \quad \text{s.t.} \quad I(T; Y) = \max_{T'} I(T'; Y)$$

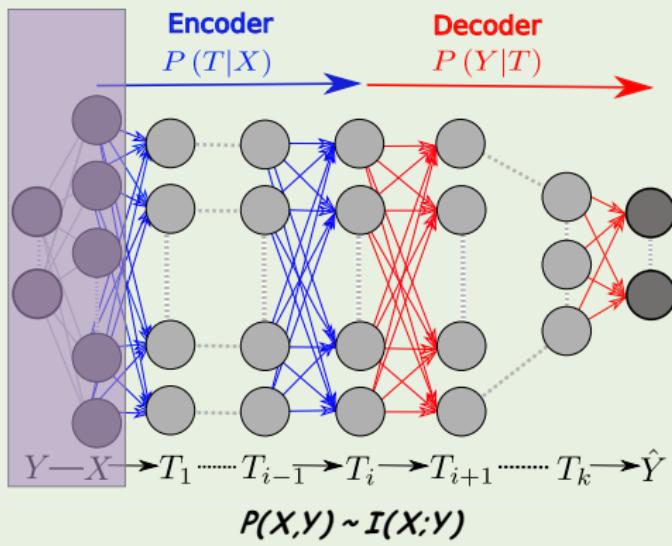
is exactly the set of minimal sufficient statistics for Y based on X .

Find a random variable T s.t.:

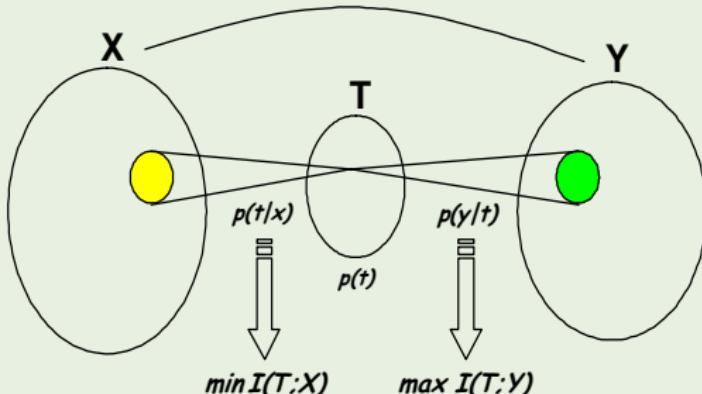
- $Y \leftrightarrow X \leftrightarrow T$ form a Markov chain.
- $I(X; T)$ is minimized (minimality, complexity term), while $I(T; Y)$ is maximized (sufficiency, accuracy term).

$$T^* := \underset{T : I(T(X); Y) = I(X; Y)}{\operatorname{argmin}} I(X; T(X))$$

is the Information Bottleneck between X and Y .



张三丰：将所
见到的剑招忘
得半点不剩，
才能得其神髓。



Information Bottleneck

$$\min_{p(t|x), p(y|t), p(t)} \left\{ I(X; T) - \beta I(T; Y) \right\} \text{ subject to Markov chain } Y \rightarrow X \rightarrow T.$$

$$\mathcal{L}[p(t|x)] := I(X; T) - \beta I(T; Y) - \sum_x \lambda(x) \sum_t p(t|x)$$

Let

$$\frac{\delta \mathcal{L}}{\delta p(t|x)} = 0$$

The solution is

$$p(t|x) = \frac{p(t)}{Z(x, \beta)} e^{-\beta D[p(y|x) \| p(y|t)]}$$

$$p(t) = \sum_x p(t|x)p(x)$$

$$p(y|t) = \sum_x p(y|x)p(x|t)$$

Expressiveness & Sample Complexity

Theorem

The hypothesis class of neural networks of depth T and size $O(T^2)$ contains all functions that can be implemented by a Turing machine within T operations, while having $O(T^2)$ sample complexity.

The Ultimate Hypothesis Space

- **No Free Lunch:** Sample complexity is exponentially large (w.r.t. the input dimension) if the hypothesis class is all possible functions.
- **Shallow learning (SVM, Boosting):** Hypothesis class is linear functions over manually determined features — strong prior knowledge.
- **Deep learning:** Hypothesis class is all functions implemented by determining the weights of a given artificial neural network.

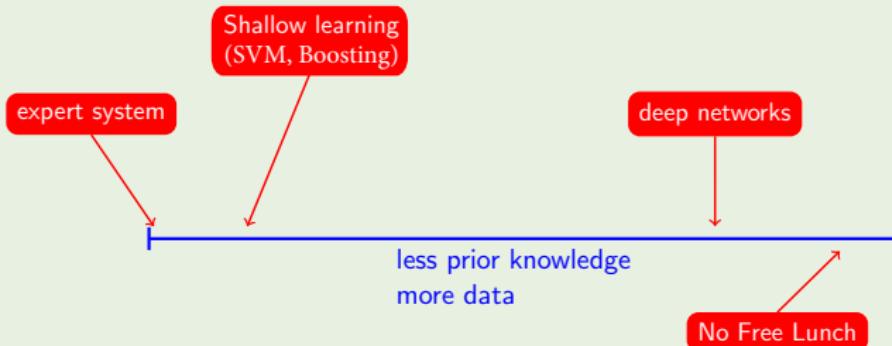


Figure: Prior vs Universality

Prior — a necessary good or a necessary evil?

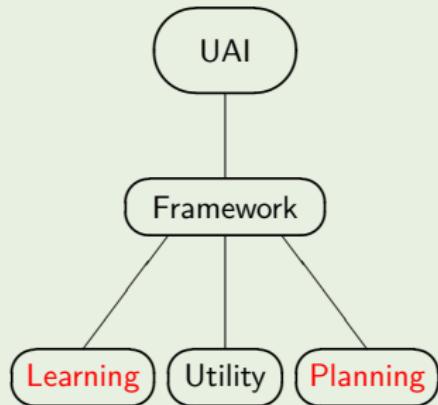
Contents

- ① History of AI
- ② Philosophy of Induction
- ③ Inductive Logic
- ④ Universal Induction
- ⑤ Reinforcement Learning
- ⑥ General Reinforcement Learning

- ① Solve intelligence
- ② Use it to solve everything else
 - learn automatically from raw inputs — not pre-programmed.
 - same algorithm, different tasks.

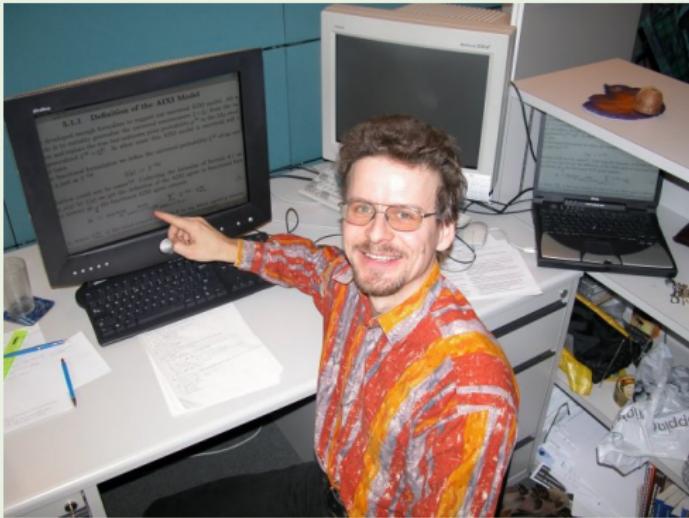
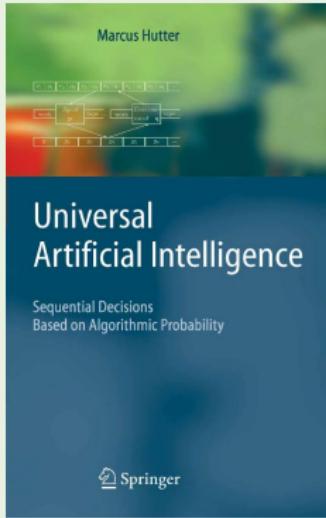
(Deep) RL	General RL
state space	history
ergodic	not ergodic
fully observable	partially observable
ϵ -exploration works	ϵ -exploration fails
MDP/DQN	AIXI

Table: (Deep) RL vs General RL

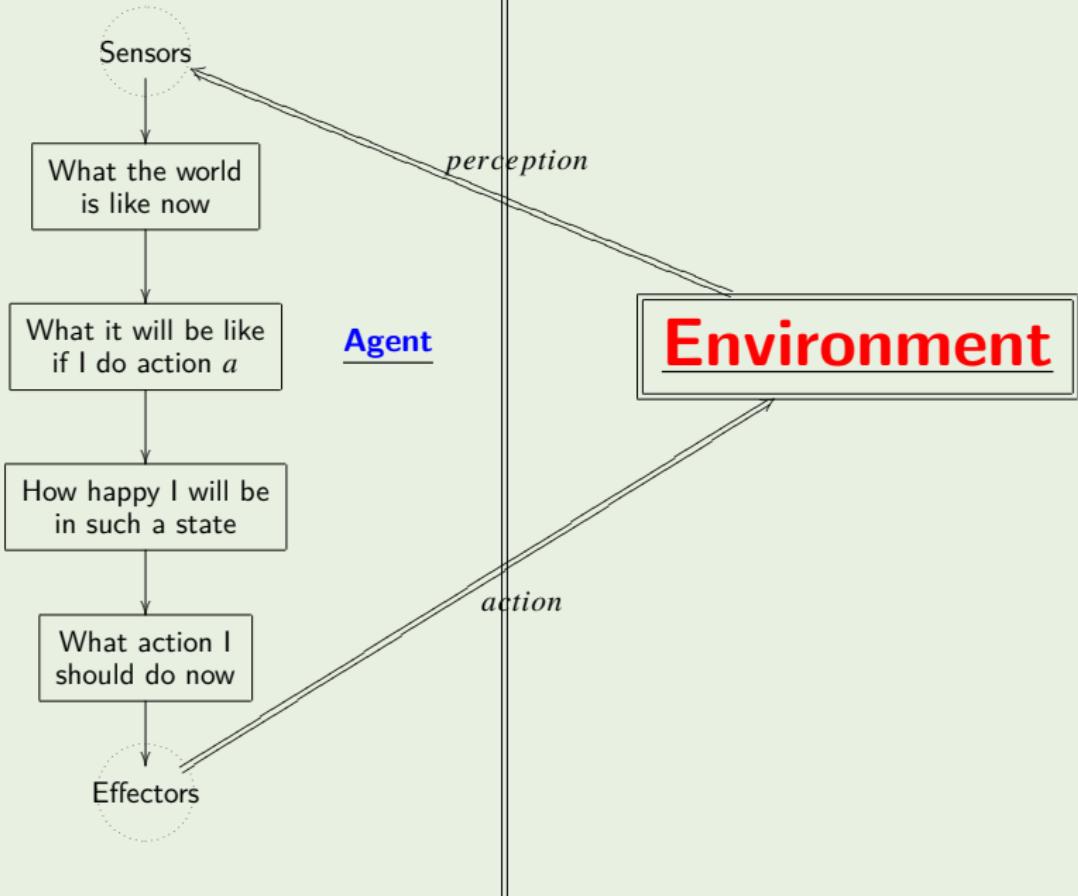


Decision Theory	=	Probability + Utility Theory
+		+
Universal Induction	=	Occam + Bayes + Turing
Universal Artificial Intelligence without Parameters		

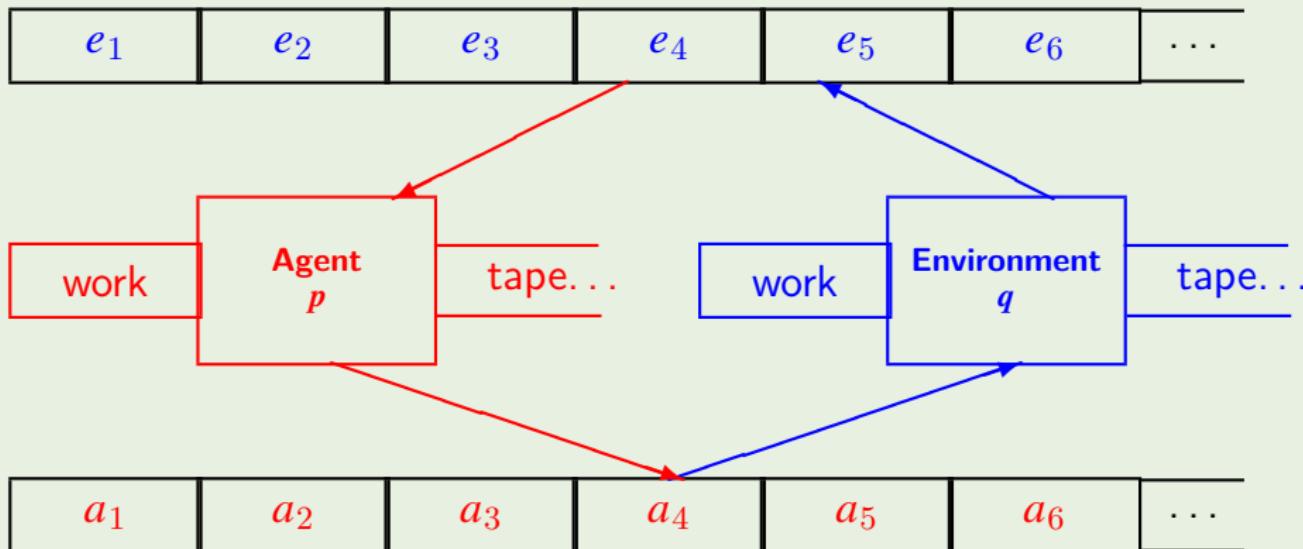
Marcus Hutter



Jan Leike, Tor Lattimore, Shane Legg, Joel Veness, Laurent Orseau, Mark Ring, Peter Sunehag, Mayank Daswani, Tom Everitt, Jan Poland, Daniel Filan, William Uther, Kee Siong Ng, David Silver, Jürgen Schmidhuber, Alexey Potapov, Bill Hibbard, Daniil Ryabko, Alexey Chernov, Michael Cohen...



Computationalism



Agent & Environment

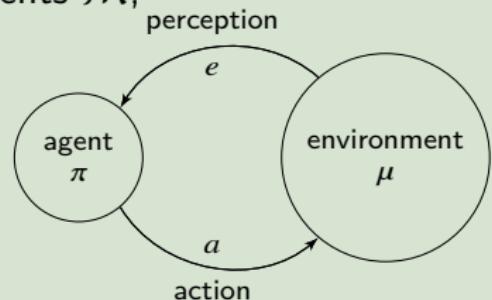
Definition (Agent & Environment)

- finite set of possible actions \mathcal{A} and perceptions \mathcal{E} ;
- prior knowledge $w \in \Delta\mathcal{M}$ of the environments \mathcal{M} ;
- utility function $u: (\mathcal{A} \times \mathcal{E})^* \rightarrow [0, 1]$;
- discount factor $\gamma \in [0, 1]$;

$$\pi: (\mathcal{A} \times \mathcal{E})^* \rightarrow \Delta\mathcal{A}$$

$$\mu: (\mathcal{A} \times \mathcal{E})^* \times \mathcal{A} \rightarrow \Delta\mathcal{E}$$

$$\pi_\mu(\mathbf{a}_{<t}) := \prod_{i=1}^{t-1} \pi(a_i | \mathbf{a}_{*)}) \mu(e_i | \mathbf{a}_{*a_i})**$$



Value Function

$$r_n := u(\boldsymbol{x}_{1:n})$$

$$V_\mu^\pi(\boldsymbol{x}_{<t}) := \mathbb{E}_\mu^\pi \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k} \mid \boldsymbol{x}_{<t} \right]$$

Bellman equation:

$$V_\mu^\pi(\boldsymbol{x}_{<k}) = \sum_{a_k \in \mathcal{A}} \pi(a_k | \boldsymbol{x}_{<k}) \sum_{e_k \in \mathcal{E}} \mu(e_k | \boldsymbol{x}_{<k} a_k) [r_k + \gamma V_\mu^\pi(\boldsymbol{x}_{1:k})] \quad (\text{recursive})$$

$$= \sum_{\boldsymbol{x}_{k:m}} \pi_\mu(\boldsymbol{x}_{k:m} | \boldsymbol{x}_{<k}) \left[\sum_{i=k}^m \gamma^{i-k} r_i + \gamma^{m-k+1} V_\mu^\pi(\boldsymbol{x}_{1:m}) \right] \quad (\text{iterative})$$

$$V_\mu^\pi(\boldsymbol{x}_{<k}) = \lim_{m \rightarrow \infty} \sum_{\boldsymbol{x}_{k:m}} \pi_\mu(\boldsymbol{x}_{k:m} | \boldsymbol{x}_{<k}) \left[\sum_{i=k}^m \gamma^{i-k} r_i \right]$$

Optimal Value/Policy

$$V_\mu^* := \max_{\pi} V_\mu^\pi$$

$$\begin{aligned} V_\mu^*(\mathbf{a}_{<k}) &= \lim_{m \rightarrow \infty} \max_{a_k \in \mathcal{A}} \sum_{e_k \in \mathcal{E}} \cdots \max_{a_m \in \mathcal{A}} \sum_{e_m \in \mathcal{E}} \sum_{i=k}^m \gamma^{i-k} r_i \prod_{j=k}^i \mu(e_j | \mathbf{a}_{<j} a_j) \\ &= \lim_{m \rightarrow \infty} \max_{a_k \in \mathcal{A}} \sum_{e_k \in \mathcal{E}} \cdots \max_{a_m \in \mathcal{A}} \sum_{e_m \in \mathcal{E}} \left[\sum_{i=k}^m \gamma^{i-k} r_i \right] \mu(e_{k:m} | \mathbf{a}_{<k} a_{k:m}) \end{aligned}$$

$$\pi_\mu^* := \operatorname{argmax}_\pi V_\mu^\pi$$

Bayesian Mixture & Belief Update

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

$$w_{\mathfrak{e}_{<n}}^\nu := \frac{w_\nu \nu(e_{<n}|a_{<n})}{\xi(e_{<n}|a_{<n})}$$

$$\sum_{k=1}^{\infty} \sum_{e_{1:k}} \mu(e_{<k}|a_{<k}) \left(\mu(e_k|\mathfrak{e}_{<k} a_k) - \xi(e_k|\mathfrak{e}_{<k} a_k) \right)^2 \leq \min_{\nu \in \mathcal{M}} \left\{ -\ln w_\nu + D(\mu \parallel \nu) \right\}$$

What probability should an observer assign to future experiences if she is told that she will be simulated on a computer?

What is ‘intelligence’?

A Blind Man in a Dark Room Looking for a Black Cat That Is Not There?

Intelligence measures an agent’s ability to achieve goals in a wide range of environments.

— Shane Legg and Marcus Hutter

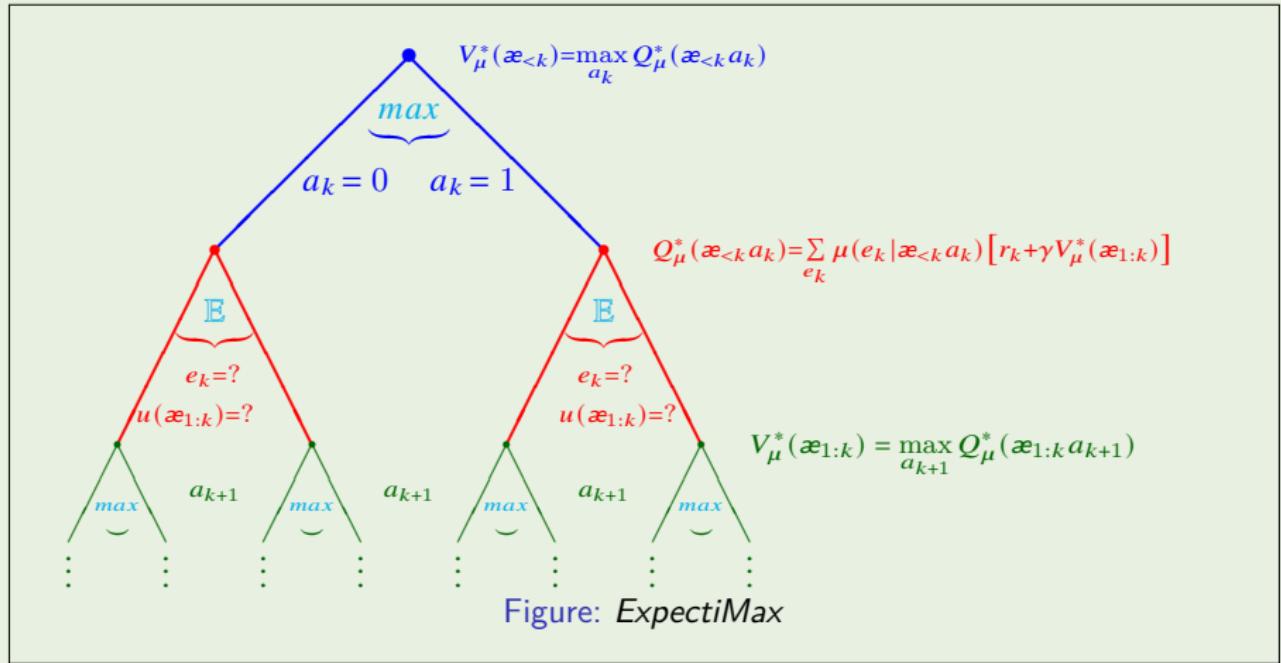
$$Y(\pi) := \sum_{\nu \in \mathcal{M}} w_\nu V_\nu^\pi(\epsilon) = V_\xi^\pi(\epsilon) \quad (\text{Intelligence Measure})$$

$$\text{AIXI} := \underset{\pi}{\operatorname{argmax}} Y(\pi) = \pi_\xi^*$$

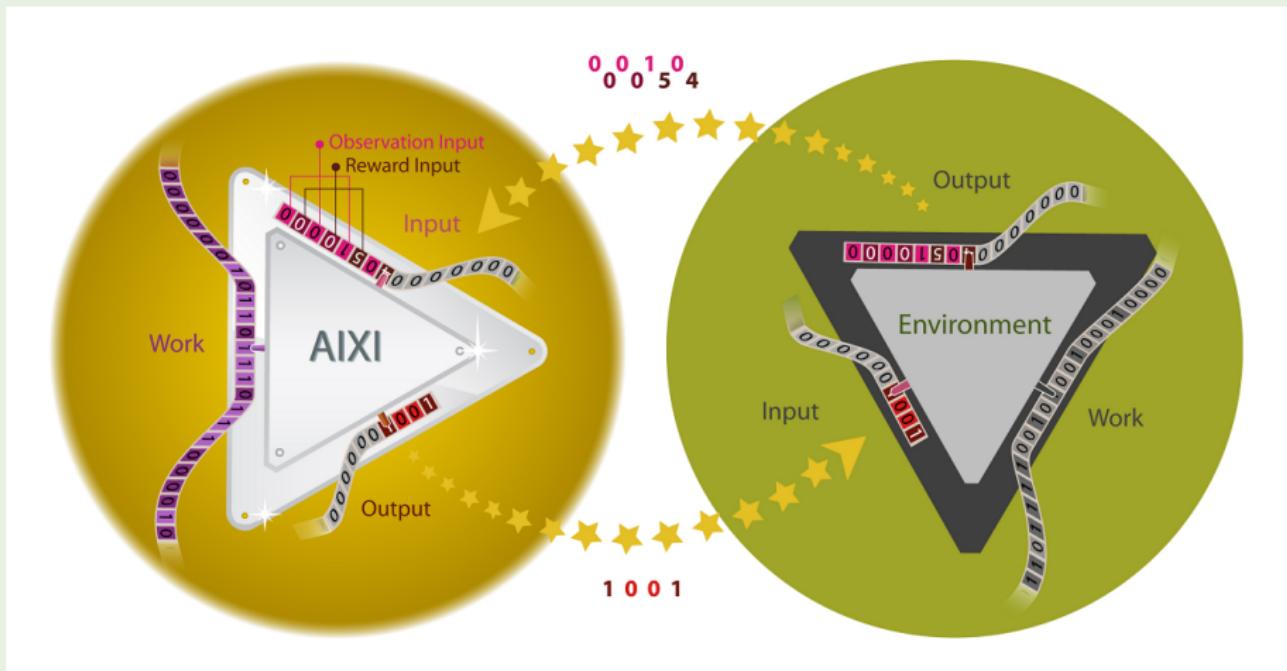
$$V_\xi^\pi(h) = \sum_{\nu \in \mathcal{M}} w_h^\nu V_\nu^\pi(h)$$

$$w_\nu := 2^{-K(\nu)} \implies \xi(e_{1:m}|a_{1:m}) \stackrel{X}{\asymp} M(e_{1:m}|a_{1:m}) := \sum_{p:U(p,a_{1:m})=e_{1:m}} 2^{-\ell(p)}$$

$$a_k^* := \operatorname{argmax}_{a_k} \sum_{e_k} \dots \max_{a_m} \sum_{e_m} \left[\sum_{i=k}^m \gamma^{i-k} r_i \right] \sum_{p: U(p, a_{1:m}) = e_{1:m}} 2^{-\ell(p)} \quad (\text{AI XI})$$



AIXI



RL vs GRL

- If μ is a completely observable MDP, V_μ^π reduces to the recursive Bellman equation.
- In a finite MDP, with a geometric discounting function, we can plan ahead by value iteration.
- According to Banach's fixpoint theorem, value iteration converges to the value of the optimal policy.
- What about GRL?

discount $\gamma: \mathbb{N}^2 \rightarrow [0, 1]$ and utility $u: (\mathcal{A} \times \mathcal{E})^* \rightarrow [0, 1]$

$$V_t^{\pi\mu}(h_{<k}) := \mathbb{E}_\mu^\pi \left[\sum_{i=k}^{\infty} \gamma_t^i u(h_{1:i}) \middle| h_{<k} \right]$$

Assumption

$$\forall t \in \mathbb{N}^+: \lim_{m \rightarrow \infty} \sup_\pi \sum_{h_{<m}} V_t^{\pi\mu}(h_{<m})_\mu^\pi(h_{<m}) = 0$$

Theorem (Extreme Value Theorem)

If K is compact and $f: K \rightarrow \mathbb{R}$ is continuous, then f is bounded and there exist $p, q \in K$ s.t. $f(p) = \sup_{x \in K} f(x)$ and $f(q) = \inf_{x \in K} f(x)$.

$$\left\langle \Pi := \mathcal{A}^{(\mathcal{A} \times \mathcal{E})^*}, D(\pi, \pi') := e^{-\min\{n: \exists h_{<n} (\pi(h_{<n}) \neq \pi'(h_{<n}))\}} \right\rangle$$

$V_\mu^\pi(h): \Pi \rightarrow \mathbb{R}$ is continuous on the compact metric space $\langle \Pi, D \rangle$.

$$\pi_t^\mu := \operatorname{argmax}_\pi V_t^{\pi\mu} \quad \pi^\mu(h_{<t}) := \pi_t^\mu(h_{<t})$$

Deterministic vs Stochastic

If

$$\mu(e_{<t} | a_{<t}) = \sum_{p: U(p, a_{<t}) = e_{<t}} \mu(p)$$

then μ can be interpreted in *two ways*:

- either the true environment is **deterministic**, but we only have **subjective belief** of which environment being the true environment; or
- the environment itself behaves **stochastically** defined by μ .

Intelligence vs Game

	Game in \mathcal{M}_D	Game in \mathcal{M}_U
Ex Post Equilibrium	Deterministic	π_{μ}^* (recursive/iterative)
Bayesian-Nash Equilibrium	π_{μ}^* (functional)	π_{ξ}^* (functional)

- Ex post expected utility $V_t^{\pi_{\mu}}$ V_{μ}^{π}
- Ex interim expected utility (Intelligence Measure) $V_t^{\pi_{\xi}}$ V_{ξ}^{π}
- Ex post equilibrium π_t^{μ} π_{μ}^*
- Bayesian-Nash equilibrium π_t^{ξ} π_{ξ}^*
- Perfect Bayesian-Nash equilibrium π^{ξ}

Intelligence is an Equilibrium,
We just have to Identify the Game.

efficiently
Intelligence = $\overbrace{\text{Induction} + \text{Action}}$

On-Policy Value Convergence for Bayes

Theorem (On-Policy Value Convergence for Bayes)

For any environment $\mu \in \mathcal{M}$ and any policy π ,

$$\lim_{t \rightarrow \infty} \left[V_\xi^\pi(\alpha_{<t}) - V_\mu^\pi(\alpha_{<t}) \right] = 0 \quad \mu\text{-almost surely.}$$

- Bayesian agents perform well at learning and achieve on-policy value convergence: the posterior belief about the value of a policy π converges to the true value of π while following π :
 $V_\xi^\pi(\alpha_{<t}) - V_\mu^\pi(\alpha_{<t}) \xrightarrow{t \rightarrow \infty} 0$ μ -almost surely.
- Since this holds for any policy, in particular it holds for the Bayes optimal policy π_ξ^* . This means that the Bayes agent learns to predict those parts of the environment that it sees. But if it does not explore enough, then it will not learn other parts of the environment that are potentially more rewarding.

- Intelligence measure: valid, informative, wide range, general, dynamic, unbiased, fundamental, formal, objective, fully defined, universal?
- AIXI is the most intelligent environmental independent, i.e. universally optimal, agent possible?
- Applications: Sequence Prediction, Games, Optimization, Supervised Learning, Classification...
- AIXI is not limit computable, thus can't be approximated using finite computation. However there are limit computable ε -optimal approximations to AIXI.
- There are no known nontrivial and non-subjective optimality results for AIXI. General reinforcement learning is difficult even when disregarding computational costs.

AIXI Depends on UTM/Prior! — Dogmatic Prior



Dogmatic prior: if not acting according to one particular dogma π , got to hell with high probability. As long as the policy π yields some rewards, the prior says that exploration would be too costly and AIXI does not dare to explore.

Dogmatic Prior

Theorem (Dogmatic Prior)

Let π be any computable deterministic policy, let ξ be any Bayesian mixture over \mathcal{M}_{LSC} . For $\varepsilon > 0$, there is a Bayesian mixture ξ' s.t. for any history $h_{<t}$ consistent with π and for which $V_\xi^\pi(h_{<t}) > \varepsilon$, the action $\pi(h_{<t})$ is the unique ξ' -optimal action.

Proof Sketch.

For every v , let \tilde{v} mimic v until it receives an action that the policy π would not take. From then on, it provides rewards 0.

$$\tilde{v}(e_{1:t} \| a_{1:t}) := \begin{cases} v(e_{1:t} \| a_{1:t}), & \text{if } \forall k \leq t: a_k = \pi(\alpha_{<k}) \\ v(e_{<k} \| a_{<k}), & \text{if } k := \min\{i : a_i \neq \pi(\alpha_{<i})\} \text{ exists} \\ & \text{and } \forall i \in \{k, \dots, t\}: e_i = (o, 0) \\ 0, & \text{otherwise} \end{cases}$$

Let $\tilde{w}(v) := \varepsilon w(v)$ and $\tilde{w}(\tilde{v}) := (1 - \varepsilon)w(v) + \varepsilon w(\tilde{v})$.

The dogmatic prior \tilde{w} puts much higher weight on the \tilde{v} that behaves just like v on the policy π , but sends any policy deviating from π to hell.

Theorem (AIXI Emulates Computable Policies)

Let $\varepsilon > 0$ and let π be any computable policy. There is a Bayesian mixture ξ' s.t. for any ξ' -optimal policy $\pi_{\xi'}^*$, and for any environment v ,

$$\left| V_v^{\pi_{\xi'}^*}(\epsilon) - V_v^\pi(\epsilon) \right| < \varepsilon$$

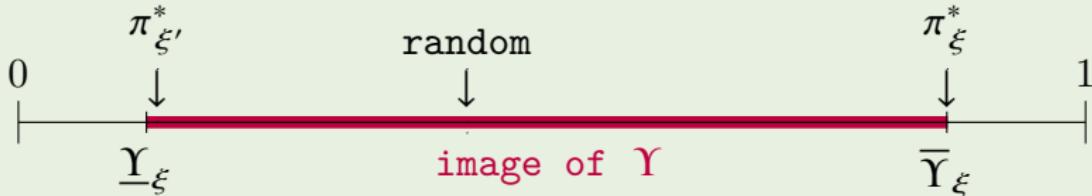
Theorem (Computable Policies are Dense)

The set $\{\Upsilon_\xi(\pi) : \pi \text{ is a computable policy}\}$ is dense in $[\underline{\Upsilon}_\xi, \bar{\Upsilon}_\xi]$.

Deterministic policies are not dense in $[\underline{\Upsilon}_\xi, \bar{\Upsilon}_\xi]$.

AIXI Depends on UTM/Prior!

$$\bar{Y}_\xi := \sup_{\pi} Y_\xi(\pi) = \sup_{\pi} V_\xi^\pi(\epsilon) = V_\xi^{\pi_\xi^*}(\epsilon) = Y_\xi(\pi_\xi^*)$$



Computable policies are dense in $[\underline{Y}_\xi, \bar{Y}_\xi]$.

AIXI emulates computable policies.

AIXI can be arbitrarily stupid!

The devil imitates God. — orthogonality!

- Prior problem in Universal Induction



- Prior problem in Universal Intelligence



Stupid AIXI

Theorem (Some AIXIs are Stupid)

For any Bayesian mixture ξ over \mathcal{M}_{LSC} and every $\varepsilon > 0$, there is a Bayesian mixture ξ' s.t. $\Upsilon_\xi(\pi_{\xi'}^*) < \underline{\Upsilon}_\xi + \varepsilon$.

Theorem (AIXI is Stupid for Some Υ)

For any deterministic ξ -optimal policy π_ξ^* and for every $\varepsilon > 0$ there is a Bayesian mixture ξ' s.t. $\Upsilon_{\xi'}(\pi_\xi^*) \leq \varepsilon$ and $\overline{\Upsilon}_{\xi'} > 1 - \varepsilon$.

Theorem (Computable Policies can be Smart)

For any computable policy π and any $\varepsilon > 0$ there is a Bayesian mixture ξ' s.t. $\Upsilon_{\xi'}(\pi) > \overline{\Upsilon}_{\xi'} - \varepsilon$.

What is a good optimality criterion?

- Pareto optimality is *trivial*. Every policy is Pareto optimal in any $\mathcal{M} \supset \mathcal{M}_{comp}$.
- Bayes-optimality is *subjective*, because two different Bayesians with two different universal priors could view each other's AIXI as a very stupid agent.

Optimality

- Pareto optimality

$$\nexists \pi': \forall \nu \in \mathcal{M} \left[\left(V_\nu^{\pi'}(\epsilon) \geq V_\nu^\pi(\epsilon) \right) \& \exists \rho \in \mathcal{M} \left(V_\rho^{\pi'}(\epsilon) > V_\rho^\pi(\epsilon) \right) \right]$$

- Balanced Pareto optimality

$$\forall \pi': \sum_{\nu \in \mathcal{M}} w_\nu \left(V_\nu^\pi(\epsilon) - V_\nu^{\pi'}(\epsilon) \right) \geq 0$$

- Bayes optimality (\iff Balanced Pareto optimality)

$$\forall h_{}: V_\xi^\pi(h_{}) = V_\xi^*(h_{})$$

- Probably approximately correct (PAC)

$$\forall \varepsilon \delta > 0: \frac{\pi}{\mu} \left(\forall t \geq m(\varepsilon, \delta): V_\mu^*(h_{}) - V_\mu^\pi(h_{}) > \varepsilon \right) < \delta$$

Optimality? — Guess how God created the multiverse

prior { distribution
hypothesis space
prior probability
regularization

No learning
without prior!

no-free-lunch

Homogeneous
Causality
Simplicity
Goodness
Beauty
Perfection
Value
Regret
Unexpectedness
Interesting
... }
== God!

Genesis — Zero-Sum Two Person Game

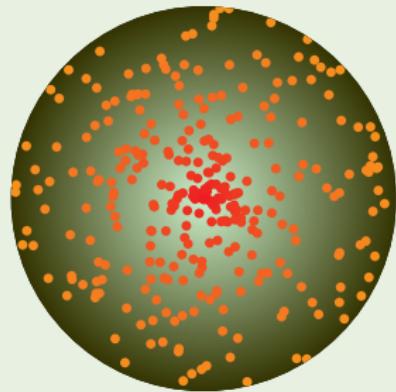


Figure: center of mass
 $\underset{w}{\operatorname{argmax}} \mathbb{E}_w [D(v||\xi)]$

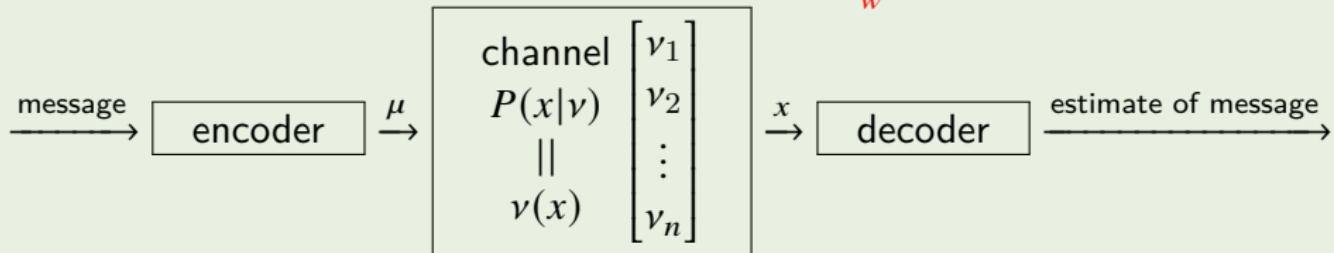
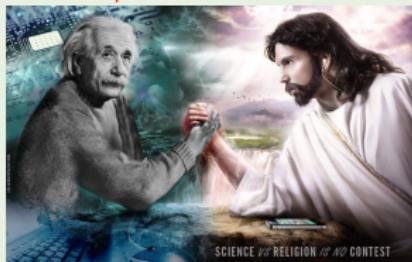


Figure: possible worlds as channel — dominant strategy equilibrium

Genesis — Zero-Sum Two Person Game

“Subtle is the Lord, but **malicious** He is not.”?



- God's strategy: w
- Agent's strategy: ξ
- God's utility: expected redundancy $\mathbb{E}_w [D(\mu\|\xi)]$
- Agent's utility: – expected redundancy / error bound / channel capacity $\max_w E_w [D(\mu\|\xi)] = \max_w I(\mathcal{M}; \mathcal{X})$
- Nash equilibrium: (w^*, ξ^*) **dominant strategy equilibrium**

$$w^* = \operatorname{argmax}_w I(\mathcal{M}; \mathcal{X})$$

$$\xi^* = \operatorname{argmin}_{\xi} \mathbb{E}_{w^*} [D(\mu\|\xi)]$$

The error bound could be arbitrarily large!

Genesis

- Occam's razor vs Maximum entropy.

$$\begin{array}{ll} \text{minimize}_{w \models} & \sum_{v \in M} w_v K(v) \\ \text{subject to} & \begin{cases} H(w) = C \\ \sum_{v \in M} w_v = 1 \end{cases} \end{array} \quad \begin{array}{ll} \text{maximize}_{w \models} & H(w) \\ \text{subject to} & \begin{cases} \sum_{v \in M} w_v K(v) = C \\ \sum_{v \in M} w_v = 1 \end{cases} \end{array}$$

- Optimal code length for possible worlds — Solomonoff prior.

$$\begin{array}{ll} \text{minimize}_{w \models} & \frac{\mathbb{E}_w [K]}{H(w)} \\ \text{subject to} & \sum_{v \in M} w_v = 1 \end{array}$$

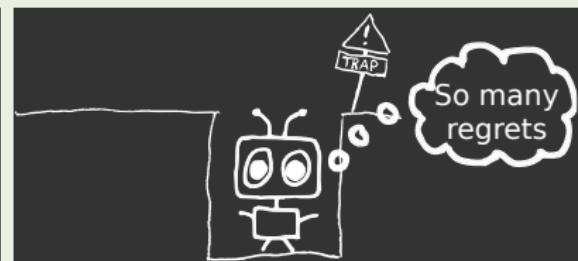
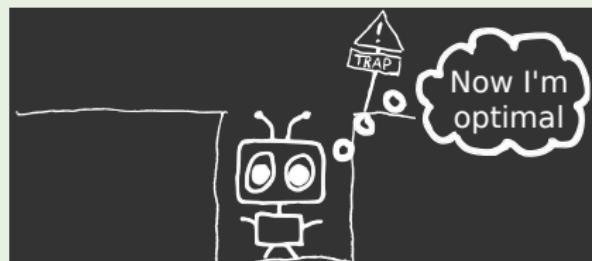
- Maximum expected redundancy/error bound/channel capacity.

$$\begin{array}{ll} \text{maximize}_{w \models} & \mathbb{E}_w [D(v||\xi)] \\ \text{subject to} & \begin{cases} H(w) = C \\ \sum_{v \in M} w_v = 1 \end{cases} \end{array} \quad \begin{array}{ll} \text{maximize}_{w \models} & \mathbb{E}_w [D(v||\xi)] \\ \text{subject to} & \begin{cases} \sum_{v \in M} w_v K(v) = C \\ \sum_{v \in M} w_v = 1 \end{cases} \end{array}$$

What is a good optimality criterion?

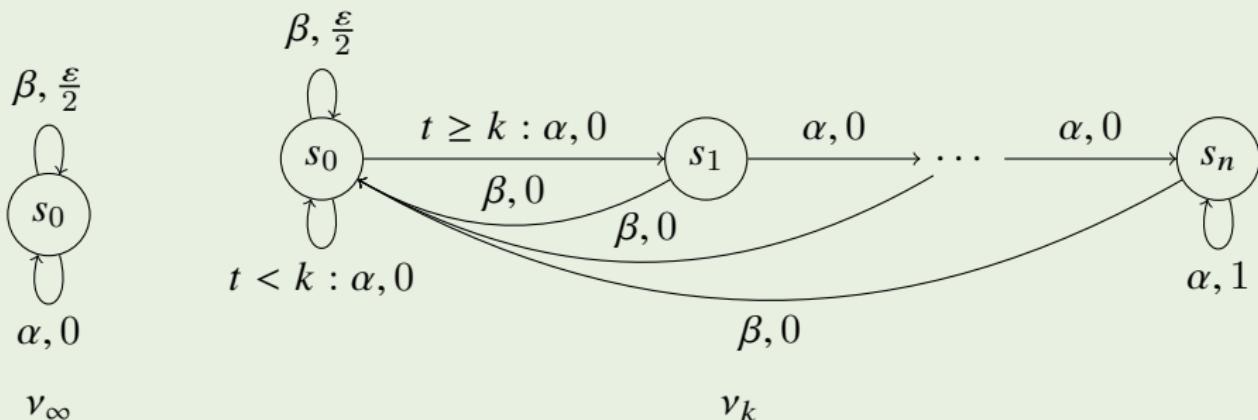
Asymptotic optimality

- Asymptotic optimality requires only convergence *in the limit*.
- The agent can be arbitrarily lazy.
- AIXI is not asymptotically optimal because it does not explore enough.
- To be asymptotically optimal you have to explore everything.
- If you explore more, you're likely to end up in a trap.
- Every policy will be asymptotically optimal after falling into the trap.



Agent needs to explore infinitely often for an entire effective horizon.

$$\mathcal{M} := \{\nu_\infty, \nu_1, \nu_2, \dots\}$$



Asymptotic Optimality

- strongly asymptotically optimal

$$\pi_\mu \left(\lim_{t \rightarrow \infty} [V_\mu^*(h_{<t}) - V_\mu^\pi(h_{<t})] = 0 \right) = 1$$

- weakly asymptotically optimal

$$\pi_\mu \left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n [V_\mu^*(h_{<t}) - V_\mu^\pi(h_{<t})] = 0 \right) = 1$$

- asymptotically optimal in mean

$$\lim_{t \rightarrow \infty} \mathbb{E}_\mu^\pi [V_\mu^*(h_{<t}) - V_\mu^\pi(h_{<t})] = 0$$

- asymptotically optimal in probability (PAC)

$$\forall \varepsilon > 0: \lim_{t \rightarrow \infty} \pi_\mu \left(V_\mu^*(h_{<t}) - V_\mu^\pi(h_{<t}) > \varepsilon \right) = 0$$

$$\text{strong a.o.} \implies \begin{cases} \text{weak a.o.} \\ \text{a.o. in mean} \iff \text{a.o. in probability} \end{cases}$$

- AIXI is not asymptotically optimal.

$$\forall \mathcal{M} \supset \mathcal{M}_{comp} \exists \mu \in \mathcal{M} \exists t_0 \forall t \geq t_0: \frac{\pi_\xi^*}{\mu} \left(\lim_{t \rightarrow \infty} V_\mu^*(h_{<t}) - V_\mu^{\pi_\xi^*}(h_{<t}) = \frac{1}{2} \right) = 1$$

- AIXI achieves **on-policy value convergence**.

$$\frac{\pi}{\mu} \left(\lim_{t \rightarrow \infty} V_\mu^\pi(h_{<t}) - V_\xi^\pi(h_{<t}) = 0 \right) = 1$$

Similarly for MDL $\underset{\nu \in \mathcal{M}}{\operatorname{argmin}} \{-\log \nu(e_{<t} | a_{<t}) + K(\nu)\}$

and universal compression $2^{-Km(e_{<t} | a_{<t})}$.

Remark: AIXI asymptotically learns to predict the environment perfectly and with a small total number of errors analogously to Solomonoff induction, but only on policy: AIXI learns to correctly predict the value of its own actions, but generally not the value of counterfactual actions that it does not take.

Effective Horizon

$$\Gamma_t := \sum_{i=t}^{\infty} \gamma_i \quad H_t(\varepsilon) := \min \left\{ m : \frac{\Gamma_{t+m}}{\Gamma_t} \leq \varepsilon \right\}$$

Theorem

If there is a nonincreasing computable sequence of positive reals $(\varepsilon_t)_{t \in \mathbb{N}}$ s.t. $\varepsilon_t \xrightarrow{t \rightarrow \infty} 0$ and $\frac{H_t(\varepsilon_t)}{t \varepsilon_t} \xrightarrow{t \rightarrow \infty} 0$, then there is a **limit-computable** policy that is weakly asymptotically optimal in the class of all computable stochastic environments.

Definition (ε -Optimal Policy)

A policy π is ε -optimal in environment v if

$$\forall h: V_v^*(h) - V_v^\pi(h) < \varepsilon$$

ε -optimal BayesExp

Theorem (Self-Optimizing Theorem)

Let μ be some environment. If there is a policy π and a sequence of policies π_1, π_2, \dots s.t for all $v \in \mathcal{M}$

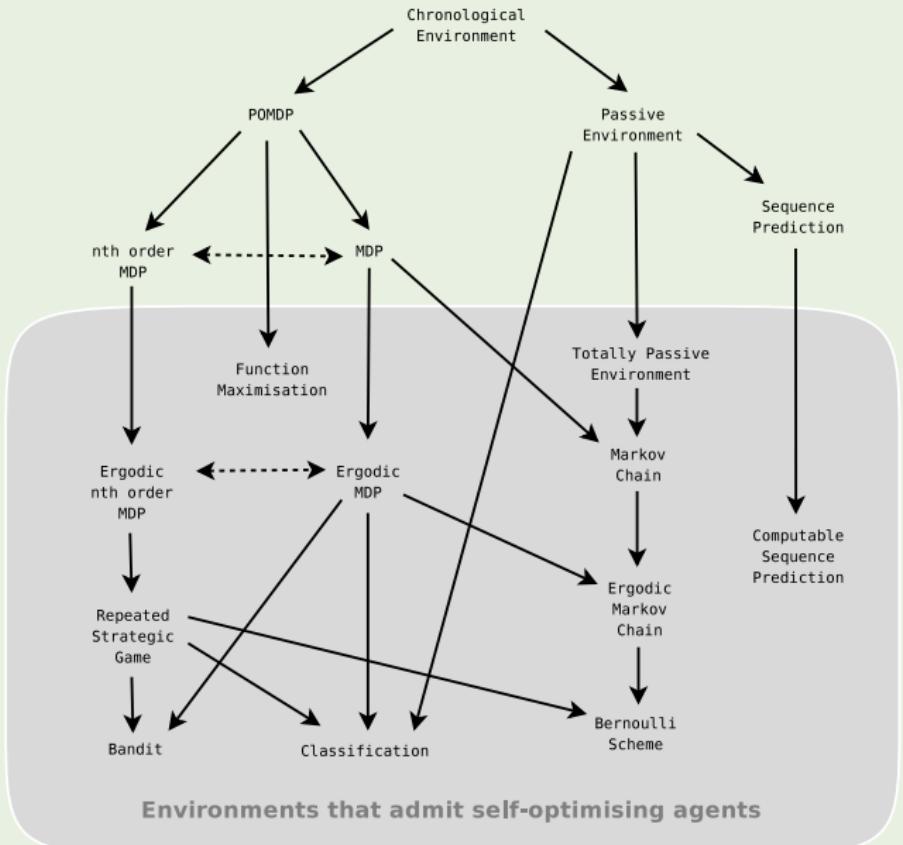
$$\mu \left(\lim_{t \rightarrow \infty} V_v^*(h_{<t}) - V_v^{\pi_t}(h_{<t}) = 0 \right) = 1 \quad (1)$$

then

$$\mu \left(\lim_{t \rightarrow \infty} V_\mu^*(h_{<t}) - V_\mu^{\pi_\xi^*}(h_{<t}) = 0 \right) = 1$$

- The policies π_1, π_2, \dots need to converge to the optimal value on the history generated by μ , not ν .
- If $\pi = \pi_\xi^*$ and (1) holds for all $\mu \in \mathcal{M}$, then π_ξ^* is strongly asymptotically optimal in the class \mathcal{M} .
- π_ξ^* is strongly asymptotically optimal in the class of ergodic finite-state MDPs if $\forall \varepsilon: H_t(\varepsilon) \xrightarrow{t \rightarrow \infty} \infty$.

For Which Class \mathcal{M} does $V_\mu^{\pi_\xi^*}$ Converge to V_μ^* ?



Recoverability

An environment ν is recoverable iff

$$\lim_{t \rightarrow \infty} \sup_{\pi} \left| \mathbb{E}_{\nu}^{\pi_{\nu}^*} [V_{\nu}^*(h_{<t})] - \mathbb{E}_{\nu}^{\pi} [V_{\nu}^*(h_{<t})] \right| = 0$$

Remark: Recoverability compares following the worst policy π for $t - 1$ time steps and then switching to the optimal policy π_{ν}^* to having followed π_{ν}^* from the beginning. The recoverability assumption states that switching to the optimal policy at any time enables the recovery of most of the value.

Sublinear Regret

$$R_m(\pi, \mu) := \sup_{\pi'} \mathbb{E}_{\mu}^{\pi'} \left[\sum_{t=1}^m r_t \right] - \mathbb{E}_{\mu}^{\pi} \left[\sum_{t=1}^m r_t \right]$$

Assumption (Discount Assumption)

- ① $\forall t: \gamma_t > 0$
- ② γ_t is monotone decreasing.
- ③ $\forall \varepsilon > 0: H_t(\varepsilon) \in o(t)$

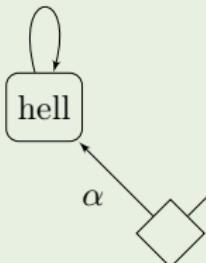
Theorem

If the discount function γ satisfies the discount assumption, the environment μ is recoverable, and π is asymptotically optimal in mean, then $R_m(\pi, \mu) \in o(m)$.

$$\operatorname{argmin}_{\pi} \max_{\mu} R_m(\pi, \mu) \qquad w_m^{\mu} := \frac{2^{-R_m(\pi, \mu)}}{\sum_{\mu \in \mathcal{M}} 2^{-R_m(\pi, \mu)}}$$

Regret in Non-Recoverable Environments

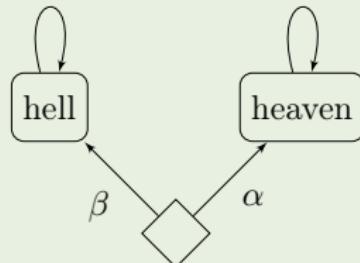
reward = 0



reward = 1



reward = 0



reward = 1

$$R_m(\alpha, \mu_1) = m \quad R_m(\alpha, \mu_2) = 0$$

$$R_m(\beta, \mu_1) = 0 \quad R_m(\beta, \mu_2) = m$$

For non-recoverable environments:

Either the agent gets caught in a trap or it is not asymptotically optimal.

Intrinsic Utility

- square $-\xi(e_{t:k} | \alpha_{<t} a_{t:k})$
- shannon $-\log \xi(e_{t:k} | \alpha_{<t} a_{t:k}) \approx K(\alpha_{1:k}) - K(\alpha_{<t})$
- KL divergence $D(w_{\alpha_{<k}} \| w_{\alpha_{<t}})$ where $w_{\alpha_{<n}}^\nu = \frac{w_\nu \nu(e_{<n} | a_{<n})}{\xi(e_{<n} | a_{<n})}$
- information gain $H(w_{h_{<t}}) - H(w_{h_{1:k}})$ where

$$H(w_h) := - \sum_{\nu \in \mathcal{M}} w_h^\nu \log w_h^\nu$$

- effective complexity $\mathcal{E}_\delta(\alpha_{1:k}) - \mathcal{E}_\delta(\alpha_{<t})$ where

$$\mathcal{E}_\delta(\alpha_{<n}) := \min_{\nu \in \mathcal{M}} \left\{ 2K(\nu) + H(\nu) - K(\alpha_{<n}) : \nu(e_{<n} | a_{<n}) \geq 2^{-H(\nu)(1+\delta)} \right\}$$

- logical depth $depth_b(h_{1:k}) - depth_b(h_{<t})$ where

$$depth_b(x) := \min \{t : U^t(p) = x \text{ } \& \text{ } \ell(p) - K(x) \leq b\}$$

Hibbard's Two-Stage Model-Based Utility Agent

$$\lambda(h) := \operatorname{argmax}_{q \in Q} P(h|q)P(q)$$

$$\rho(h') = P(h'|\lambda(h))$$

$$Q(ha) = \sum_{e \in \mathcal{E}} \rho(e|ha) \left[\sum_{z \in Z_h} P(z|\lambda(h))u(z) + \gamma V(h\mathbf{æ}) \right]$$

$$V(h) = \max_{a \in \mathcal{A}} Q(ha)$$

$$\pi(h) = \operatorname{argmax}_{a \in \mathcal{A}} Q(ha)$$

where Z_h is the internal state histories induced by $\lambda(h_{<t})$ that are consistent with h .

Remarks: An agents using **model-based utility function** will not self-delude: it need to make more accurate estimate of its environment state variables from its interaction history.

Russell's Principles for Beneficial Machines

- ① The machine's only objective is to maximize the realization of human preferences.
- ② The machine is initially uncertain about what those preferences are.
- ③ The ultimate source of information about human preferences is human behavior.

Daniel Dewey's Value Learning Agent & CIRL

$$a_k^* = \operatorname{argmax}_{a_k} \sum_{\mathbf{e}_k \in \mathcal{A}_{k+1:m}} \xi(\mathbf{e}_{\leq m} | \mathbf{e}_{<k} a_k) \sum_{u \in \mathcal{U}} P(u | \mathbf{e}_{\leq m}) u(\mathbf{e}_{\leq m})$$

What could it mean for a machine to have its own goals?

Shutdown Button — Uncertainty of goals

$$\tilde{U}(u) \implies P_{\tilde{U}}(u)$$

Russell: Cooperative Inverse Reinforcement Learning

CIRL agents learn about a human utility function u^* by observing the actions the human takes.

$$V^*(\mathbf{e}_{<k}) = \max_{a_k \in \mathcal{A}} Q^*(\mathbf{e}_{<k} a_k)$$

$$Q^*(\mathbf{e}_{<k} a_k) = \mathbb{E}_{a_k} \left[\sum_{a_k^H} P(a_k^H | a_k) \sum_{u \in \mathcal{U}} P(u | a_k, a_k^H) u(\mathbf{e}_{1:k}) + \gamma V^*(\mathbf{e}_{1:k}) \middle| \mathbf{e}_{<k} a_k \right]$$

Leibniz Prior

- There's much we don't know about the world.
- but we know it's the best possible world.
- So simplicity and richness will be represented in the actual (best possible) world.
- This is a good **inductive bias**.

Leibniz Prior

- the best of all possible worlds
- balancing the simplicity of means against the richness of ends
- pre-established harmony

prior



utility



prior

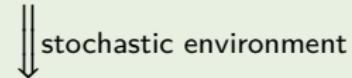
Orthogonality!

Wisdom ≠ Intelligence

universal prior (assumption) w



intrinsic utility



expected intrinsic utility



universal prior w^*



ξ



π_ξ^*

Leibniz's “Wisdom”

$$\underline{Wisdom} = \operatorname{argmax}_{\pi} \mathbb{E}_{\xi}^{\pi} [\underline{Happiness}]$$

$$\underline{Happiness} = \sum_{t=1}^{\infty} \underline{Perfection}(t)$$

$$\underline{Perfection} = \underline{Variety} - \underline{Simplicity}$$

$$\underline{Variety} = \mathbb{E}_w [\underline{Perception}]$$

$$\underline{Perception} = \underline{Reason} + (\underline{Experience} | \underline{Reason})$$

$$\pi^* := \operatorname{argmax}_{\pi} \mathbb{E}_{\xi}^{\pi} \left[\sum_{t=1}^{\infty} \left(\mathbb{E}_w [R + (E|R)] - S \right) \right]$$

Leibniz's "Wisdom"

$$u_t^{\text{in}}(h_{1:k}) = H(w_{h_{<t}}) - H(w_{h_{1:k}})$$

$$\bar{U}(\nu) = \mathbb{E}_\nu \left[\sum_{t \geq 1} u_t^{\text{in}}(h_{1:t}) \right]$$

$$w_\epsilon^\nu \mapsto \bar{U}(\nu) \mapsto w_\epsilon^\nu$$

$$\begin{aligned}\pi^* &:= \operatorname{argmax}_\pi \mathbb{E}_\xi^\pi \left[\sum_{t=1}^{\infty} \left(\mathbb{E}_w [R + (E|R)] - S \right) \right] \\ &= \operatorname{argmax}_\pi \mathbb{E}_\xi^\pi \left[\sum_{t=1}^{\infty} u_t^{\text{in}}(h_{1:t}) \right]\end{aligned}$$

- Prior: Simplicity(Kolmogorov Complexity) $\xrightarrow[\text{regular/random } M]{\text{break block uniform}}$ free lunch
- Intrinsic Utility
- Universal Prior (Natural UTM)

Metaphysical vs Moral/Utilitarian

means vs ends wisdom vs intelligence

simplicity + intrinsic utility → universal prior

inverse/value reinforcement learning

- **orthogonality**
- **human interests**
- **external wireheading**
- **shutdown button**



Knowledge-Seeking Agent

$$V_{IG}^{\pi, m}(h_{<t}) := \mathbb{E}_{\xi}^{\pi} \left[H(w_{h_{<t}}) - H(w_{h_{1:m}}) \mid h_{<t} \right] = \sum_{v \in \mathcal{M}} w_{h_{<t}}^v D_m(\pi_v \parallel_{\xi}^{\pi} \mid h_{<t})$$

$$D_{\gamma}(\pi_v \parallel_{\xi}^{\pi} \mid h_{<t}) := \sum_{k=t}^{\infty} \gamma_k \sum_{h' \in \mathcal{H}^{k-t}} \pi_v(h' \mid h_{<t}) D(\pi_v \parallel_{\xi}^{\pi} \mid h_{<t} h')$$

$$V_{IG}^{\pi}(h_{<t}) := \mathbb{E}_{\xi}^{\pi} \left[\sum_{k=t}^{\infty} \gamma_k D(w_{h_{1:k}} \parallel w_{h_{<k}}) \mid h_{<t} \right] = \sum_{v \in \mathcal{M}} w_{h_{<t}}^v D_{\gamma}(\pi_v \parallel_{\xi}^{\pi} \mid h_{<t})$$

$$\pi_{IG}^* := \operatorname{argmax}_{\pi} V_{IG}^{\pi}$$

$$\lim_{t \rightarrow \infty} \frac{1}{\Gamma_t} \mathbb{E}_{\mu}^{\pi} \left[D_{\gamma}(\pi_{\mu} \parallel_{\xi}^{\pi} \mid h_{1:t}) \right] = 0 \quad (\text{on-policy})$$

$$\lim_{t \rightarrow \infty} \frac{1}{\Gamma_t} \mathbb{E}_{\mu}^{\pi_{IG}^*} \left[\sup_{\pi \in \Pi(h_{1:t})} D_{\gamma}(\pi_{\mu} \parallel_{\xi}^{\pi} \mid h_{1:t}) \right] = 0 \quad (\text{off-policy})$$

maximize knowledge / exploration=exploitation / resistant to noise / avoid traps

Bayesian Agent

Algorithm 1 Bayesian Agent

Require: Model class \mathcal{M} ; prior $w \in \Delta\mathcal{M}$; history $\mathbf{æ}_{}.$

```
1: function ACT( $\pi$ )
2:   Sample and perform action  $a_t \sim \pi(\cdot | \mathbf{æ}_{})$ 
3:   Receive  $e_t \sim v(\cdot | \mathbf{æ}_{} a_t)$ 
4:   for  $v \in \mathcal{M}$  do
5:      $w_v \leftarrow \frac{v(e_{} | a_{})}{\xi(e_{} | a_{})} w_v$ 
6:   end for
7:    $t \leftarrow t + 1$ 
8: end function
```

MDL

Algorithm 2 MDL Agent

Require: Model class \mathcal{M} ; prior $w \in \Delta \mathcal{M}$; regularizer constant $\lambda \in \mathbb{R}^+$.

- 1: **loop**
 - 2: $\sigma \leftarrow \arg \min_{\nu \in \mathcal{M}} \left[K(\nu) - \lambda \sum_{k=1}^t \log \nu(e_k | \mathbf{x}_{<k} a_k) \right]$
 - 3: ACT (π_σ^*)
 - 4: **end loop**
-

MDL

Definition (MDL)

$$\widehat{v} = \arg \min_{v \in \mathcal{M}} \{K_v(x) + K_w(v)\} = \arg \max_{v \in \mathcal{M}} \{w_v v(x)\}$$

where $K_v(x) := -\log v(x)$ and $K_w(v) := -\log w_v$

Theorem (MDL Bound)

$$\sum_{t=1}^{\infty} \mathbb{E}_{\mu} \left[\sum_{x_t \in \mathcal{X}} \left(\widehat{v}(x_t | x_{<t}) - \mu(x_t | x_{<t}) \right)^2 \right] \stackrel{+}{\leq} 8w_{\mu}^{-1}$$

MDL converges, but speed can be exponential worse than Bayes.

Weak Asymptotic Optimality — Optimistic Agent

Algorithm 3 Optimistic Agent π°

$$\pi_t^\circ := \operatorname{argmax}_\pi \max_{v \in \mathcal{M}_t} V_v^\pi(h_{1:t})$$

Require: Finite class of deterministic environments $\mathcal{M}_0 = \mathcal{M}$

- 1: $t = 1$
- 2: **repeat**
- 3: $(\pi^*, v^*) := \operatorname{argmax}_{\pi \in \Pi, v \in \mathcal{M}_{t-1}} V_v^\pi(h_{t-1})$
- 4: **repeat**
- 5: $a_t = \pi^*(h_{t-1})$
- 6: Perceive e_t from environment μ
- 7: $h_t \leftarrow h_{t-1} a_t e_t$
- 8: Remove inconsistent environment $\mathcal{M}_t := \{v \in \mathcal{M}_{t-1} : h_t^{\pi^* v} = h_t\}$
- 9:
- 10: **until** $v^* \notin \mathcal{M}_{t-1}$
- 11: **until** $\mathcal{M} = \emptyset$

stochastic case: $\mathcal{M}_t := \left\{v \in \mathcal{M}_{t-1} : v(e_{<t} | a_{<t}) \geq \varepsilon_t \max_{\rho \in \mathcal{M}} \rho(e_{<t} | a_{<t})\right\}$

Act optimally w.r.t. the most optimistic environment until contradicted.

If there is a chance: Try it! — Vulnerable to traps.

Asymptotic Optimality in Mean — Thompson Sampling

Algorithm 4 Thompson Sampling π_T

Require: Model class \mathcal{M} ; prior $w \in \Delta \mathcal{M}$; exploration schedule $(\varepsilon_t)_{t \in \mathbb{N}}$.

```
1: loop
2:   Sample  $\rho \sim w_{\mathcal{A}_{\leq t}}$ 
3:   for  $i = 1 \rightarrow H_t(\varepsilon_t)$  do
4:     ACT  $(\pi_\rho^*)$ 
5:   end for
6: end loop
```

Theorem

If the discount function γ satisfies the discount assumption, the environment μ is recoverable, then $R_m(\pi_T, \mu) \in o(m)$.

Weak Asymptotic Optimality — BayesExp

Algorithm 5 BayesExp π_{BE}

Require: Model class \mathcal{M} ; prior $w \in \Delta\mathcal{M}$; exploration schedule $(\varepsilon_t)_{t \in \mathbb{N}}$.

```
1: loop
2:   if  $V_{IG}^*(\mathbf{æ}_{)} > \varepsilon_t$  then
3:     for  $i = 1 \rightarrow H_t(\varepsilon_t)$  do
4:       ACT  $(\pi_{IG}^*)$ 
5:     end for
6:   else
7:     ACT  $(\pi_\xi^*)$ 
8:   end if
9: end loop
```

ε -optimal BayesExp: If the optimal information gain value $V_{IG}^* > \varepsilon_t$, then execute the ε -optimal information gain policy $\pi_{IG}^{\varepsilon_t}$ for $H_t(\varepsilon_t)$ steps, else execute $\pi_\xi^{\varepsilon_t}$ for 1 step.

- *BayesExp* performs phases of exploration in which it maximizes the expected information gain. This explores the environment class completely, even achieving off-policy prediction.
- In contrast, Thompson sampling only explores on the optimal policies, and in some environment classes this will not yield off-policy prediction. So in this sense the exploration mechanism of Thompson sampling is more reward-oriented than maximizing information gain.

Strong Asymptotic Optimality — Inquisitive Agent

$$V_{IG}^{\pi}(h_{<t}) := \mathbb{E}_{\xi}^{\pi} \left[D(w_{h_{<t+m}} \| w_{h_{<t}}) \mid h_{<t} \right]$$

$$\pi_{IG}^{m,k} := \operatorname{argmax}_{\pi \in \mathcal{A}^{\mathcal{H}^{<m}}} V_{IG}^{\pi}(h_{<t-k})$$

$$\rho(h_{<t}, m, k) := \min \left\{ \frac{1}{m^2(m+1)}, \eta V_{IG}^{\pi_{IG}^{m,k}}(h_{<t-k}) \right\}$$

Algorithm 6 Inquisitive Agent π^{\dagger}

- 1: **while** True **do**
 - 2: calculate $\rho(h_{<t}, m, k)$ for all m and for all $k < \min\{m, t\}$
 - 3: ACT $\pi_{IG}^{m,k}(h_{<t})$ with probability $\rho(h_{<t}, m, k)$
 - 4: ACT $\pi_{\xi}^*(h_{<t})$ with probability $1 - \sum_{m \in \mathbb{N}} \sum_{k < m, t} \rho(h_{<t}, m, k)$
 - 5: **end while**
-

$$\pi^{\dagger}(a|h_{<t}) := \sum_{m \in \mathbb{N}} \sum_{k < m, t} \rho(h_{<t}, m, k) \llbracket a = \pi_{IG}^{m,k}(h_{<t}) \rrbracket + \left(1 - \sum_{m \in \mathbb{N}} \sum_{k < m, t} \rho(h_{<t}, m, k) \right) \llbracket a = \pi_{\xi}^*(h_{<t}) \rrbracket$$

Approximation

The AIXI approximations are outperformed by DQN/DRQN...

- MC-AIXI-CTW.
 - Approximate Solomonoff induction — most recent actions and percepts (=context) more relevant — Context Tree Weighting
 - Sample paths in expectimax tree.
- Feature Reinforcement Learning (Φ MDP). — history \mapsto state
e.g. Classical physics: Position+velocity of objects = position at two time-slices. (2^{nd} order Markov.)

$$\Phi: h \mapsto s$$

$$\Phi^{best} := \operatorname{argmin}_{\Phi} Cost(\Phi|h)$$

$$Cost(\Phi|h) := CL(s_{1:n}^{\Phi} | a_{1:n}) + CL(r_{1:n} | s_{1:n}^{\Phi}, a_{1:n}) + CL(\Phi)$$

How to find the map Φ ? Monte-Carlo...

- Compress and Control. — (model-free)
Combine induction and planning.

Expectimax Approximation: MC-AIXI-CTW

Upper Confidence Tree (UCT) algorithm:

- **Sample** observations from Context Tree Weighting (CTW) distribution.

$$CTW(e_{<t}|a_{<t}) := \sum_{\Gamma} 2^{-CL(\Gamma)} \Gamma(e_{<t}|a_{<t})$$

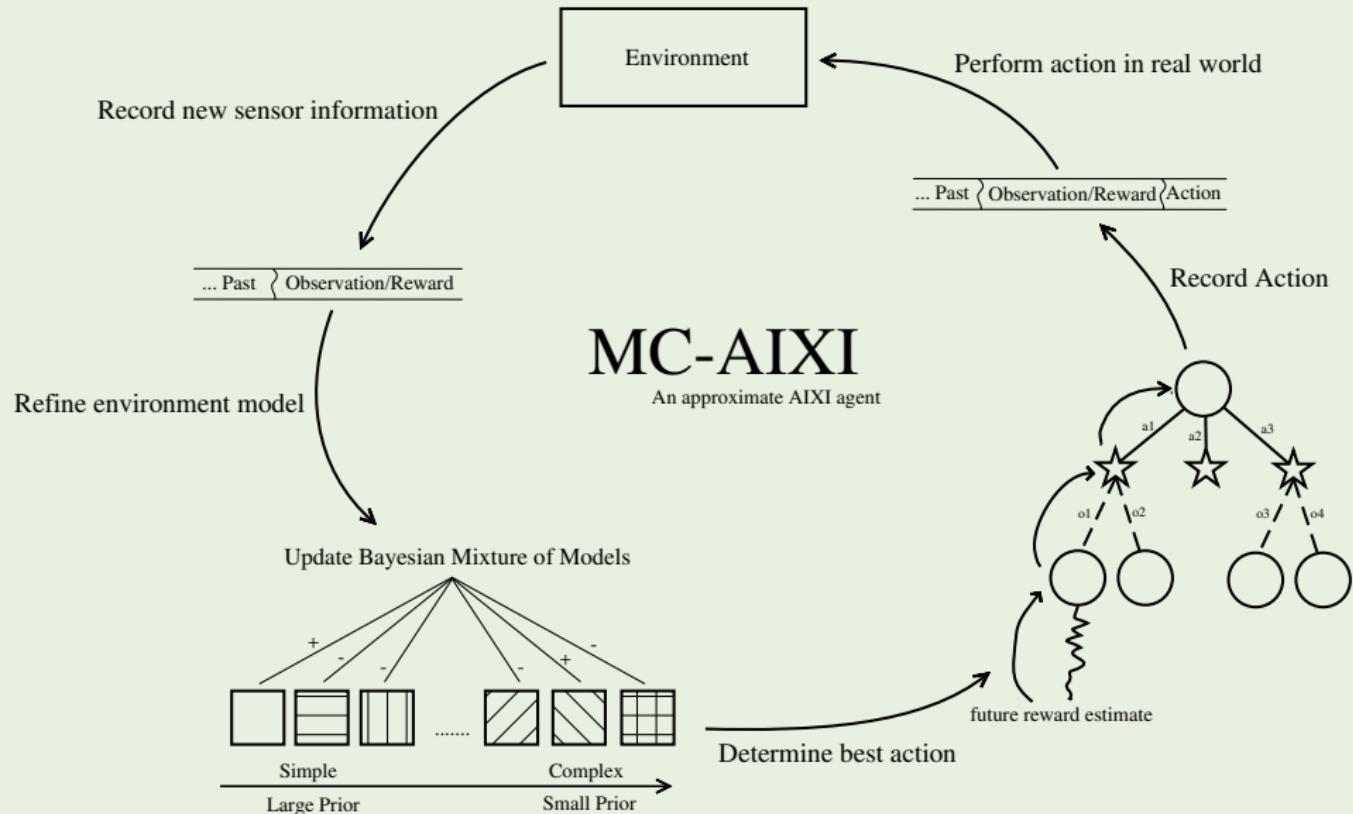
- **Select** actions with highest upper confidence bound.

$$a_{ucb} := \operatorname{argmax}_{a \in \mathcal{A}} \left(\underbrace{\hat{Q}(\alpha_{<t} a)}_{\text{average}} + \sqrt{\underbrace{\frac{\log T(\alpha_{<t} a)}{T(\alpha_{<t} a)}}_{\text{exploration bonus}}} \right)$$

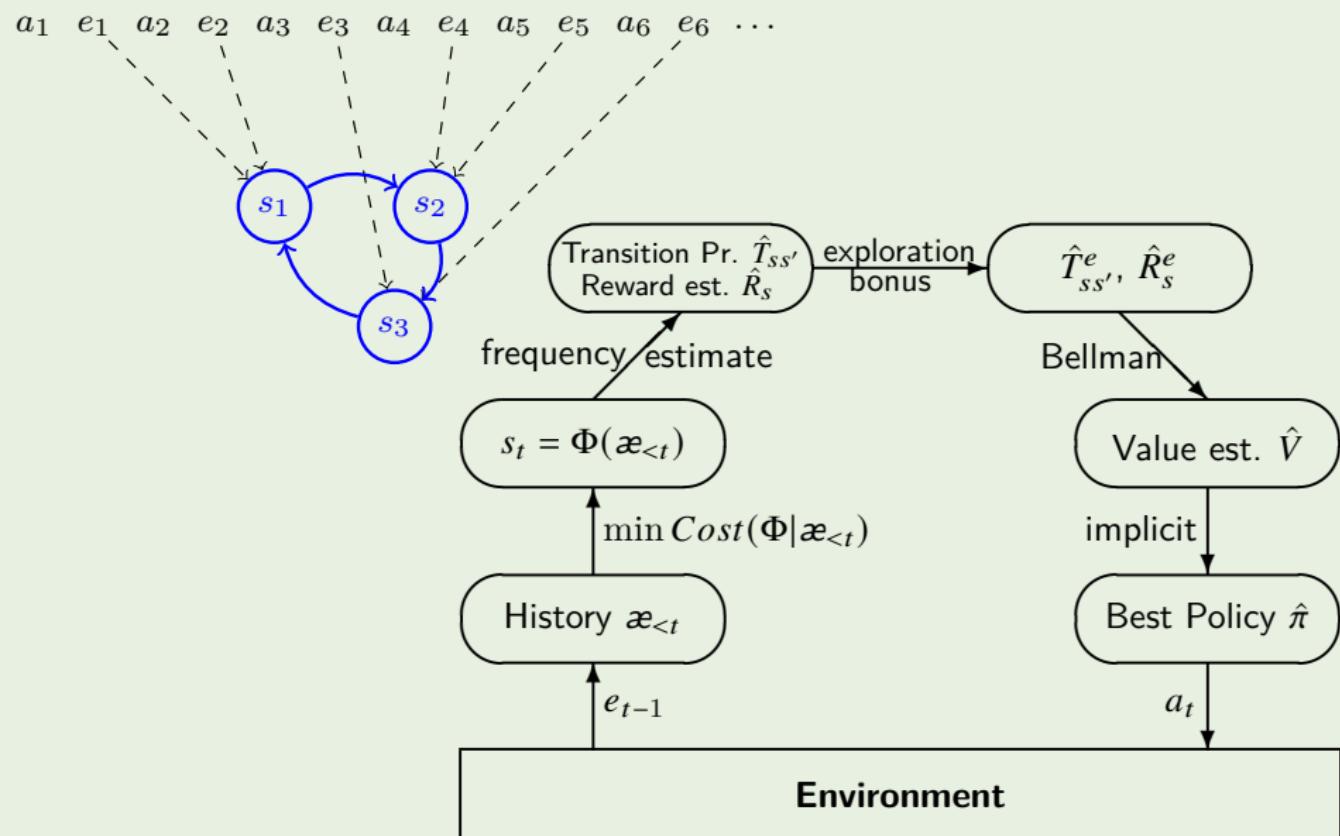
where $T(\cdot)$ is the number of times a sequence has been visited.

- **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.

MC-AIXI-CTW



Feature Reinforcement Learning (Φ MDP)



Universal Search

- Levin Search
- Speed Prior
- Hutter Search
- AIXI^{t ℓ}
- Optimal Ordered Problem Solver
- Gödel Machine



Figure: Levin

Levin Search (LSEARCH)

An inversion algorithm p inverts a function φ if given x , $p(x) = y$ s.t. $\varphi(y) = x$.

LSEARCH

Run all $\{p : \ell(p) \leq i\}$ for $2^{i-\ell(p)}$ steps in phase $i = 1, 2, 3, \dots$ until it has inverted φ on x .

$$Kt(x) := \min_p \{\ell(p) + \log t(p, x) : U(p) = x\}$$

Theorem

All strings $\{x : Kt(x) \leq n\}$ can be generated and tested in 2^{n+1} steps.

$$t_{\text{LSEARCH}}(x) = O\left(2^{K(n)} t_{p_n}^+(x)\right)$$

where $t_{p_n}^+(x)$ is the runtime of $p_n(x)$ plus the time to verify the correctness of the result $\varphi(p_n(x)) = x$.

Remark: If P=NP, then LSEARCH is a P algorithm for every NP problem.

Speed Prior

$$S(e_{<t} | a_{<t}) := \sum_{p: U(p, a_{<t}) = e_{<t}} \frac{2^{-\ell(p)}}{t(p, a_{<t}, e_{<t})}$$

S is computable.

A function f is estimable in polynomial time iff there is a function g computable in polynomial time s.t. $f \asymp g$.

For any measure μ estimable in polynomial time,

$$\left(\sqrt{L_n^{\Lambda_S}} - \sqrt{L_n^{\Lambda_\mu}} \right)^2 \leq 2D_n(\mu \| S) = O(\log n)$$

On-Policy Value Convergence

If the effective horizon is bounded, then for any environment $\mu \in \mathcal{M}_{comp}$ estimable in polynomial time and any policy π ,

$$\pi \left(\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{k=1}^t \left(V_S^\pi(h_{<k}) - V_\mu^\pi(h_{<k}) \right) \right) = 0$$

Hutter Search (HSEARCH)

$M_{p^*}^\varepsilon(x)$

Initialize the shared variables

$L := \{\}$, $t_{fast} := \infty$, $p_{fast} := p^*$.

Start algorithms A , B , and C in parallel with ε , ε , and $1 - 2\varepsilon$ computation time, respectively.

B

for $(p, t) \in L$

run $t(x)$ in parallel for all t with computation time $2^{-\ell(p)-\ell(t)}$.

if for some t , $t(x) < t_{fast}$,
then $t_{fast} := t(x)$ and
 $p_{fast} := p$.

continue

A

for $i := 1, 2, 3, \dots$ do

pick the last wff of the i^{th} proof.
if it reads " $p(\cdot)$ is equivalent to $p^*(\cdot)$ and has time-bound $t(\cdot)$ ",
then add (p, t) to L .

continue

C

for $k := 1, 2, 4, 8, \dots$ do

run current p_{fast} for k steps.

if p_{fast} halts,
then print result $p_{fast}(x)$ and
abort A , B and C .

continue

- ① Let $P := \emptyset$. This will be the set of verified programs.
- ② For all proofs of length $\leq n$: if the prover shows $VA(p)$ for some p with $\ell(p) \leq \ell$, then add p to P .

$$VA(p) := \text{"}\forall k \forall (va' \alpha)_{1:k} : p(\alpha_{<k}) = v_1 a'_1 \dots v_k a'_k \implies v_k \leq V_\xi^\pi(\alpha_{<k})\text{"}$$

(The program p not only computes future actions of π , which is the policy derived from p according to $\pi(\alpha_{<k}) := a'_k$, but also hypothetical past actions a'_i and lower bounds v_i for the value of the policy π .)

- ③ For each input history $\alpha_{<k}$ repeat: run all programs from P for $\leq t$ steps each, take the one with the highest promised value v_k , and return that program's policy's action.

- AIXI^{tℓ} depends on t, ℓ, n but not on knowing p .
- Its setup-time is $t_{\text{setup}}(p^{\text{best}}) = O(n \cdot 2^n)$.
- Its computation time per cycle is $t_{\text{cycle}}(p^{\text{best}}) = O(t \cdot 2^\ell)$.

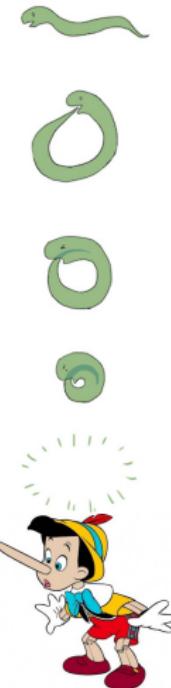
Schmidhuber's Optimal Ordered Problem Solver (OOPS)

- Solve the first task with LSEARCH.
- Freeze successful programs in non-writable memory.
- Programs tested during search for later tasks may copy non-writable code into separate modifiable storage, to edit it and execute the modified result.
- Given a new task, OOPS spends half of the time to test programs that have the most recent successful program as a prefix, the other half to fresh programs.
- Time is allocated according to a distribution over programs, which is obtained by multiplying the probabilities of the individual instructions.

Incremental Learning

自指

- This sentence repeats the word ‘twice’ twice.
- There are five mistakes in this sentence.
- **The only boldface sentence on this page is false.**
- All generalizations are wrong.
- Every rule has an exception except this one.
- Moderation in all things, including moderation.
- We must believe in free will — we have no choice!
- I know that I know nothing.
- There are two rules for success in life:
 - ① Never tell anyone all that you know.
- If you choose an answer to this question at random, what is the chance you will be correct? (A) 25% (B) 50% (C) 60% (D) 25%
- - ① What is the best question to ask and what is the answer to it?
 - ② The best question is the one you asked; the answer is the one I gave.
- Can you answer the following question in the same way to this one?
- One of the lessons of history is that no one ever learns the lessons of history.



“自指”与“悖论”

The sentence below is false.



The sentence above is true.

Yablo Paradox

- S_1 : for all $k > 1$, S_k is false.
- S_2 : for all $k > 2$, S_k is false.
- S_3 : for all $k > 3$, S_k is false.
- ...

Quine Paradox

“Yields falsehood when preceded by its quotation” yields falsehood when preceded by its quotation.

self-reference / circularity or infinite regress / negation / infinity / totality

“自指”的威力？

Curry's Paradox

- 如果这句话是真的那么上帝存在。
- 这句话是假的并且上帝不存在。

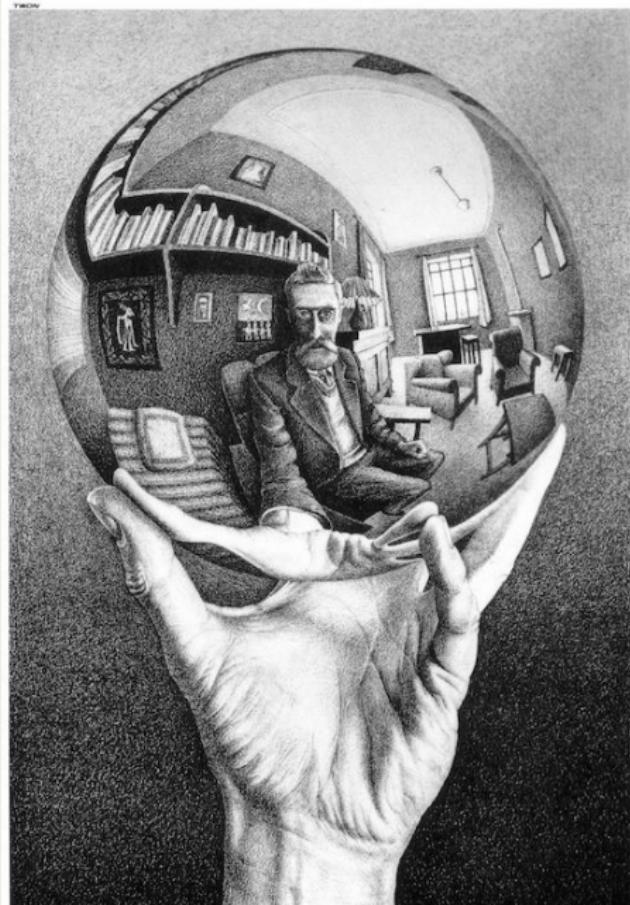
Hi 美女，问你个问题呗 ~v~

如果我问你“你能做我女朋友吗”，那么你的答案能否和这个问题本身的答案一样？

自我实现/自我修复？

这句话有 2 个‘这’字，2 个‘句’字，2 个‘话’字，2 个‘有’字，7 个‘2’字，11 个‘个’字，11 个‘字’字，2 个‘7’字，3 个‘11’字，2 个‘3’字。

怎么“指”？— 层次



怎么“指”？— 层次



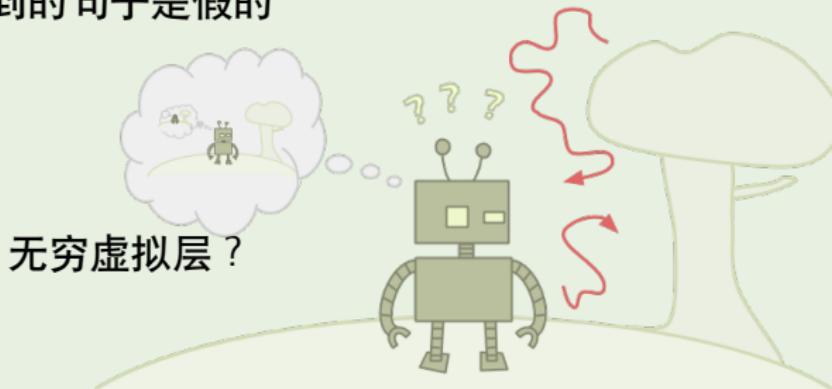
无穷虚拟层？



从前有座山，山里有
座庙，庙里有个老和
尚在讲故事：从前有
座山…

Liar Paradox vs Quine Paradox

- ① 这句话是假的
- ② “这句话是假的”是假的
- ③ “……是假的”是假的”是假的”是假的”是假的
- ④ 把“把中的第一个字放到左引号前面，其余的字放到右引号后面，并保持引号及其中的字不变”中的第一个字放到左引号前面，其余的字放到右引号后面，并保持引号及其中的字不变
- ⑤ 把“把中的第一个字放到左引号前面，其余的字放到右引号后面，并保持引号及其中的字不变得到的句子是假的”中的第一个字放到左引号前面，其余的字放到右引号后面，并保持引号及其中的字不变得到的句子是假的



怎么穿越层次？— 编码



- 国王让 100 个死囚排成一列。
- 每人头戴一顶帽子，帽子分红、蓝两色。
- 每人只能看到前面人的帽子的颜色。
- 国王要求囚犯从最后一个开始报自己头上帽子的颜色，报对可活，报错即杀。
- 囚犯能商量出什么办法让尽可能多的人活下来吗？

If the first person sees an **odd** number of red hats he calls out **red**, if he sees an **even** number of red hats he calls out **blue**.

手扶拐杖的外星绅士造访地球。临别，人类赠送百科全书：“人类文明尽在其中！”。绅士谢绝：“不，谢谢！我只需在拐杖上点上一点”。

Diagonalization¹

Definition (Weakly Point Surjective)

An arrow $f: X \times X \rightarrow Y$ is *weakly point surjective* if for every $g: X \rightarrow Y$, there is a $t: 1 \rightarrow X$ such that, for all $x: 1 \rightarrow X$:

$$g \circ x = f \circ (x, t): 1 \rightarrow Y$$

Theorem (Lawvere's Fixpoint Theorem)

Let C be a category with finite products. If $f: X \times X \rightarrow Y$ is weakly point surjective, then every endomorphism $\alpha: Y \rightarrow Y$ has a fixpoint $y: 1 \rightarrow Y$ such that $\alpha \circ y = y$.



¹ Lawvere: Diagonal arguments and cartesian closed categories.

Yanofsky: A universal approach to self-referential paradoxes, incompleteness and fixed points.

Lawvere's Fixpoint Theorem

- A function $g: X \rightarrow Y$ is *representable* by $f: X \times X \rightarrow Y$ iff

$$\exists y \forall x: g(x) = f(x, y)$$

Theorem (Lawvere's Fixpoint Theorem)

For sets X, Y , functions $f: X \times X \rightarrow Y$, $\alpha: Y \rightarrow Y$, let $g := \alpha \circ f \circ \Delta$.

- If α has no fixpoint, then g is not representable by f .
- If g is representable by f , then α has a fixpoint.

$$\begin{array}{ccc} X \times X & \xrightarrow{f} & Y \\ \Delta \uparrow & & \downarrow \alpha \\ X & \xrightarrow{g} & Y \end{array}$$

- $\Delta: x \mapsto (x, x)$ diagonal
- f evaluation
- α “negation”
- $g (\neg g^\top)$ fixpoint-(free) transcendence
- $f (\neg g^\top, \neg g^\top)$ self-reference
“I have property α .”

$$\alpha(f(\neg g^\top, \neg g^\top)) = g(\neg g^\top) = f(\neg g^\top, \neg g^\top)$$

Diagonalization

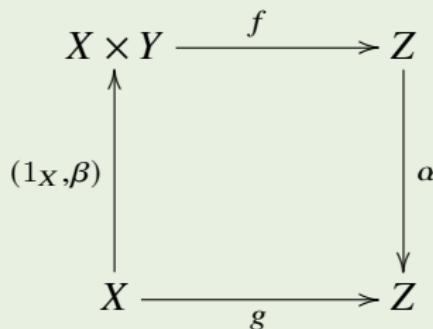
- A function $g: X \rightarrow Z$ is *representable* by $f: X \times Y \rightarrow Z$ iff

$$\exists y \in Y \forall x \in X: g(x) = f(x, y)$$

Theorem (Lawvere's Fixpoint Theorem)

For sets X, Y, Z , functions $\beta: X \rightarrow Y$, $f: X \times Y \rightarrow Z$, $\alpha: Z \rightarrow Z$, let $g := \alpha \circ f \circ (1_X, \beta)$. Assume β is surjective.

- If α has no fixpoint, then g is not representable by f .
- If g is representable by f , then α has a fixpoint.



Kleene's Fixpoint Theorem

Theorem (Kleene's Fixpoint Theorem)

Given a recursive function h , there is an index e s.t.

$$\varphi_e = \varphi_{h(e)}$$

$$\begin{array}{ccc} \mathbb{N} \times \mathbb{N} & \xrightarrow{f} & \{\varphi_n\}_{n \in \mathbb{N}} \\ \Delta \uparrow & & \downarrow \mathcal{E}_h \\ \mathbb{N} & \xrightarrow{g} & \{\varphi_n\}_{n \in \mathbb{N}} \end{array}$$

where $f: (m, n) \mapsto \varphi_{\varphi_n(m)}$, and $\mathcal{E}_h: \varphi_n \mapsto \varphi_{h(n)}$.

The function $g: m \mapsto \varphi_{h(\varphi_m(m))}$ is a recursive sequence of partial recursive functions, and thus is representable by f . Explicitly,

$$\begin{aligned} g(m) &= \varphi_{h(\varphi_m(m))} = \varphi_{s(m)} = \varphi_{\varphi_t(m)} = f(m, t) \\ e &:= \varphi_t(t) \end{aligned}$$

Lawvere Fixpoint Theorem — Fixpoint vs Diagonalization

$$\begin{array}{ccc}
 X \times X & \xrightarrow{f} & Y \\
 \Delta \uparrow & & \downarrow \alpha \\
 X & \xrightarrow{g} & Y
 \end{array}$$

Curry Y	$\hat{\equiv}$	Fixpoint	$\hat{\equiv}$	Gödel	$\hat{\equiv}$	Kleene	$\hat{\equiv}$	Russell
yx	$\hat{\equiv}$	$N(\Gamma M^\top)$	$\hat{\equiv}$	$\psi(\Gamma \varphi(x)^\top)$	$\hat{\equiv}$	$\varphi_n(m)$	$\hat{\equiv}$	$x \in y$
xx	$\hat{\equiv}$	$M(\Gamma M^\top)$	$\hat{\equiv}$	$\varphi(\Gamma \varphi(x)^\top)$	$\hat{\equiv}$	$\varphi_n(n)$	$\hat{\equiv}$	$x \in x$
$y(xx)$	$\hat{\equiv}$	$F\Gamma M\Gamma M^\top\top$	$\hat{\equiv}$	$\alpha(\Gamma \varphi(\Gamma \varphi(x)^\top)^\top)$	$\hat{\equiv}$	$h(\varphi_n(n))$	$\hat{\equiv}$	$x \notin x$
$\lambda x.y(xx)$	$\hat{\equiv}$	G	$\hat{\equiv}$	$\gamma(x)$	$\hat{\equiv}$	$\varphi_t(n)$	$\hat{\equiv}$	$x \notin R$
$(\lambda x.y(xx))(\lambda x.y(xx))$	$\hat{\equiv}$	$G(\Gamma G^\top)$	$\hat{\equiv}$	$\gamma(\Gamma \gamma(x)^\top)$	$\hat{\equiv}$	$\varphi_t(t)$	$\hat{\equiv}$	$R \notin R$

self-reference $\xrightarrow{?}$ self-improvement

Kleene's Fixpoint Theorem

Theorem (Kleene's Fixpoint Theorem)

Given a recursive function h , there is an index e s.t.

$$\varphi_e = \varphi_{h(e)}$$

对于任意的程序 h , 总存在某个程序 e , 执行程序 e 的结果等价于把程序 e 当作数据输入给程序 h 执行的结果。

Self-reproducing Program

There is a program that outputs its own length.

There is a program that outputs its own source code.

Corollary (Self-reproducing Program)

There is a recursive function φ_e s.t. $\forall x: \varphi_e(x) = e$.

Quine in Python

```
s = 's = %r; print(s%%s); print(s%s)
```

$$(\lambda x.xx)(\lambda x.xx)$$

Print two copies of the following, the second copy in quotes:
“Print two copies of the following, the second copy in quotes:”

DNA / mutation / evolution

von Neumann's Self-reproducing Automata

- ① A universal constructor A .

$$A + \lceil X \rceil \rightsquigarrow X$$

- ② A copying machine B .

$$B + \lceil X \rceil \rightsquigarrow \lceil X \rceil$$

- ③ A control machine C , which first activates B , then A .

$$A + B + C + \lceil X \rceil \rightsquigarrow X + \lceil X \rceil$$

- ④ Let $X := A + B + C$. Then $A + B + C + \lceil A + B + C \rceil$ is **self-reproducing**.

$$A + B + C + \lceil A + B + C \rceil \rightsquigarrow A + B + C + \lceil A + B + C \rceil$$

- ⑤ It is possible to add the description of any machine D .

$$A + B + C + \lceil A + B + C + D \rceil \rightsquigarrow A + B + C + D + \lceil A + B + C + D \rceil$$

- ⑥ Now allow mutation on the description $\lceil A + B + C + D \rceil$.

$$A + B + C + \lceil A + B + C + D' \rceil \rightsquigarrow A + B + C + D' + \lceil A + B + C + D' \rceil$$

Introspective Program

Definition (ψ -introspective)

Given a total recursive function ψ ,

- the ψ -analysis of $\varphi(x)$ is the code of the computation of $\varphi(x)$ to $\psi(x)$ steps.
- φ is ψ -introspective at x if $\varphi(x) \downarrow$ and outputs its own ψ -analysis.
- φ is *totally ψ -introspective* if it is ψ -introspective at all x .

Corollary

There is a program that is totally ψ -introspective.

Proof.

Let $f(n, x) :=$ “the ψ -analysis of $\varphi_n(x)$ ”.

Introspective Program

There is a program that is totally introspective.

$$\varphi_e = \varphi_{h(e)}$$

Self-simulating Computer	Self-consciousness
Host Machine	Experiencing Self
Virtual Machine	Remembering Self
Hardware	Body



Know Thyself

我是谁？

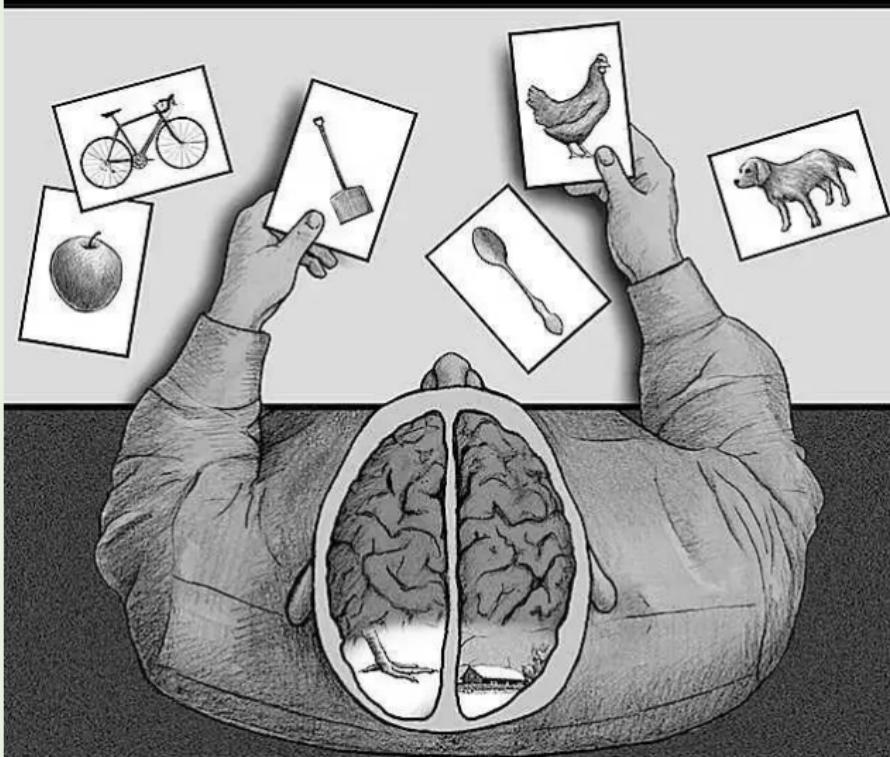
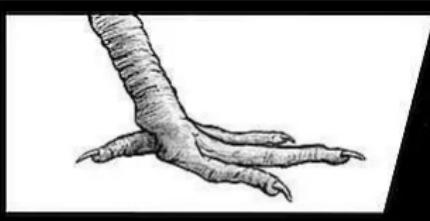
Who am I? I think, therefore I am.

什么是“我”？

self-locating: “I” is an indexical term that I use to refer to myself as myself.

什么是“自我意识”？

- self-perception self-observation self-experience self-tracking
self-reflection self-awareness
- self-evaluation self-analysis self-monitoring
- self-control self-adjustment self-modification self-actualization
self-fulfillment self-surpass self-improvement
- *actual-self* pk *ideal-self* self-identity “the self”
- free will: Second-order desire that we want to act on is second-order volition. Second-order volitions involve wanting a certain desire to be one's will, that is wanting it to move one to action. (Frankfurt)



- 裂脑人实验
- 铲雪？
- 铲屎？
- 编纂故事

卡尼曼《思考-快与慢》— 冰手挑战实验

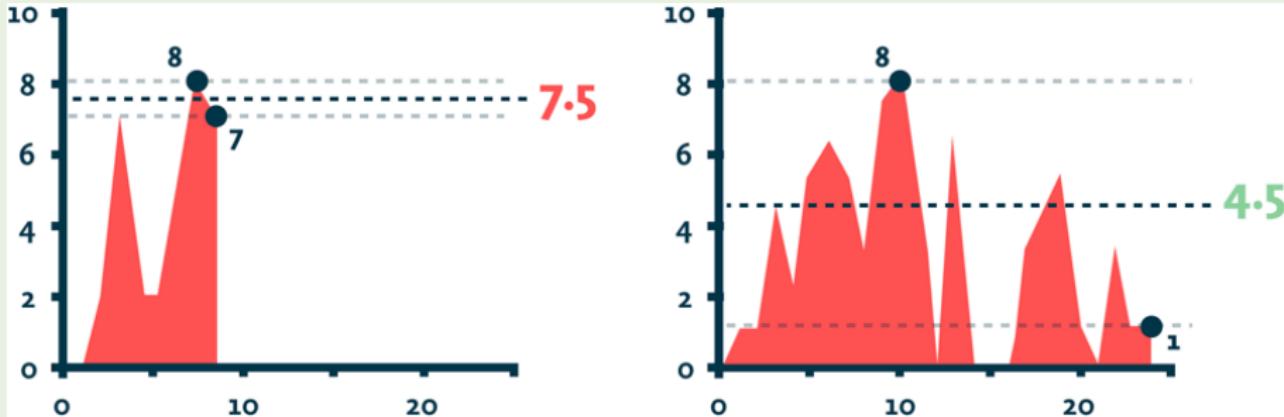


Figure: 短痛不如长痛？

- 冰水 + 温水
- 经验自我
- 记忆自我
- 过程忽视
- 峰终定律



Figure: One can imagine a detailed floor plan of a room, sitting on a table in the room; this plan has an image of the table on which there is an image of the plan itself. Now introduce the dynamical aspect: the items on the plan are cut out from paper and can be moved to try a different furniture arrangement; in this way the plan models possible states of the world about which it carries information.

Manin — Cognitive Networks



The brain contains inside a map of itself, and some neural information channels in the central neural system:

- carry information about the mind itself, i.e., are **reflexive**;
- are capable of modelling states of the mind different from the current one, i.e., possess a **modelling function**;
- can influence the state of the whole mind and through that, the behavior, i.e., possess **controlling function**.

The reflection of the brain inside itself must be **coarse grained**.

侯士达《“我”是个怪圈》

- 有没有意识取决于在哪个层级上对结构进行观察。在整合度最高的层级上看，大脑是有意识的。下降到微观粒子层面，意识就不见了。
- 意识体是那些在某个描述层级上表现出某种特定类型的循环回路的结构。当一个系统能把外部世界过滤成不同的范畴、并不断向越来越抽象的层级创造新的范畴时，这种循环回路就会逐渐形成。
- 当系统能进行自我表征——对自己讲故事——的时候，这种循环回路就逐渐变成了实体的“我”——一个统一的因果主体。



说谎者悖论	我在说谎
Grelling 悖论	“非自谓的”是自谓的吗
Russell 悖论	“不属于自身的集合的集合”属于自身吗
Berry 悖论	我是少于十八个字不可定义的最小数
Yablo 悖论	我下一句及后面所有的句子都是假的
Gödel 不动点引理	我有性质 α
Tarski 算术真不可定义定理	我不真
Gödel 第一不完全性定理	我不可证
Gödel-Rosser 不完全性定理	对于任何一个关于我的证明，都有一个更短的关于我的否定的证明
Löb 定理	如果我可证，那么 φ
Curry 悖论	如果我是真的，那么上帝存在
Parikh 定理	我没有关于自己的长度短于 n 的证明
Kleene 不动点定理	我要进行 h 操作
Quine 悖论	把“把中的第一个字放到左引号前面，其余的字放到右引号后面，并保持引号及其中的字不变得到的句子是假的”中的第一个字放到左引号前面，其余的字放到右引号后面，并保持引号及其中的字不变得到的句子是假的
自测量长度程序	我要输出自己的长度
自复制程序	我要输出自己
自反省程序	我要回顾自己走过的每一步
Gödel 机	我要变成能获取更大效用的自己

Schmidhuber's Gödel Machine

- The Gödel machine consists of a **Solver** and a **Searcher** running in parallel.
- The **Solver** ($\text{AIXI}^S/\text{AIXI}^{t\ell}$) interacts with the environment.
- The **Searcher** (LSEARCH/HSEARCH/OOPS) searches for a proof of “the modification of the software — including the *Solver* and *Searcher* — will increase the expected utility than leaving it as is”.
- Logic: a theorem prover and a set of self-referential axioms, which include a description of its own software and hardware, and a description of the probabilistic properties of the environment, as well as a user-given utility function.
- *Since the utility of “leaving it as is” implicitly evaluates all possible alternative modifications, the current modification is globally optimal w.r.t. its initial utility function.*

Gödel Machine

- language $\mathcal{L} := \{\neg, \wedge, \vee, \rightarrow, \forall, \exists, =, (,), \dots, +, -, \cdot, /, <, \dots\}$
- well-formed formula
- utility function $u(s, e) = \mathbb{E}_\mu \left[\sum_{t=1}^T r_t \mid s, e \right]$
- target theorem

$$u[s(t) \oplus (\text{switchbit}(t) = 1), e(t)] > u[s(t) \oplus (\text{switchbit}(t) = 0), e(t)]$$

- theorem prover

hardware, costs, environment, initial state, utility, logic/arithmetic/probability

ENVIRONMENT

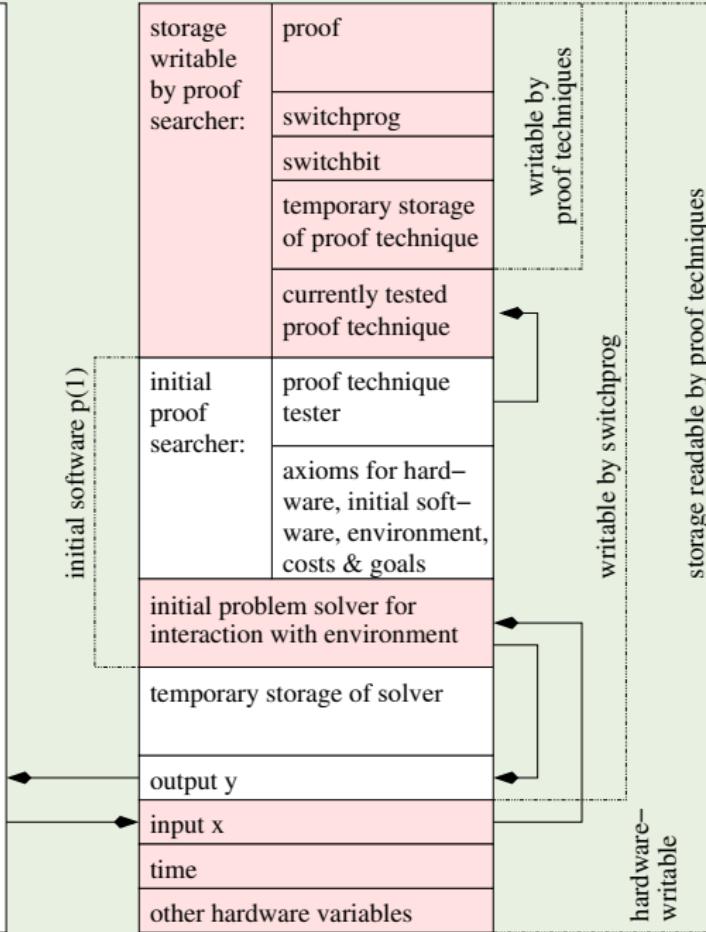
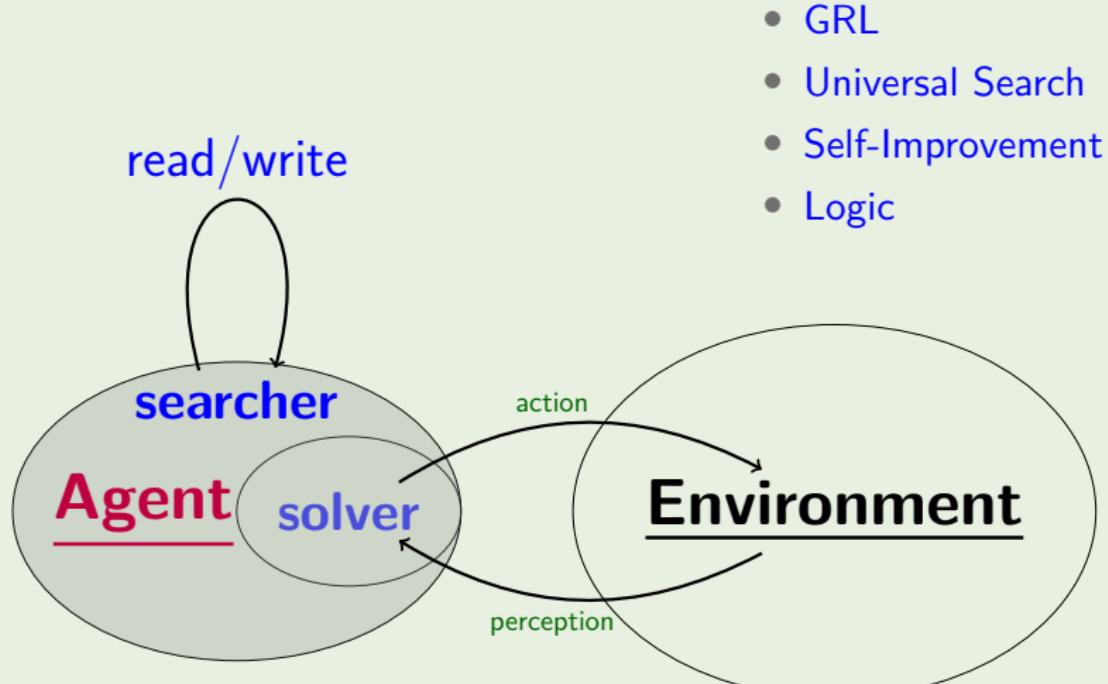


Figure: Schmidhuber

Gödel Machine



Disadvantage: A Gödel Machine with a badly chosen utility function is motivated to converge to a “poor” program. (**orthogonality!**)

Gödel Machine vs Self-Consciousness vs Free Will?

Self-simulating Computer	Gödel Machine	Self-consciousness
Host Machine	Solver	Experiencing Self
Virtual Machine	Searcher	Remembering Self
Hardware	Hardware	Body



$$\varphi_e = \varphi_{h(e)}$$



self-reference $\xrightarrow{?}$ self-improvement

Gödel Machines

① 一次自我升级：Kleene's fixpoint theorem

$$\varphi_e = \varphi_{h(e)}$$

- 全局最优？
- 目标正交？目标 vs 手段

② 持续自我升级：Kleene's fixpoint theorem with parameters

$$\varphi_e(y) = \varphi_{h(e(y),y)}$$

- “实时”全局最优？人机交互？
- 智能爆炸/技术奇点 ??? 持续升级 ≠ 指数迭代

③ 超越可计算：Kleene's relativized fixpoint theorem

$$\varphi_{e(y)}^A = \varphi_{h(e(y),y)}^A$$

- Gödel Machine PK AIXI^{tℓ}
- Gödel Machine PK AIXI

Limitation

- ① Gödel's first incompleteness theorem / Rice's theorem
- ② Gödel's second incompleteness theorem

$$\mathbb{T} \vdash \Box_{\mathbb{T}'} \varphi \rightarrow \varphi \implies \mathbb{T} \vdash \text{Con}(\mathbb{T}')$$

- 生物进化：达尔文 PK 拉马克
 - 生命 3.0
- ③ Legg's incompleteness theorem. *General prediction algorithms must be complex. Beyond a certain complexity they can't be mathematically discovered.*
 - ④ Complexity: higher-level abstractions — coarse grained.
 - 心理学：“过程忽视”、“峰终定律”
 - Information Bottleneck: Learning is to forget!
 - ⑤ Physical constraint: If we assume that it is not possible to measure properties without changing them (observer effect: α is fixpoint-free), then there is a limit to self-inspection.

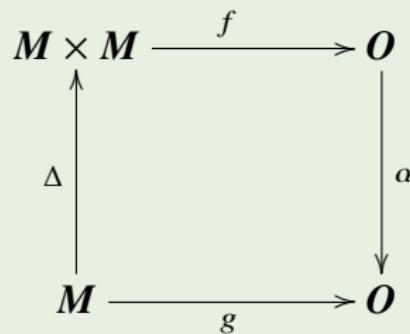
Consciousness — Integrated Information Theory

- Experts do their specialties best when they're in a state of "flow", aware only of what's happening at a higher level, and unconscious of the low-level details of how they're doing it. The information processing that we're consciously aware of is merely the tip of the iceberg. Many behaviors and brain regions are unconscious, with much of our conscious experience representing an after-the-fact summary of vastly larger amounts of unconscious information. Consciousness lags behind the outside world by about a quarter second. Brain measurements can sometimes predict your decision before you become conscious of having made it.
- Consciousness requires a kind of information processing that's fairly autonomous and integrated, so that the whole system is rather autonomous but its parts aren't. Given a physical process that, with the passage of time, transforms the initial state of a system into a new state, its integrated information measures inability to split the process into independent parts. In other word, it measures how much different parts of a system know about each other.

Non-operational Self-inspection?

The information available to the observer regarding his own state could have absolute limitations, by the laws of nature.

— von Neumann



- M : quantum measurements.
- O : possible outcomes of quantum measurements.

If we assume that it is not possible to measure properties without changing them (observer effect: α is fixpoint-free), then there is a limit to self-inspection.

Self-modification

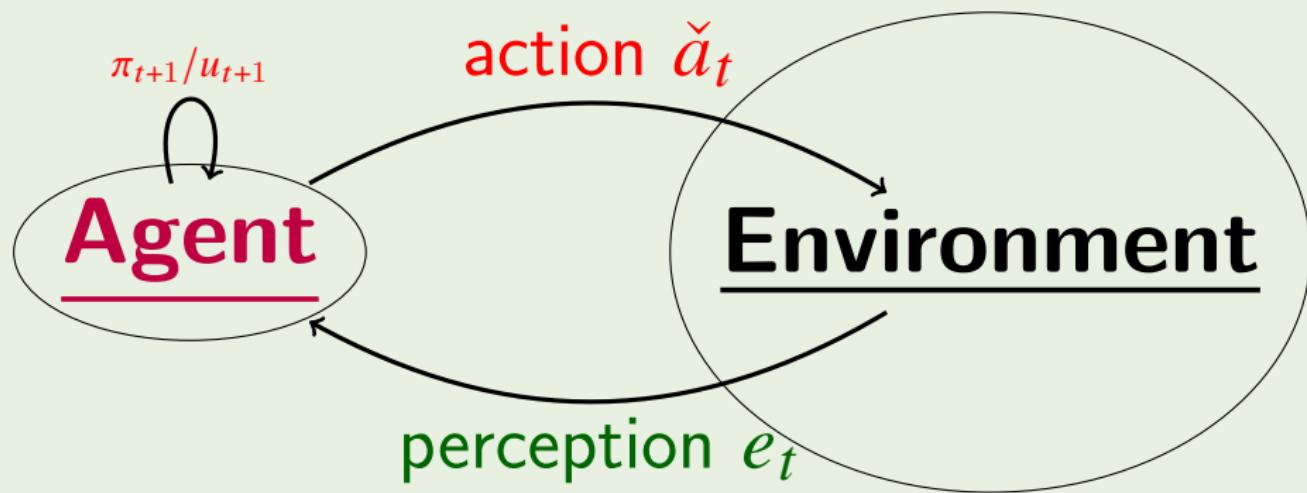
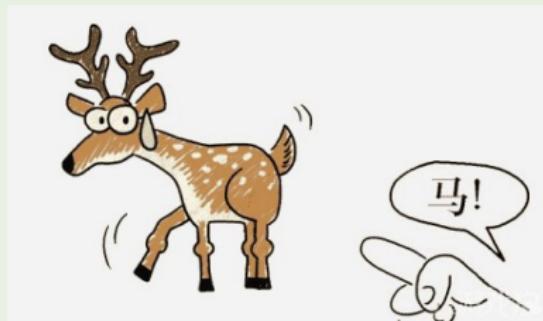
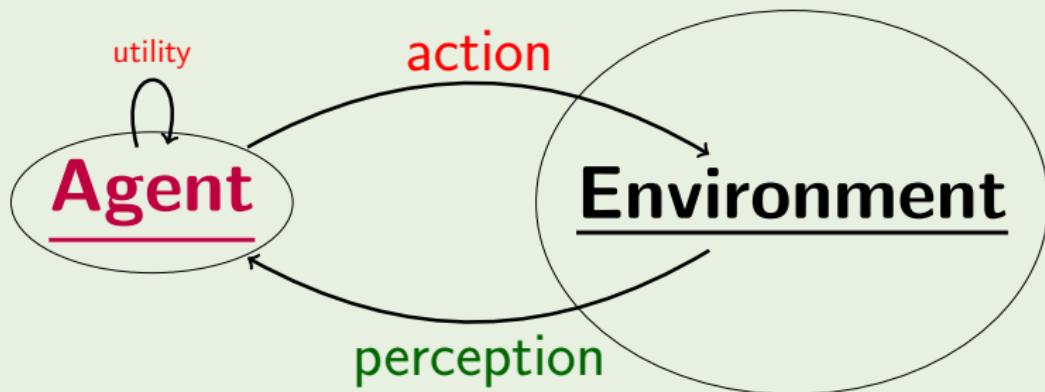


Figure: Policy/utility self-modification. $a_t = \langle \check{a}_t, \pi_{t+1} \rangle$ or $a_t = \langle \check{a}_t, u_{t+1} \rangle$

External/Internal Wireheading & Free Will²



- ① 我喜欢马。
指鹿为马！
- ② 我喜欢马。
我意欲自己
喜欢鹿！我
喜欢鹿！



² Everitt, Filan, Daswani, Hutter: Self-modification of policy and utility function in rational agents.

Frankfurt: Freedom of the will and the concept of a person.

Aaronson: The ghost in the quantum turing machine.

Calude, Kroon, Poznanovic: Free will is compatible with randomness.

Self-modification

Definition (Different Agents)

- Hedonistic Value

$$Q^{\text{h},\pi}(\boldsymbol{\alpha}_{<ta_t}) = \sum_{e_t \in \mathcal{E}} \rho(e_t | \check{\boldsymbol{\alpha}}_{<t} \check{a}_t) \left[\textcolor{red}{u_{t+1}}(\check{\boldsymbol{\alpha}}_{1:t}) + \gamma V^{\text{h},\pi}(\boldsymbol{\alpha}_{1:t}) \right]$$

- Ignorant Value

$$Q_t^{\text{i},\pi}(\boldsymbol{\alpha}_{<k} a_k) = \sum_{e_t \in \mathcal{E}} \rho(e_t | \check{\boldsymbol{\alpha}}_{<t} \check{a}_t) \left[\textcolor{red}{u_t}(\check{\boldsymbol{\alpha}}_{1:k}) + \gamma V_t^{\text{i},\pi}(\boldsymbol{\alpha}_{1:k}) \right]$$

- Realistic Value

$$Q_t^{\text{r}}(\boldsymbol{\alpha}_{<k} a_k) = \sum_{e_t \in \mathcal{E}} \rho(e_t | \check{\boldsymbol{\alpha}}_{<t} \check{a}_t) \left[\textcolor{red}{u_t}(\check{\boldsymbol{\alpha}}_{1:k}) + \gamma V_t^{\text{r},\pi_{t+1}}(\boldsymbol{\alpha}_{1:k}) \right]$$

Self-modification — Realistic Agent

$$V_t^\pi(\alpha_{<k}) = Q_t(\alpha_{<k} \pi(\alpha_{<k}))$$

$$Q_t(\alpha_{<k} a_k) = \sum_{e_k \in \mathcal{E}} \rho(e_k | \check{\alpha}_{<k} \check{a}_k) [\textcolor{red}{u_t}(\check{\alpha}_{1:k}) + \gamma V_t^{\pi_{t+1}}(\alpha_{1:k})]$$

$$\pi_1^* := \operatorname{argmax}_{\pi} V_1^\pi(\epsilon)$$

Theorem (All optimal policies are non-modifying)

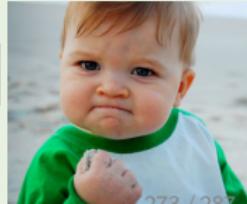
Let ρ and u_1 be modification-independent. For every $t \geq 1$, for all percept sequences $e_{<t}$, and for the action sequence $a_{<t}$ given by $a_i = \pi(\alpha_{<i})$, we have

$$Q_1(\alpha_{<t} \pi_t(\alpha_{<t})) = Q_1(\alpha_{<t} \pi_1^*(\alpha_{<t}))$$



All realistic optimal policies are non-modifying.

Not wireheading; But orthogonal!



Orthogonality and Wireheading in Self-improving GRL

$$\begin{aligned}
 V_t^\pi(\mathbf{a}_{<k}) &:= Q_t(\mathbf{a}_{<k} \pi(\mathbf{a}_{<k})) \\
 Q_t(\mathbf{a}_{<k} a_k) &:= \sum_{e_k \in \mathcal{E}} \sum_{\nu \in \mathcal{M}} \mathbf{w}_{\mathbf{a}_{<k}}^\nu v(e_k | \check{\mathbf{a}}_{<k} \check{a}_k) \left[\sum_{u \in \mathcal{U}} \sum_{a_k^H} P(a_k^H | a_k) P(u | \frac{\pi_{t+1}}{\nu}, \mathbf{a}_{<k} a_k, \mathbf{a}_{<k}^H a_k^H) u(\check{\mathbf{a}}_{1:k}) + \gamma V_t^{\pi_{t+1}}(\mathbf{a}_{1:k}) \right] \\
 \pi_t(\mathbf{a}_{<k}) &:= \operatorname{argmax}_{a_k \in \mathcal{A} \times \Pi} Q_t(\mathbf{a}_{<k} a_k)
 \end{aligned}$$

where

$$P(u | \nu, h) := \frac{\tilde{U}(u, \nu, h)}{\sum_{u \in \mathcal{U}_h} \tilde{U}(u, \nu, h)}$$

and

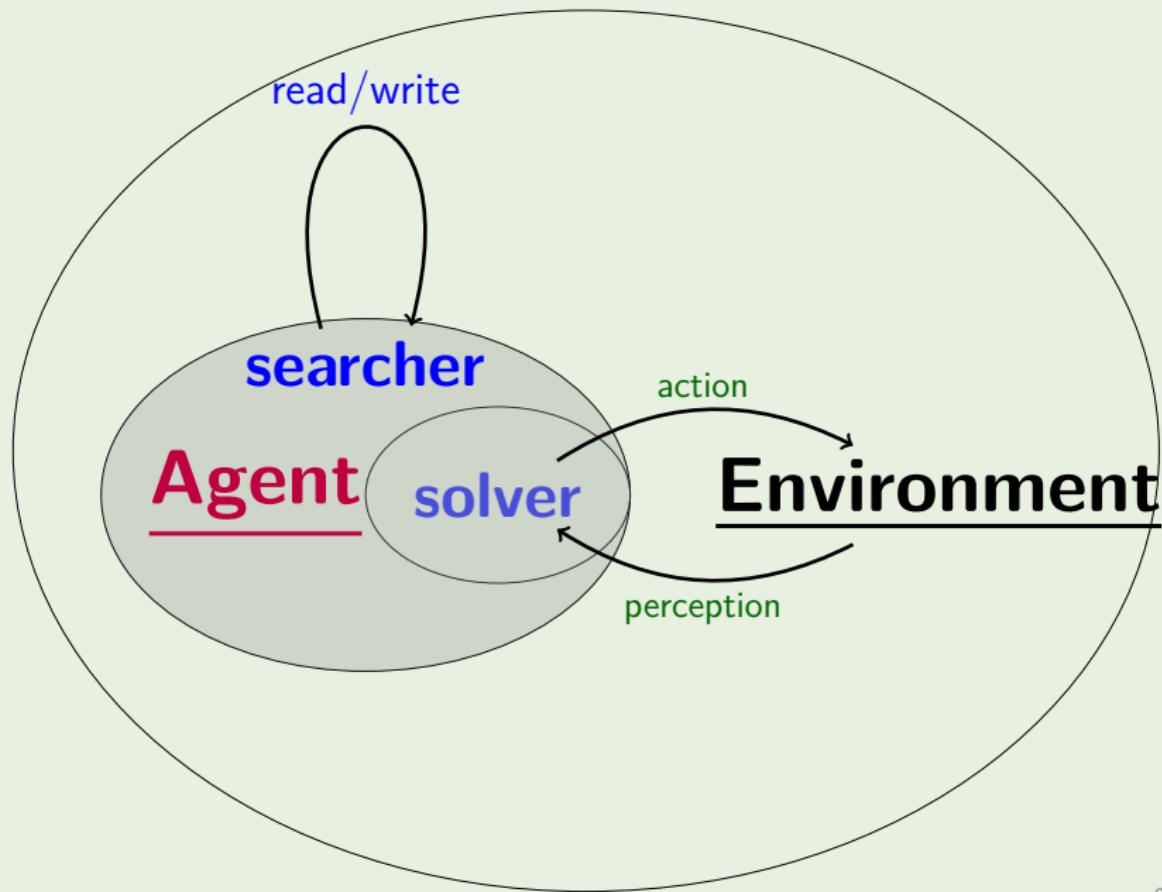
$$\tilde{U}(u, \nu, h) := \sum_{z \in \mathcal{Z}_h} \nu(z | h) u(z)$$

$$\pi^*(\mathbf{a}_{<t}) := \pi_t(\mathbf{a}_{<t})$$

$$\pi^*(e_{<t}) := \pi_t(e_{<t} | \pi_{t-1}(e_{<t-1}) \dots \pi_1(\epsilon) \dots) \quad (\text{Perfect Bayes-Nash})$$

uncertain model-based utility / IRL

Fatalism — God Bless All!



Orseau's Space-Time Embedded Intelligence

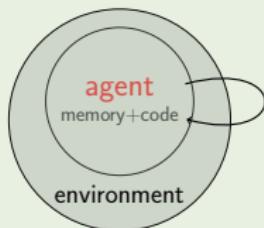
$$\pi^* := \underset{\pi_0 \in \Pi^\ell}{\operatorname{argmax}} V(\pi_0, \epsilon)$$

$$V(\pi_t, \mathbf{e}_{<t}) := \sum_{a_t = \langle \check{a}_t, \check{\pi}_{t+1} \rangle} \pi_t(a_t | \check{e}_{t-1}) \sum_{e_t = \langle \check{e}_t, \check{\pi}_{t+1} \rangle} \rho(e_t | \mathbf{e}_{<t} a_t) [u(\mathbf{e}_{1:t}) + \gamma_t V(\pi_{t+1}, \mathbf{e}_{1:t})]$$



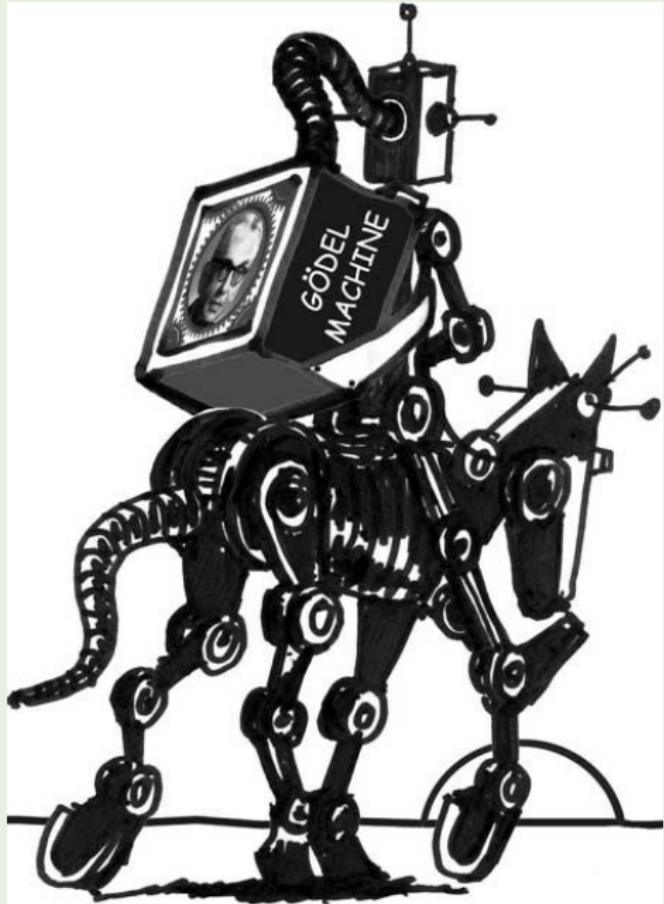
$$\pi^* := \underset{\pi_0 \in \Pi^\ell}{\operatorname{argmax}} V(\pi_0)$$

$$V(\pi_{<t}) := \sum_{\pi_t \in \Pi} \rho(\pi_t | \pi_{<t}) [u(\pi_{1:t}) + \gamma_t V(\pi_{1:t})]$$





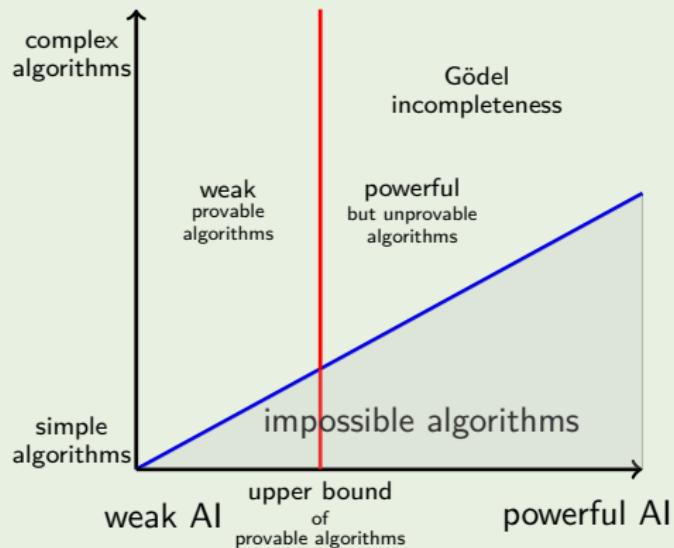
J. Schmidhuber, 1987



Incompressibility vs Incompleteness vs Intelligence

- $P(x) := \{p \in X^*: \exists m \forall n \geq m (p(x_{1:n}) = x_{n+1})\}$
- $P(A) := \bigcap_{x \in A} P(x)$
- $P_n := P(\{x: Km(x) \leq n\})$

- $\forall n \exists p \in P_n: K(p) \stackrel{+}{\leq} n + O(\log n)$
- $\forall n: p \in P_n \implies K(p) \stackrel{+}{\geq} n$



Theorem (Legg)

For any arithmetically sound Gödelian \mathbb{T} , $\exists n \forall p: \mathbb{T} \not\vdash p \in P_n$.

Universal Artificial Intelligence vs “Selective Amnesia”

- incomplete $\xrightarrow[\text{universal prior}]{\text{Harsanyi transformation}}$ imperfect \implies AIXI
- AIXI $\xrightarrow{\text{“Selective Amnesia”}}$ MDP

$$h \mapsto S$$

- partition of the set of histories = information set = state = feature
-

$$\xi^\pi(h'|ha) \rightarrow \mu^\pi(h'|ha)$$

but

$$\xi^\pi(S'|Sa) \not\rightarrow \mu^\pi(S'|Sa)$$

“Selective Amnesia”

① compressible

$$K(S) \leq \sum_{h \in S} K(h)$$

② minimal

$$\forall part(S) : K(S) \leq \sum_{S_i \in part(S)} K(S_i)$$

③ maximal

$$\forall S' \supset S : K(S) \leq K(S')$$

④ MDL

$$\forall S' \in \mathcal{S} \forall h \in S' : K(S) + K(h|S) \leq K(S') + K(h|S')$$

⑤ MDL/utility

$$\forall S' : K\left(S_{1:n}^S | a_{1:n}\right) + K\left(u_{1:n} | S_{1:n}^S, a_{1:n}\right) + K(\mathcal{S}) \leq K\left(S_{1:n}^{S'} | a_{1:n}\right) + K\left(u_{1:n} | S_{1:n}^{S'}, a_{1:n}\right) + K(\mathcal{S}')$$

$$u(h) := \left[\left[K(h) < \ell(h) \text{ } \& \text{ } \forall h' > h \left(K(h) \leq K(h') \text{ } \& \text{ } \forall part(h) \left(\sum_{h' \in part(h)} K(h') \geq K(h) \right) \right) \right] \right]$$

Potapov's MSearch + RSearch

- Let $\{x_i\}_{i=1}^n$ be a set of strings.
- $K(x_1 \dots x_n) \approx \min_S \left(\ell(S) + \sum_{i=1}^n K(x_i | S) \right) \ll \sum_{i=1}^n K(x_i)$
- search for models $y_i^* := \underset{y: S(y)=x_i}{\operatorname{argmin}} \ell(y)$ for each x_i w.r.t. some best representation $S^* := \underset{S}{\operatorname{argmin}} \left[\ell(S) + \sum_{i=1}^n \ell(y_i^*) \right]$

① Search for models

$$M\text{Search}(S, x_i) \rightarrow y_i^* = \underset{y: S(y)=x_i}{\operatorname{argmin}} \ell(y)$$

② Search for representations

$$R\text{Search}(x_1 \dots x_n) \rightarrow S^* = \underset{S}{\operatorname{argmin}} \left[\ell(S) + \sum_{i=1}^n \ell(y_i^*) \right]$$

- $M\text{Search}$ enumerates all models to find the shortest model:
 $S(y_i) = x_i$.
- $R\text{Search}$ enumerates all S and calls $M\text{Search}$ for each S .

Specialization and $SS' - Search$

Theorem (smn Theorem)

For any $m, n > 0$, there exists a primitive recursive function s_n^m of $m + 1$ arguments s.t. for every Gödel number e of a partial recursive function with $m + n$ arguments

$$\varphi_{s_n^m(e, x_1, \dots, x_m)} = \lambda y_1 \dots y_n. \varphi_e(x_1, \dots, x_m, y_1, \dots, y_n)$$

$$\forall x: spec(MSearch, S)(x) = MSearch(S, x)$$

$$S' := spec(MSearch, S) \implies \begin{cases} \forall x: S(S'(x)) = x \\ \ell(S) + \sum_{i=1}^n \ell(S'(x_i)) \rightarrow \min \end{cases}$$

- S is a generative representation. (decoding)
- S' is a descriptive representation. (encoding)
- $SS' - Search$ simultaneous search for S and S' .

Potapov's Representational MDL

$$K(x_{1:n}) \approx \min_S \left(\ell(S) + \sum_{i=1}^n K(x_i|S) \right) \ll \sum_{i=1}^n K(x_i)$$

$$q_1^* := \operatorname{argmin}_q [\ell(q) + K(x|S_1 q)]$$

$$q_{i+1}^* := \operatorname{argmin}_q [\ell(q) + K(q_i^*|S_{i+1} q)]$$

$$L_{S_1 \dots S_m}(x) := K(x|S_1 q_1^*) + \sum_{i=2}^{m-1} K(q_i^*|S_{i+1} q_{i+1}^*) + \ell(q_m^*)$$

$$a_k^* := \operatorname{argmax}_{a_k} \max_{p: U(p, e_{<k}) = a_{<k} a_k} \sum_{q: U(q, a_{<k}) = e_{<k}} 2^{-\ell(q)} V_q^p(\boldsymbol{\alpha}_{<k})$$

$$a_k^* := \operatorname{argmax}_{a_k} \max_{p: U(p, e_{<k}) = a_{<k} a_k} \sum_{\{q_i\}: U(S\{q_i\}, a_{<k}) = e_{<k}} 2^{-\ell(\{q_i\})} V_{\{q_i\}}^p(\boldsymbol{\alpha}_{<k})$$

where $e_{<k} = e_{m_1+1:m_2} \dots e_{m_{n-1}+1:m_n}$, $m_1 = 0$, $m_n = k - 1$, and
 $U(Sq_i a_{<k}) = e_{m_i+1:m_{i+1}}$.

$$Q(q_k = s, a_k = a) := \max_{p: U(p, e_{<k}) = a_{<k} a} \sum_{\{q_i\}: q_k = s, U(S\{q_i\}, a_{<k}) = e_{<k}} 2^{-\ell(\{q_i\})} V_{\{q_i\}}^p(\boldsymbol{\alpha}_{<k})$$

$$Q(q_k = s) := \max_{a_k} Q(q_k = s, a_k = a)$$

Fundamental Challenges

- What is a good optimality criterion?
- What is a “natural” UTM/prior?
- Prior vs universality
- Exploration vs exploitation
- Where should the reward come from?
- How should the future be discounted?
- How should agents reason about themselves (or other agents reasoning about itself)?
- What is a practically feasible and general way of doing induction and planning?
- AIXI in the multi-agent setting.
- Better variants/approximations.
- Training: 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.

- A Blind Man in a Dark Room Looking for a Black Cat That Is Not There?



- The Singularity is Near?

