
OPTIMALITY OF UNIVERSAL BAYESIAN PREDICTION FOR GENERAL LOSS AND

Marcus Hutter

Istituto Dalle Molle di Studi sull'Intelligenza
IDSIA, Galleria 2, CH-6928 Manno-Lugano, S
marcus@idsia.ch, <http://www.idsia.ch/~m>

2000 – 2002

Universal Induction = Ockham + Epicurus + B

$$\frac{\text{Loss}(\text{Universal Prediction Scheme})}{\text{Loss}(\text{Any other Prediction Scheme})} \leq 1 + o(1)$$

Table of Contents

- The Philosophical Dilemma of Predicting t
- (Conditional) Probabilities and their Interp
- Probability that the Sun will rise Tomorrow
- Kolmogorov Complexity
- Universal Probability Distribution
- Universal Sequence Prediction
- Loss Bounds & Optimality
- Application to Games of Chance (Bound o
- Generalization: Continuous Probability Cla
- Generalization: The Universal $AI\xi$ Model
- Further Generalizations, Outlook, Conclusi

Problem Setup

- Every induction problem can be phrased as a sequence prediction problem.
- Classification is a special case of sequence prediction.
(With some tricks the other direction is also true)
- I'm interested in maximizing profit (minimizing loss).
I'm not (primarily) interested in finding a (true/predictive/consistent) hypothesis.
- Separating noise from data is *not* necessary in this setting!

My Position to Occam

- Most of us believe in or at least use the axioms of logic, probability theory, and the natural numbers when doing science, without questioning their validity.
- We should/must add Occam's razor in some quantified form as an axiom because it is the foundation of machine learning and science.
- There is (yet) no mathematical proof of Occam's razor, and it is not an independent axiom, but there is lots of evidence that this is a good principle.

On the Foundations of Machine Learning

- Example: **Algorithm/complexity theory**: The goal is to find hard problems and to show lower bounds on their computation time.
rigorously defined: algorithm, Turing machine, problem class, ...
- Most **disciplines** start with an informal way of attacking a subject and get **more and more formalized** often to a point where they are fully formalized.
Examples: set theory, logical reasoning, proof theory, probability theory, infinitesimal calculus, quantum field theory, ...
- **Machine learning**: Tries to build and understand systems which learn from data, to make good prediction, which are able to generalize.
Many terms only **vaguely defined** or there are many alternatives.

Occam to the Rescue

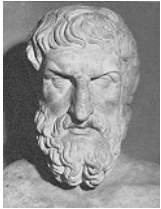
- Is it possible to give machine learning a rigorous mathematical framework/definition?
- Yes! Use Occam's razor, quantified in terms of Kolmogorov complexity, combine it with Bayes, and possibly sequential decision theory.
- There is at the moment no alternative suggestion of how to formalize machine learning rigorously.

My view of (future) Machine learning

- Application = Solve learning tasks by approximating Kolmogorov complexity (MML, MDL, SRM, and much more specific ones, like SVM)
- Theory = Proof theorems, especially on convergence and approximation
- Non-standard ML = Modify "Occam's axiom" with the goal of making it work better.

Induction = Predicting the Future

Extrapolate past observations to the future
 but how can we know something about the future?



Epicurus' principle of multiple explanations

If more than one theory is consistent with the observations...



Ockham's razor (simplicity) principle

Entities should not be multiplied beyond necessity.



Hume's negation of Induction The only form of induction as the conclusion is already logically contained in the premises.



Bayes' rule for conditional probabilities

Given the prior belief/probability one can predict a future event.



Solomonoff's universal prior

Solves the question of how to choose the prior if no other information is available.

Strings and Conditional Probab

Strings: $x = x_1x_2\dots x_n$ with $x_t \in \mathcal{X}$ and $x_{1:m} := x_1x_2\dots x_{m-1}x_m$

$\rho(x_1\dots x_n)$ is the probability that an (infinite) sequence starts with

Heavy use of Bayes' rule in the following forms:

$$\rho(x_n|x_{<n}) = \rho(x_{1:n})/\rho(x_{<n}),$$

$$\rho(x_1\dots x_n) = \rho(x_1) \cdot \rho(x_2|x_1) \cdot \dots \cdot \rho(x_n|x_1\dots x_{n-1})$$

If the true prior probability $\mu(x_1\dots x_n)$ is known, then the optimal strategy to minimize the μ - expected loss.

Interpretation of Probabilities

Frequentist: Probabilities come from experiments.

Objectivist: Probabilities are real aspects of the world.

Subjectivist: Probabilities describe ones believe.

Probability of Sunrise Tomorrow

What is the probability that the sun will rise tomorrow? It is $\mu(\text{sun will rise tomorrow})$ where μ is the probability distribution over the lifetime of the sun in days.

1 = sun raised. 0 = sun will not raise.

- The probability is undefined, because there has never been an experiment that has tested the existence of the sun *tomorrow* (reference class problem).
- The probability is 1, because in all experiments that have been conducted (in the past d days) the sun raised.
- The probability is $1 - \epsilon$, where ϵ is the proportion of stars in the universe that explode in a supernova per day.
- The probability is $(d + 1)/(d + 2)$ (Laplace estimate by assuming a Bernoulli process with uniformly distributed raising prior probability p).
- The probability can be derived from the type, age, size and distance of the sun, even though we never have observed another star with similar characteristics.

Solomonoff solved the problem of unknown prior μ by introducing a universal probability distribution ξ related to Algorithmic Information Theory.

Kolmogorov Complexity

The Kolmogorov Complexity of a string x is the length of the shortest program producing x .

$$K(x) := \min_p \{l(p) : U(p) = x\} \quad , \quad U = \text{universal}$$

The definition is "nearly" independent of the choice of U

$$|K_U(x) - K_{U'}(x)| < c_{UU'}, \quad K_U(x) \stackrel{\pm}{\approx} K_{U'}$$

$\stackrel{\pm}{\approx}$ indicates equality up to a constant $c_{UU'}$ independent of x .

K satisfies most properties an information measure should satisfy

$$K(xy) \stackrel{+}{\leq} K(x) + K(y).$$

$K(x)$ is not computable, but only co-enumerable (semi-computable)

Universal Probability Distribution

The universal semimeasure is the probability that output of U starting from input x is provided with fair coin flips

$$\xi(x) = \sum_{\mu_i \in \mathcal{M}} w_{\mu_i} \cdot \mu_i(x) \stackrel{\times}{=} \sum_{p : U(p)=x^*} 2^{-l(p)}, \quad \text{e.g. } w_{\mu_i}$$

[Solomonoff 64]

Universality property of ξ : ξ dominates every computable probability measure

$$\xi(x) \geq w_{\mu_i} \cdot \mu_i(x) \quad \forall \mu_i \in \mathcal{M}$$

Furthermore, the μ expected squared distance sum between ξ and any computable μ

$$\sum_{t=1}^{\infty} \sum_{x_{1:t}} \mu(x_{<t}) (\xi(x_t|x_{<t}) - \mu(x_t|x_{<t}))^2 \leq \ln 2$$

[Solomonoff 78] (for binary alphabet)

$$\Rightarrow \xi(x_n|x_{<n}) \xrightarrow{n \rightarrow \infty} \mu(x_n|x_{<n}) \text{ with } \mu \text{ probability } 1 \quad \Rightarrow \quad \xi \text{ is } \mu\text{-almost surely correct}$$

Convergence Theorem

The universal conditional probability $\xi(x_t|x_{<t})$ of the next symbol related to the true conditional probability $\mu(x_t|x_{<t})$ in the following

$$i) \quad \sum_{t=1}^n \mathbf{E} \left[\sum_{x_t} \left(\mu(x_t|x_{<t}) - \xi(x_t|x_{<t}) \right)^2 \right] \equiv S_n \leq D_n \leq$$

$$ii) \quad \sum_{x_t} \left(\mu(x_t|x_{<t}) - \xi(x_t|x_{<t}) \right)^2 \equiv s_t(x_{<t}) \leq d_t(x_{<t}) \rightarrow$$

$$iii) \quad \xi(x'_t|x_{<t}) \rightarrow \mu(x'_t|x_{<t}) \quad \text{for } t \rightarrow \infty \text{ with } \mu \text{ probability } 1$$

$$iv) \quad \sum_{t=1}^n \mathbf{E} \left[\left(\sqrt{\frac{\xi(x_t|x_{<t})}{\mu(x_t|x_{<t})}} - 1 \right)^2 \right] \leq D_n \leq \ln w_\mu^{-1} < \infty$$

$$v) \quad \frac{\xi(x_t|x_{<t})}{\mu(x_t|x_{<t})} \rightarrow 1 \quad \text{for } t \rightarrow \infty \text{ with } \mu \text{ probability } 1$$

where $d_t = \sum_{x_n} \mu(x_n|x_{<n}) \ln \frac{\mu(x_n|x_{<n})}{\xi(x_n|x_{<n})}$ and $D_n = \sum_{x_{1:n}} \mu(x_{1:n})$ relative entropies, and w_μ is the weight of μ in ξ .

Universal Sequence Prediction

A prediction is very often the basis for some decision. The decision which itself leads to some reward or loss. Let $\ell_{x_t y_t} \in [0, 1]$ be the loss of taking action $y_t \in \mathcal{Y}$ and $x_t \in \mathcal{X}$ is the t^{th} symbol of the sequence. Example: decision $\mathcal{Y} = \{\text{umbrella, sunglasses}\}$ based on weather forecasts.

Loss	sunny	rainy
umbrella	0.3	0.1
sunglasses	0.0	1.0

The goal is to minimize the μ -expected loss. More generally we consider a prediction scheme

$$y_t^{\Lambda_\rho} := \arg \min_{y_t \in \mathcal{Y}} \sum_{x_t} \rho(x_t | x_{<t}) \ell_{x_t y_t}$$

which minimizes the ρ -expected loss. The actual μ -expected loss of the t^{th} symbol and the total μ -expected loss in the first n predictions

$$l_{t\Lambda_\rho}(x_{<t}) := \sum_{x_t} \mu(x_t | x_{<t}) \ell_{x_t y_t^{\Lambda_\rho}} \quad , \quad L_n^{\Lambda_\rho} := \sum_{t=1}^n \sum_{x_{<t}} \mu(x_{<t}) l_{t\Lambda_\rho}(x_{<t})$$

Loss Bounds (Main Theorem)

$L_n^{\Lambda_\mu}$ made by the informed scheme Λ_μ ,

$L_n^{\Lambda_\xi}$ made by the universal scheme Λ_ξ ,

L_n^Λ made by **any** (causal) prediction scheme Λ .

i) $L_n^{\Lambda_\mu} \leq L_n^\Lambda$ for **any** (causal) prediction scheme Λ .

ii) $0 \leq L_n^{\Lambda_\xi} - L_n^{\Lambda_\mu} \leq 2D_n + 2\sqrt{L_n^{\Lambda_\mu} D_n}$

iii) if $L_{\infty\Lambda_\mu}$ is finite, then $L_{\infty\Lambda_\xi}$ is finite

iv) $L_n^{\Lambda_\xi} / L_n^{\Lambda_\mu} = 1 + O((L_n^{\Lambda_\mu})^{-1/2}) \xrightarrow{L_n^{\Lambda_\mu} \rightarrow \infty} 1$

v) $\sum_{t=1}^n \mathbf{E}[(l_{t\Lambda_\xi}(x_{<t}) - l_{t\Lambda_\mu}(x_{<t}))^2] \leq 2D_n \leq 2 \ln w_\mu^{-1}$

vi) $0 \leq l_{t\Lambda_\xi}(x_{<t}) - l_{t\Lambda_\mu}(x_{<t}) \leq \begin{cases} \sqrt{2d_t(x_{<t})} \\ 2d_t(x_{<t}) + 2\sqrt{l_{t\Lambda_\mu}(x_{<t})} \end{cases}$

where $D_n := \sum_{x_{1:n}} \mu(x_{1:n}) \ln \frac{\mu(x_{1:n})}{\xi(x_{1:n})} \leq \ln \frac{1}{w_\mu} = \ln 2 \cdot K(\mu)$

Remark: The bound is valid for **any loss function** $\in [0, 1]$ with n i.i.d., Markovian, stationary, ergodic, ...) on the structure of the

Example Application

A dealer has two dice, one with 2 white and 4 black faces, the other with 4 white and 2 black faces. He chooses a die according to some deterministic strategy. We bet $s = \$3$ on white or black and receive $r = \$5$ for every correct prediction.

If we know μ , i.e. the die the dealer chooses, we should predict the correct side and win money. Expected Profit (= -Loss): $P_{n\Lambda_\mu}/n = \frac{1}{3}$.

If we don't know μ we can use Solomonoff prediction scheme Λ_ξ and achieve the same profit:

$$P_{n\Lambda_\xi}/P_{n\Lambda_\mu} = 1 - O(n^{-1/2})$$

Bound on Winning Time

Estimate of the number of rounds before reaching the winning strategy.

$$P_{n\Lambda_\xi} > 0 \quad \text{if} \quad L_n^{\Lambda_\xi} < 0 \quad \text{if} \quad n > 330 \ln 2 \cdot K(\mu) + O(1)$$

Λ_ξ is asymptotically optimal with rapid convergence.

General Bound for Winning T

For every (passive) game of chance for which there exists a winning strategy, you can make money by using Λ_ξ even if you don't know the underlying process/algorithm.

Λ_ξ finds and exploits every regularity.

The time n needed to reach the winning zone is

$$n \leq \left(\frac{2p_\Delta}{\bar{p}_{n\Lambda_\mu}} \right)^2 \cdot \ln \frac{1}{w_\mu}, \quad \bar{p}_{n\Lambda_\mu} := \frac{1}{n} \sum_{t=1}^n p_{t\Lambda_\mu}, \quad p_\Delta =$$

Generalization: Continuous Probability

In statistical parameter estimation one often has a continuous h Bernoulli(θ) process with unknown $\theta \in [0, 1]$).

$$\mathcal{M} := \{\mu_\theta : \theta \in \mathbb{R}^d\}, \quad \xi(x_{1:n}) := \int_{\mathbb{R}^d} d\theta w(\theta) \cdot \mu_\theta(x_{1:n}),$$

The only property of ξ needed was $\xi(x_{1:n}) \geq w_{\mu_i} \cdot \mu_i(x_{1:n})$ which dropping the sum over μ_i . Here, restrict the integral over \mathbb{R}^d to θ . For sufficiently smooth μ_θ and $w(\theta)$ we expect

$$\xi(x_{1:n}) \gtrsim |N_{\delta_n}| \cdot w(\theta) \cdot \mu_\theta(x_{1:n}) \implies D_n \lesssim \ln \frac{1}{w_\mu}$$

The average Fisher information \bar{j}_n measures the curvature (par $\ln \mu_\theta$. Under some weak regularity conditions on \bar{j}_n one can show

$$D_n := \sum_{x_{1:n}} \mu(x_{1:n}) \ln \frac{\mu(x_{1:n})}{\xi(x_{1:n})} \leq \ln \frac{1}{w_\mu} + \frac{d}{2} \ln \frac{n}{2\pi} + \frac{1}{2}$$

i.e. D_n grows only logarithmically with n .

Optimality of the Universal Prior

- There are \mathcal{M} and $\mu \in \mathcal{M}$ and weights w_μ for which the loss
- The universal prior ξ is **pareto-optimal**, in the sense that the $\mathcal{F}(\mu_a, \rho) \leq \mathcal{F}(\mu_a, \xi)$ for all $\mu_a \in \mathcal{M}$ and strict inequality for $\rho \neq \xi$ where \mathcal{F} is the instantaneous or total squared distance s_t, S_t, d_t, D_n , or error e_t, E_n , or loss l_t, L_n .
- ξ is **elastic pareto-optimal** in the sense that by accepting a slight decrease in some environments one can only achieve a slight increase in other environments.
- Within the set of enumerable weight functions with short program weights $w_\nu = 2^{-K(\nu)}$ lead to the smallest performance bound (to $\ln w_\mu^{-1}$) constant in all enumerable environments.

Does all this justify Occam's razor ?

Larger & Smaller Environmental C

- all finitely computable probability measures
($\xi \notin \mathcal{M}$ in no sense computable)
- all enumerable (approximable from below) semi-measures [S
($\xi \in \mathcal{M}$ enumerable)
- all cumulatively enumerable semi-measures [Schmidhuber 01
(distribution enumerable and $\in \mathcal{M}$)
- all approximable (asymptotically computable) measures [Sch
($\xi \notin \mathcal{M}$ in no sense computable)
- Speed prior related to Levin complexity and Levin search [S
(which distributions are dominated?)
- finite-state automata instead of general Turing machines [F
to Lempel-Ziv data compression ($\xi \notin \mathcal{M}$)

Generalization: The Universal AI

$$\text{Universal AI} = \text{Universal Induction} + \text{Sequential Dec}$$

Replace μ^{AI} in decision theory model $AI\mu$ by an appropriate gen

$$\xi(yx_{1:t}) := \sum_{q:q(y_{1:t})=x_{1:t}} 2^{-l(q)}$$

$$y_t = \arg \max_{y_t} \sum_{x_t} \max_{y_{t+1}} \sum_{x_{t+1}} \dots \max_{y_m} \sum_{x_m} (r(x_t) + \dots + r(x_m)) \cdot \xi(x_{t:m}|y_{t:m})$$

Claim: $AI\xi$ is the most intelligent environmental independent, i.e. agent possible.

Applications: Strategic Games, Function Minimization, Supervis
Examples, Sequence Prediction, Classification.

[Proceedings of ECML-2001] and [<http://www.hutt>]

Further Generalizations:

- Time and history dependent loss function in general intervals.
- Infinite (countable and uncountable) action/decision space.
- Partial Sequence Prediction.
- Independent Experiments & Classification.

Outlook:

- Infinite (prediction) alphabet \mathcal{X} .
- Delayed and Probabilistic Sequence Prediction.
- Unification with (Lossbounds for) aggregating strategies.
- Determine suitable performance measures for universal $AI\xi$.
- Study learning aspect of Λ_ξ and $AI\xi$.
- Information theoretic interpretation of winning time.
- Implementation and application of Λ_ξ for specific finite \mathcal{M} .
- Downscale theory and results to MDL approach.

Conclusions

- Solomonoff's prediction scheme, which is related to Kolomo formally solves the general problem of induction.
- We proved convergence and loss-bounds for Solomonoff prediction is well suited, even for difficult prediction problems.
- We proved several optimality properties for Solomonoff prediction.
- We made no structural assumptions on the probability distribution.
- The bounds are valid for any bounded loss function.
- We proved a bound on the time to win in games of chances.
- Discrete and continuous probability classes have been considered.
- Generalizations to active agents with reinforcement feedback.
- At least all this is a lot of evidence that Occam's razor is a

See [<http://www.idsia.ch/~marcus>] for details.