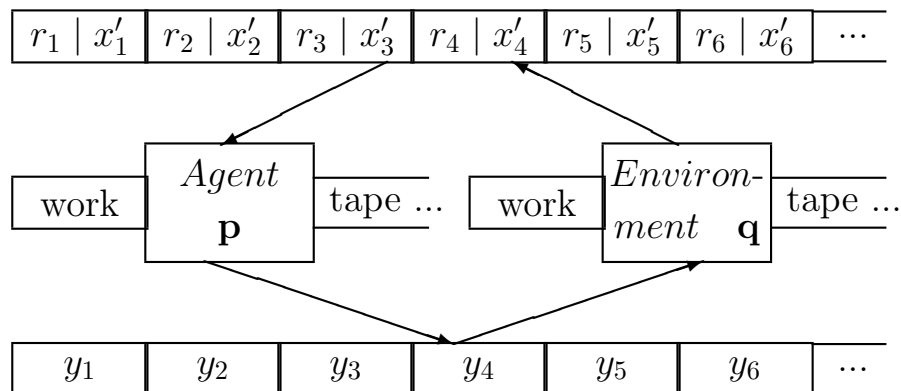


OPTIMAL SEQUENTIAL DECISIONS

BASED ON

ALGORITHMIC PROBABILITY



Marcus Hutter

OPTIMAL SEQUENTIAL DECISIONS BASED ON ALGORITHMIC PROBABILITY

Habilitationsschrift vorgelegt von

Marcus Hutter

IDSIA, Galleria 2, CH-6928 Manno-Lugano

zur Erlangung der *venia legendi* in Informatik
an der Technischen Universität München

Juni 2003

$$\begin{array}{rcl} \text{Decision Theory} & = & \text{Probability} + \text{Utility Theory} \\ + & & + \\ \text{Universal Induction} & = & \text{Occam} + \text{Epicurus} + \text{Bayes} \\ \parallel & & \parallel \\ \text{Universal Artificial Intelligence without Parameters} \end{array}$$

Preface

Personal motivation. The dream of creating artificial devices which reach or outperform human intelligence is an old one. It is also one of the two dreams of my youth, which have never let me go (the other is finding a physical theory of everything). What makes this challenge so interesting? A solution would have enormous implications on our society, and there are reasons to believe that the AI problem can be solved in my expected lifetime. So it's worth sticking to it for a lifetime, even if it will take 30 years or so to reap the benefits.

The AI problem. The science of Artificial Intelligence (AI) may be defined as the construction of intelligent systems and their analysis. A natural definition of a *system* is anything which has an input and an output stream. Intelligence is more complicated. It can have many faces like creativity, solving problems, pattern recognition, classification, learning, induction, deduction, building analogies, optimization, surviving in an environment, language processing, knowledge, and many more. A formal definition incorporating every aspect of intelligence, however, seems difficult. Most, if not all known facets of intelligence can be formulated as goal driven or, more precisely, as maximizing some utility function. It is, therefore, sufficient to study goal driven AI; e.g. the (biological) goal of animals and humans is to survive and spread. The goal of AI systems should be to be useful to humans. The problem is that, except for special cases, we know neither the utility function, nor the environment in which the agent will operate, in advance. The major goal of this thesis is to develop a theory which solves these problems.

The nature of this thesis. The thesis is theoretical in nature. For most parts we assume availability of unlimited computational resources. The first important observation is that this makes the AI problem not trivial. Playing chess optimally or solving NP-hard problems become trivial, but driving a car or surviving in nature do not. The reason being, that it is a challenge itself to well-define these problems, not to mention to present an algorithm. In other words: The AI problem has not yet been well-defined. One may view the thesis as a suggestion for (and discussion of) such a mathematical definition of AI.

Extended abstract. The *Goal* of this thesis is to develop a universal theory of sequential decision making akin to Solomonoff's celebrated universal theory of induction. Solomonoff derived an optimal way of predicting future data, given previous

observations, provided the data is sampled from a computable probability distribution. We extend this approach to derive an optimal rational reinforcement learning agent, called AIXI, embedded in an unknown environment. The *main idea* is to replace the unknown environmental distribution μ in the Bellman equations by a suitably generalized universal Solomonoff distribution ξ . The state space is the space of complete histories. AIXI is a universal theory without adjustable parameters, making no assumptions about the environment except that it is sampled from a computable distribution. From an algorithmic complexity perspective, the AIXI model generalizes optimal passive universal induction to the case of active agents. From a decision theoretic perspective, AIXI is a suggestion of a new (implicit) “learning” algorithm, which may overcome all (except computational) problems of previous reinforcement learning algorithms.

Chapter 1: We start with a survey of the contents and main results of this work.

Chapter 2: How and in which sense induction is possible at all has been subject to long philosophical controversies. Highlights are Epicurus’ principle of multiple explanations, Occam’s razor and the chain rule for conditional probabilities. Solomonoff elegantly unified all these aspects into one formal theory of inductive inference based on a universal probability distribution ξ , which is closely related to Kolmogorov complexity $K(x)$, the length of the shortest program computing x . We classify the (non)existence of universal priors for several generalized computability concepts.

Chapter 3: We prove rapid convergence of ξ to the unknown true environmental distribution μ and tight loss bounds for arbitrary bounded loss functions and finite alphabet. We show Pareto-optimality of ξ in the sense that there is no other predictor which performs better or equal in all environments and strictly better in at least one. Finally, we give an Occam’s razor argument showing that predictors based on ξ are optimal. We apply the results to games of chance and compare them to predictions with expert advice. All together this shows that Solomonoff’s induction scheme represents a universal (formal) solution to all *passive* prediction problems.

Chapter 4: Sequential decision theory provides a framework for finding optimal reward-maximizing strategies in a *reactive* environment (e.g. chess playing as opposed to weather forecasting), assuming the environmental probability distribution μ is known. We present this theory in a very general form (called $AI\mu$ model) in which actions and observations may depend on arbitrary past events. We clarify the connection to the Bellman equations and discuss minor parameters including (the size of) the I/O spaces and the lifetime of the agent and their universal choice which we have in mind. Optimality of $AI\mu$ is obvious by construction.

Chapter 5: Reinforcement learning algorithms are usually used in the case of unknown μ . They can succeed if the state space is either small or has effectively been made small by generalization techniques. The algorithms work only in restricted (e.g. Markovian) domains, have problems with optimally trading off exploration versus exploitation, have nonoptimal learning rate, are prone to diverge, or are otherwise ad hoc. The formal solution proposed in this thesis is to generalize Solomonoff’s

universal prior ξ to include conditions and replace μ by ξ in the $\text{AI}\mu$ model, resulting in the AIXI model, which we claim to be universally optimal. We investigate what we can expect from a universally optimal agent and clarify the meanings of *universal*, *optimal*, etc. We show that (a variant of) AIXI is self-optimizing and Pareto-optimal.

Chapter 6: We show how a number of AI problem classes fit into the general AIXI model. They include sequence prediction, strategic games, function minimization, and supervised learning. We first formulate each problem class in its natural way (for known μ) and then construct a formulation within the $\text{AI}\mu$ model and show their equivalence. We then consider the consequences of replacing μ by ξ . The main goal is to understand in which sense the problems are solved by AIXI.

Chapter 7: The major drawback of AIXI is that it is incomputable, or more precisely, only asymptotically computable, which makes an implementation impossible. To overcome this problem, we construct a modified model $\text{AIXI}t_l$, which is still superior to any other time t and length l bounded algorithm. The computation time of $\text{AIXI}t_l$ is of the order $t \cdot 2^l$. A way of overcoming the large multiplicative constant 2^l is presented at the expense of an (unfortunately even larger) additive constant. The constructed algorithm M is capable of solving all well-defined problems p as quickly as the fastest algorithm computing a solution to p , save for a factor of $1 + \varepsilon$ and lower-order additive terms. The solution requires an implementation of first order logic, the definition of a universal Turing machine within it and a proof theory system.

Chapter 8: Finally we discuss and remark on some otherwise unmentioned topics of general interest. We also critically review what has been achieved in this thesis, including assumptions, problems, limitations, performance, and generality of AIXI in comparison to other approaches to AI. We conclude the thesis with some less technical remarks on various philosophical issues.

Prerequisites. I tried to make the thesis as self-contained as possible. Especially I provide all necessary background knowledge on algorithmic information theory in Chapter 2 and sequential decision theory in Chapter 4. Nevertheless, some prior knowledge in these areas could be of some help. The chapters have been designed to be readable independently of one another (after having read Chapter 1). This necessarily implies minor repetitions. Most of the issues to be addressed in the thesis can already be found scattered in various reports and publications of mine, available at <http://www.idsia.ch/~marcus/ai>. Feedback is welcome (e.g. errors, typos, proofs, ...).

Problem classification. There are problems included at the end of each chapter of different motivation and difficulty. We use Knuth’s rating scheme for exercises [Knu73] in slightly adapted form (applicable if the material in the corresponding chapter has been understood). In-between values are possible.

- C00 *Very easy*. Solvable from the top of your head.
- C10 *Easy*. Needs 15 minutes to think, possibly pencil and paper.
- C20 *Average*. May take 1–2 hours to answer completely.
- C30 *Moderately difficult or lengthy*. May take several hours to a day.
- C40 *Quite difficult or lengthy*. Often a significant research result.
- C50 *Open research problem*. An obtained solution should be published.

The rating is possibly supplemented by the following qualifier(s):

- i* Especially *interesting/instructive* problem.
- m* Requires more/higher *math* than used or developed here.
- o* *Open* (unsolved) problem, could be worth publishing.
- s* *Solved* problem with published solution.
- u* *Unpublished* result by the author.

The problems with *open* and especially those with *unpublished* solutions represent a significant original contribution to this thesis. They have been placed at the end of each chapter in order to keep the main text better focused.

Acknowledgements. I would like to thank all those people who in one way or another have contributed to the success of this thesis. After having spent over 4 years working in industry, I did not believe that I would find my way back to academia. Ray Solomonoff’s positive response to my first paper [Hut01c] in this area encouraged me to return to academia. Jürgen Schmidhuber was the one who gave me that chance. He strongly believed in me and my ideas from the very beginning and organized an SNF grant (2000-61847.00) which funded this project for the last 2 years. I enjoyed the many stimulating, sometimes philosophical, discussions with him. Wilfried Brauer was so kind to agree to be the promoter of my thesis and supported me in various ways. For interesting discussions I am indebted to Leonid Levin, Peter Gács, Paul Vitányi, Richard Sutton, Leslie Kaelbling, Peter van Emde Boas, and many others. Shane Legg, Jan Poland, Viktor Zhumatiy, Douglas Eck, Ivo Kwee, Philippa Hutter, Paul Vitányi, and Jürgen Schmidhuber read and gave valuable feedback on (parts of) drafts of the thesis. Thanks also collectively to all other IDSIA’ies for the pleasant working atmosphere and their support. Thanks to my father, who taught me to think sharply and to my mother who taught me to do what one enjoys. Finally, I would like to apologize to my wife and my daughter who suffered most from my decision to leave Munich and to do a thesis.

Contents

0	Meta Contents	1
	Preface	1
	Table of Contents	4
	List of Tables, Figures, Theorems,	5
	List of Notation	7
1	A Short Tour through the Thesis	101
1.1	Introduction	102
1.2	Simplicity & Uncertainty	103
1.2.1	Introduction	103
1.2.2	Algorithmic Information Theory	104
1.2.3	Uncertainty & Probabilities	105
1.2.4	Algorithmic Probability & Universal Induction	105
1.2.5	Generalized Universal (Semi)Measures	106
1.3	Universal Sequence Prediction	107
1.3.1	Setup & Convergence	107
1.3.2	Loss Bounds	108
1.3.3	Optimality Properties	109
1.3.4	Miscellaneous	109
1.4	Rational Agents in known Probabilistic Environments	110
1.4.1	The Agent Model	110
1.4.2	Value Functions and Optimal Policies	111
1.4.3	Sequential Decision Theory & Reinforcement Learning	112
1.5	The Universal Algorithmic Agent AIXI	112
1.5.1	The Universal AIXI Model	113
1.5.2	On the Optimality of AIXI	113
1.5.3	Value Related Optimality Results	114
1.5.4	Markov Decision Processes	116
1.5.5	The Choice of the Horizon	117
1.6	Important Environmental Classes	118
1.6.1	Introduction	118
1.6.2	Sequence Prediction (SP)	118
1.6.3	Strategic Games (SG)	118

1.6.4	Function Minimization (FM)	119
1.6.5	Supervised Learning from Examples (EX)	119
1.6.6	Other Aspects of Intelligence	119
1.7	Computational Aspects	119
1.7.1	The Fastest & Shortest Algorithm for All Problems	120
1.7.2	Time Bounded AIXI Model	122
1.8	Discussion	124
1.9	History & References	126
2	Simplicity & Uncertainty	201
2.1	Introduction	202
2.1.1	Ockham, Epicurus, Hume, Bayes, Solomonoff	202
2.1.2	Problem Setup	203
2.2	Algorithmic Information Theory	204
2.2.1	Definitions and Notation	204
2.2.2	Turing Machines	205
2.2.3	Kolmogorov Complexity	207
2.2.4	Computability Concepts	209
2.3	Uncertainty & Probabilities	211
2.3.1	Frequency Interpretation / Counting	212
2.3.2	Objective Interpretation: Probabilities for Uncertain Events	212
2.3.3	Subjective Interpretation: Probabilities for Degrees of Belief	214
2.3.4	Determining Priors	215
2.4	Algorithmic Probability & Universal Induction	215
2.4.1	The Universal Prior M	216
2.4.2	Universal Sequence Prediction	217
2.4.3	Universal (Semi)Measures	218
2.4.4	Martin-Löf Randomness	224
2.5	History & References	225
2.6	Problems	229
3	Universal Sequence Prediction	301
3.1	Introduction	303
3.1.1	Induction	303
3.1.2	Universal Sequence Prediction	304
3.2	Setup and Convergence	304
3.2.1	Random sequences	304
3.2.2	Universal Prior Probability Distribution	305
3.2.3	Universal Posterior Probability Distribution	306
3.2.4	Convergence of Random Sequences	307
3.2.5	Distance Measures between Probability Distributions	308
3.2.6	Convergence of ξ to μ	310
3.2.7	Convergence in Martin-Löf Sense	312

3.2.8	The case where $\mu \notin \mathcal{M}$	316
3.2.9	Probability Classes \mathcal{M}	316
3.3	Error Bounds	317
3.3.1	Bayes-Optimal Predictors	317
3.3.2	Total Expected Numbers of Errors	318
3.3.3	Proof of Theorem 3.36	319
3.4	Loss Bounds	321
3.4.1	Unit Loss Function	321
3.4.2	Loss Bound of Merhav & Feder	323
3.4.3	Example Loss Functions	324
3.4.4	Proof of Theorem 3.48	324
3.4.5	Convergence of Instantaneous Losses	326
3.4.6	General Loss	327
3.5	Application to Games of Chance	328
3.5.1	Introduction	328
3.5.2	Games of Chance	328
3.5.3	Example	329
3.5.4	Information-Theoretic Interpretation	330
3.6	Optimality Properties	330
3.6.1	Lower Error Bound	330
3.6.2	Pareto Optimality of ξ	333
3.6.3	Balanced Pareto Optimality of ξ	335
3.6.4	On the Optimal Choice of Weights	336
3.6.5	Occam's razor versus No Free Lunches	337
3.7	Miscellaneous	337
3.7.1	Multi-Step Predictions	337
3.7.2	Continuous Probability Classes \mathcal{M}	339
3.7.3	Further Applications	341
3.7.4	Prediction with Expert Advice	341
3.7.5	Outlook	343
3.8	Summary	344
3.9	Technical Proofs	345
3.9.1	How to Deal with $\mu=0$	345
3.9.2	Entropy Inequalities (3.11)	346
3.9.3	Error Inequality (3.36)	347
3.9.4	Binary Loss Inequality for $z \leq \frac{1}{2}$ (3.57)	349
3.9.5	Binary Loss Inequality for $z \geq \frac{1}{2}$ (3.58)	350
3.9.6	General Loss Inequality (3.53)	350
3.10	History & References	351
3.11	Problems	352
4	Agents in Known Probabilistic Environments	401
4.1	The $\text{AI}\mu$ Model in Functional Form	402

4.1.1	The Cybernetic Agent Model	402
4.1.2	Strings	403
4.1.3	AI model for Known Deterministic Environment	405
4.1.4	AI Model for Known Prior Probability	406
4.2	The $AI\mu$ Model in Recursive and Iterative Form	408
4.2.1	Probability Distributions	408
4.2.2	Explicit Form of the $AI\mu$ Model	409
4.2.3	Equivalence of Functional and Explicit AI model	410
4.3	Special Aspects of the $AI\mu$ Model	412
4.3.1	Factorizable Environments	412
4.3.2	Constants and Limits	414
4.3.3	Sequential Decision Theory	415
4.4	Problems	416
5	The Universal Algorithmic Agent AIXI	501
5.1	The Universal AIXI Model	502
5.1.1	Definition of the AIXI Model	503
5.1.2	Universality of ξ^{AI}	504
5.1.3	Convergence of ξ^{AI} to μ^{AI}	505
5.1.4	Intelligence Order Relation	506
5.2	On the Optimality of AIXI	507
5.3	Value Bounds and Separability Concepts	509
5.3.1	Introduction	509
5.3.2	(Pseudo) Passive μ and the HeavenHell Example	509
5.3.3	The OnlyOne Example	510
5.3.4	Asymptotic Learnability	511
5.3.5	Uniform μ	512
5.3.6	Other Concepts	512
5.3.7	Summary	512
5.4	Value Related Optimality Results	513
5.4.1	The $AI\rho$ Models: Preliminaries	513
5.4.2	Pareto Optimality of $AI\xi$	514
5.4.3	Self-optimizing Policy p^ξ w.r.t. Average Value	515
5.5	Discounted Future Value Function	518
5.6	Markov Decision Processes (MDP)	524
5.7	The Choice of the Horizon	527
5.8	Outlook	530
5.9	Conclusions	531
5.10	Functions \leadsto Chronological Semimeasures	531
5.11	Proof of the Entropy Inequality	533
5.12	History & References	535
5.13	Problems	535

6	Important Environmental Classes	601
6.1	Repetition of the $AI\mu/\xi$ Models	602
6.2	Sequence Prediction (SP)	603
6.2.1	Using the $AI\mu$ Model for Sequence Prediction	604
6.2.2	Using the $AI\xi$ Model for Sequence Prediction	606
6.3	Strategic Games (SG)	607
6.3.1	Introduction	608
6.3.2	Strictly Competitive Strategic Games	608
6.3.3	Using the $AI\mu$ Model for Game Playing	609
6.3.4	Games of Variable Length	610
6.3.5	Using the $AI\xi$ Model for Game Playing	611
6.4	Function Minimization (FM)	612
6.4.1	Applications/Examples	612
6.4.2	The Greedy Model $FMG\mu$	613
6.4.3	The General $FM\mu/\xi$ Model	614
6.4.4	Is the General Model Inventive?	616
6.4.5	Using the AI models for Function Minimization	617
6.4.6	Remark	618
6.5	Supervised Learning from Examples (EX)	619
6.5.1	Applications/Examples	619
6.5.2	Supervised Learning with the $AI\mu/\xi$ Model	619
6.6	Other Aspects of Intelligence	621
6.7	Problems	622
7	Computational Aspects	701
7.1	The Fastest & Shortest Algorithm for All Problems	702
7.1.1	Introduction & Main Result	702
7.1.2	Levin Search	704
7.1.3	Fast Matrix Multiplication	705
7.1.4	Applicability of the Fast Algorithm $M_{p^*}^\varepsilon$	706
7.1.5	The Fast Algorithm $M_{p^*}^\varepsilon$	707
7.1.6	Time Analysis	708
7.1.7	Assumptions on the Machine Model	709
7.1.8	Algorithmic Complexity and the Shortest Algorithm	710
7.1.9	Generalizations	711
7.1.10	Summary & Outlook	712
7.2	Time Bounded AIXI Model	713
7.2.1	Introduction	713
7.2.2	Time Limited Probability Distributions	714
7.2.3	The Idea of the Best Vote Algorithm	715
7.2.4	Extended Chronological Programs	716
7.2.5	Valid Approximations	716
7.2.6	Effective Intelligence Order Relation	717

7.2.7	The Universal Time Bounded AIXI \tilde{t} Agent	718
7.2.8	Limitations and Open Questions	719
7.2.9	Remarks	719
8	Discussion	801
8.1	What has been Achieved	802
8.1.1	Results	802
8.1.2	Comparison to other Approaches	803
8.2	General Remarks	804
8.2.1	Miscellaneous	805
8.2.2	Prior Knowledge	806
8.2.3	Universal Prior Knowledge	806
8.2.4	How AIXI(t) Deals with Encrypted Information	807
8.2.5	Mortal Embodied Agents	807
8.3	Personal Remarks	808
8.3.1	On the Foundations of Machine Learning	809
8.3.2	In a World without Occam	810
8.4	Outlook & Open Questions	810
8.5	Assumptions, Problems, Limitations	812
8.5.1	Assumptions	812
8.5.2	Problems	813
8.5.3	Limitations	814
8.6	Philosophical Issues	814
8.6.1	Turing Test	814
8.6.2	On the Existence of Objective Probabilities	815
8.6.3	Free will versus Determinism	815
8.6.4	The Big Questions	817
8.7	Conclusions	818
	Appendix	901
	Bibliography	901
	Index	912

List of Tables, Figures, Theorems,

...

Table 2.2 ((Prefix) coding natural numbers and strings)	205
Thesis 2.3 (Turing)	206
Thesis 2.4 (Church)	206
Assumption 2.5 (Short compiler)	206
Definition 2.6 (Prefix/Monotone Turing machine)	206
Theorem 2.7 (Universal prefix/monotone Turing machine)	207
Definition 2.9 (Kolmogorov complexity)	208
Theorem 2.10 (Information properties of Kolmogorov complexity)	208
Definition 2.12 (Computable functions)	209
Figure 2.11 (Kolmogorov Complexity)	210
Theorem 2.13 ((Non)computability of Kolmogorov complexity)	211
Axioms 2.14 (Kolmogorov's axioms of probability theory)	212
Definition 2.15 (Conditional probability)	213
Theorem 2.16 (Bayes rule)	213
Axioms 2.17 (Cox's axioms for beliefs)	214
Theorem 2.18 (Cox's theorem)	214
Definition 2.20 ((Semi)measures)	216
Theorem 2.21 (Universality of M)	216
Theorem 2.26 (Universal (semi)measures)	219
Table 2.27 (Existence of universal (semi)measures)	220
Theorem 2.28 (Martin-Löf random sequences)	224
Definition 2.30 (μ/ξ -random sequences)	224
Definition 3.8 (Convergence of random sequences)	307
Lemma 3.9 (Relations between random convergence criteria)	308
Lemma 3.11 (Entropy Inequalities)	308
Theorem 3.19 (Convergence of ξ to μ)	310
Theorem 3.22 (μ/ξ -convergence of ξ to μ)	312
Theorem 3.36 (Error bound)	318
Theorem 3.48 (Unit loss bound)	322
Corollary 3.49 (Unit loss bound)	322

Theorem 3.59 (Instantaneous Loss Bound)	326
Theorem 3.60 (General loss bound)	327
Theorem 3.63 (Time to Win)	329
Theorem 3.64 (Lower Error Bound)	331
Definition 3.65 (Pareto Optimality)	333
Theorem 3.66 (Pareto Optimality)	333
Theorem 3.69 (Balanced Pareto Optimality w.r.t. L)	335
Theorem 3.70 (Optimality of universal weights)	337
Theorem 3.74 (Continuous Entropy Bound)	340
Definition 4.1 (The Agent Model)	402
Table 4.2 (Notation and emphasis in AI versus control theory)	404
Definition 4.4 (The $\text{AI}\mu$ model)	406
Definition 4.5 (The μ /true/generating value function)	406
Figure 4.13 (Expectimax Tree/Algorithm for $\mathcal{X}' = \mathcal{Y} = \mathcal{B}$)	410
Theorem 4.20 (Equivalence of functional and explicit AI model)	410
Theorem 4.25 (Factorizable environments μ)	413
Assumption 4.28 (Finiteness)	414
Claim 5.12 (We expect AIXI to be universally optimal)	506
Definition 5.14 (Intelligence order relation)	506
Definition 5.18 (ρ -Value function)	513
Definition 5.19 (Functional $\text{AI}\rho$ model)	513
Theorem 5.20 (Iterative $\text{AI}\rho$ model)	513
Theorem 5.21 (Linearity and convexity of V_ρ in ρ)	513
Definition 5.22 (Pareto Optimality)	514
Theorem 5.23 (Pareto Optimality)	514
Theorem 5.24 (Balanced Pareto Optimality)	515
Lemma 5.27 (Value difference relation)	516
Lemma 5.28 (Convergence of averages)	516
Theorem 5.29 (Self-optimizing policy p^ξ w.r.t. average value)	517
Definition 5.30 (Discounted $\text{AI}\rho$ model and value)	519
Theorem 5.31 (Linearity and convexity of V_ρ in ρ)	519
Theorem 5.32 (Pareto Optimality)	520
Lemma 5.33 (Value difference relation)	520
Theorem 5.34 (Self-optimizing policy p^ξ w.r.t. discounted value)	520
Theorem 5.35 (Continuity of discounted value)	521
Theorem 5.36 (Convergence of universal to true Value)	523
Definition 5.37 (Ergodic Markov Decision Processes)	524
Theorem 5.38 (Self-optimizing policies for ergodic MDPs)	524
Corollary 5.40 ($\text{AI}\xi$ is self-optimizing for ergodic MDPs)	526
Table 5.41 (Effective horizons)	528
Theorem 7.1 (The fastest algorithm)	703

Theorem 7.2 (The fastest & shortest algorithm)	711
Definition 7.8 (Effective intelligence order relation)	717
Theorem 7.9 (Optimality of AIXItl)	718
Table 8.1 (Properties of learning algorithms)	805

List of Notation

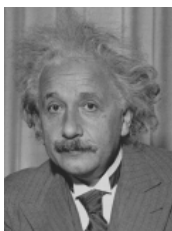
The following is a list of commonly used notation. The first entry is the symbol itself, followed by its meaning or name (if any), and the page number where the definition appears. Some standard symbols like \mathbb{R} are not defined in the text. There is a * in place of the page number for these symbols.

Symbol	Explanation	Page
∞	infinity	*
$\{a, \dots, z\}$	set containing elements a, b, \dots, y, z	*
$[a, b)$	interval on the real line, closed at a and open at b	*
[C35s]	Classification of Problems	4
[LV97]	Paper, book or other reference	*
(5.3)	Reference to formula	*
$\cap, \cup, \setminus, \in$	set intersection, union, difference, membership	*
\wedge, \vee, \neg	Boolean conjunction (and), disjunction (or), negation (not)	*
\subseteq, \subset	subset, proper subset	*
\Rightarrow	implies	*
\Leftrightarrow	equivalence, if and only if, iff	*
\square	q.e.d. (Latin), which was to be demonstrated	*
\forall, \exists	for all, there exists	*
$\approx, \lesseqgtr, \gtrless$	approximately equal, less equal, greater equal	205
\ll, \gg	much small/greater than	*
\equiv	equivalent, identical, equal by definition	*
\cong	isomorphic	*
$:=$	define as	*
\doteq	corresponds to, informal equality	*
\sim	asymptotic equality	205
\propto	proportional to	*
$=, \neq$	equal to, not equal to	*
$+, -, \cdot, /$	standard arithmetic operations: sum, difference, product, ratio	*
$\sqrt{}$	square root	*
$\leq, \geq, <, >$	standard inequalities	*

$ \mathcal{S} , a $	size/cardinality of set \mathcal{S} , absolute value of a	*
\rightarrow	mapping, approaches, Boolean implication	*
\rightarrow	converge to each other	205
$\lim_{n \rightarrow \infty}$	limes value of argument for n tending to infinity	*
\leadsto	replace with	*
$\lceil x \rceil$	ceiling of x : smallest integer larger or equal than x	*
$\lfloor x \rfloor$	floor of x : largest integer smaller or equal than x	*
$\sum_{k=1}^n$	summation from $k=1$ to n	*
$\prod_{k=1}^n$	product from $k=1$ to n	*
$\int, \int_a^b dx$	Lebesgue integral, integral from a to b over x	*
i, k, n	natural numbers	205
x, y, z	strings	205
min/max	minimal/maximal element of set: $\min_{x \in \mathcal{X}} f(x) = \min\{f(x) : x \in \mathcal{X}\}$	*
argmin	$\operatorname{argmin}_x f(x)$ is the x that minimizes $f(x)$ (ties broken arbitrarily)	*
l.h.s.	left hand side	*
r.h.s.	right hand side	*
w.r.t.	with respect to	*
e.g.	exempli gratia (Latin), for example	*
i.e.	id est (Latin), that is	*
etc.	et cetera (Latin), and so forth	*
cf.	confer (Latin, imperative of conferre), compare with	*
ibid.	ibidem (Latin), in the same place	*
viz.	videlicet (Latin), that is to say, namely	*
et al.	et alii (Latin), and others	*
q.e.d.	quod erat demonstrandum (Latin), which was to be demonstrated	*
i.i.d.	independent identically distributed (random variables)	*
iff	if and only if	*
w.p.1/i.p.	with probability 1 / in probability	307
i.m./i.m.s.	in the mean / in mean sum	307
log	logarithm to some basis (often 2)	*
\log_b	logarithm to basis b	*
ln	natural logarithm to basis $e=2.71828\dots$	*
e	base of natural logarithm $e=2.71828\dots$	*
\mathbb{R}	Set of real numbers	*
\mathbb{N}	Set of natural numbers $\{1, 2, 3, \dots\}$	204
\mathbb{N}_0	Set of natural numbers including zero $\{0, 1, 2, 3, \dots\}$	204

\mathcal{Q}	Set of rational numbers $\{\frac{n}{d}\}$	*
$\mathcal{B}=\{0,1\}$	Binary alphabet.	*
$y_t \in \mathcal{Y}$	Action (output of agent) in cycle t , followed by ...	403
$x_t \in \mathcal{X}$	Perception (feedback/input to agent) in cycle t .	215, 403
$x'_t \in \mathcal{X}'$	informative input/perception in cycle t .	403
$x'_t \in \mathcal{X}$	Potential perception in future cycle t .	506
$r_t \in \mathcal{R} \subset \mathbb{R}$	reward in cycle t .	403
ε	some small positive real number	*
ϵ	empty string	204
$*$	wildcard for some string (prefix, finite, or infinite).	204
$x_{1:n}$	$=x_1...x_n$ = String of length n .	215, 304, 403
$x_{<t}$	$=x_1...x_{t-1}$ = String of length $t-1$.	215, 304, 403
$y x_{k:n}$	Action-perception sequence $y_k x_k ... y_n x_n$.	403
$\dot{y} \dot{x}_{<k}$	Actually realized action-perception sequence $\dot{y}_1 \dot{x}_1 ... \dot{y}_{k-1} \dot{x}_{k-1}$.	406
ω	infinite string/sequence	204
S, Ω	Sample space	213, 304
$\Gamma_{x_{1:n}}$	$=\{\omega: \omega_{1:n}=x_{1:n}\}$ = cylinder set	216, 304
$l(x)$	length of string x	204
$\langle o \rangle$	Coding of object o	204
$\langle x, y \rangle$	uniquely decodable pairing of x and y	204
x'	prefix coding of x	204
$O(), o()$	big and small oh-notation	205
$a \stackrel{+}{\leq} b$	less within an additive constant, i.e. $a \leq b + O(1)$. Similarly $\stackrel{\pm}{\leq}$.	204
$a \stackrel{\times}{\leq} b$	less within a multiplicative constant, i.e. $a = O(b)$. Similarly $\stackrel{\times}{\leq}$.	204
$K(x)$	prefix Kolmogorov complexity of string x .	208
$Km(x_{1:n})$	monotone (Kolmogorov) complexity of string $x_{1:n}$	217, 606
$K(o_1 o_2)$	Kolmogorov complexity of object o_1 from object o_2 .	208
$M \stackrel{\times}{\sim} \xi_U$	Solomonoff's universal semimeasure	216
$\mathcal{M}=\{\nu\}$	(Usually countable or finite) set of probability measures.	218, 316
EC	$\in \{\text{AI, SP, FM, EX, SG, ...}\}$ is an Environmental Class.	*
AI	Algorithmic Intelligence (Most general computable env. class).	*
SP	Sequence Prediction.	603
CF	Classification.	341
SG	Strategic 2 player informed zero-sum Games.	607
FM	Function Minimization.	612
EX	Supervised Learning (by Examples).	619

pd	Probability density function / distribution / measure.	*
$\rho(x_{1:n})$	Probability of string/sequence starting with $x_{1:n}$.	216, 304
$\mu \in \mathcal{M}$	True generating environmental pd.	304
E	Expectation value, usually w.r.t. the true distribution μ .	304
P	Probability, usually w.r.t. the true distribution μ .	304
$\mu(x_1 x_2 x_3 x_4)$	μ probability that the 2^{nd} and 4^{th} symbols of a string are x_2 and x_4 , given the 1^{st} and 3^{rd} symbols are x_1 and x_3 .	408
$\nu \in \mathcal{M}$	Any pd in \mathcal{M} .	306
ρ	Any pd not necessarily in \mathcal{M} usually specifying a policy.	304
ξ	$= \sum_{\nu \in \mathcal{M}} w_\nu \nu =$ mixture (belief) pd.	218, 306
w_ν	Prior degree of belief in ν – or – weight of ν .	218, 306
ρ^{EC}	Pd of environmental argument type EC.	601
ξ^{EC}	Mixture distribution of type EC for class EC.	601
$\ell_{x_t y_t}$	Incurred loss when predicting/acting y_t and x_t is next symbol.	321
$l_{t\nu}^\Lambda$	ν -expected instant. loss in step t achieved by predictor Λ .	334, 322
$L_{n\nu}^\Lambda$	ν -expected cumulative loss of steps $1...n$ achieved by predictor Λ .	334
Λ_ρ	Predictor which minimizes the ρ -expected loss.	322
$e_{t\nu}^\Theta$	ν probability that Θ -predictor errs in step t .	318
$E_{n\nu}^\Theta$	ν expected # of errors in steps $1...n$ achieved by predictor Θ .	318
$L_n^\Lambda \equiv L_{n\mu}^\Lambda$...Abbreviation for true μ -expected loss, ...	321
$V_{km}^{p\nu}(\dot{y}_{<k})$	Value of policy p in environment ν given history $\dot{y}_{<k}$.	513
y_t^Λ	prediction/decision/action of predictor Λ	321
y_k^p	action of policy p	*
γ_k	discounting sequence	519
Γ_k	value function normalization ($\sum_{i=k}^\infty \gamma_i$)	519
m	Agents lifespan / horizon.	405
p	Agents policy.	402
q	Deterministic environment.	402
p^ν	Policy which maximizes value V_ν^p .	406
$V_\mu^* \equiv V_{1m}^{p^\mu \mu}$	True or generating value.	406
$V_\xi^* \equiv V_{1m}^{p^\xi \xi}$	Universal value.	506
$D_n \equiv D_{n\mu}^\xi$	Relative entropy between μ and ξ for the first n cycles.	309



Albert Einstein
(1879-1955)

“I have no particular talent. I am merely inquisitive.” (Albert Einstein)

Chapter 1

A Short Tour through the Thesis

1.1	Introduction	102
1.2	Simplicity & Uncertainty	103
1.2.1	Introduction	103
1.2.2	Algorithmic Information Theory	104
1.2.3	Uncertainty & Probabilities	105
1.2.4	Algorithmic Probability & Universal Induction	105
1.2.5	Generalized Universal (Semi)Measures	106
1.3	Universal Sequence Prediction	107
1.3.1	Setup & Convergence	107
1.3.2	Loss Bounds	108
1.3.3	Optimality Properties	109
1.3.4	Miscellaneous	109
1.4	Rational Agents in known Probabilistic Environments	110
1.4.1	The Agent Model	110
1.4.2	Value Functions and Optimal Policies	111
1.4.3	Sequential Decision Theory & Reinforcement Learning	112
1.5	The Universal Algorithmic Agent AIXI	112
1.5.1	The Universal AIXI Model	113
1.5.2	On the Optimality of AIXI	113
1.5.3	Value Related Optimality Results	114
1.5.4	Markov Decision Processes	116
1.5.5	The Choice of the Horizon	117

1.6	Important Environmental Classes	118
1.6.1	Introduction	118
1.6.2	Sequence Prediction (SP)	118
1.6.3	Strategic Games (SG)	118
1.6.4	Function Minimization (FM)	119
1.6.5	Supervised Learning from Examples (EX)	119
1.6.6	Other Aspects of Intelligence	119
1.7	Computational Aspects	119
1.7.1	The Fastest & Shortest Algorithm for All Problems	120
1.7.2	Time Bounded AIXI Model	122
1.8	Discussion	124
1.9	History & References	126

This Chapter represents a short tour through the thesis. It is not meant as a gentle introduction for novices, but as a condensed presentation of the most important concepts and results of the thesis. The price for this brevity is that in this chapter we mostly forgo mathematical rigor, subtleties, proofs, discussions, references and comparisons to other work. More seriously some sections demand high background knowledge. Readers unfamiliar with algorithmic information theory should first read Chapter 2 or consult the textbook [LV97]. Readers unfamiliar with sequential decision theory should first read Chapter 4 or consult the textbooks [BT96, SB98]. Before discouraged by the complexity of some of the sections, it is better to skip them completely.

1.1 Introduction

Artificial Intelligence. The science of Artificial Intelligence (AI) might be defined as the construction of intelligent systems and their analysis. A natural definition of a *system* is anything which has an input and an output stream. Intelligence is more complicated. It can have many faces like creativity, solving problems, pattern recognition, classification, learning, induction, deduction, building analogies, optimization, surviving in an environment, language processing, knowledge and many more. A formal definition incorporating every aspect of intelligence, however, seems difficult. Further, intelligence is graded, there is a smooth transition between systems, which everyone would agree to be not intelligent and truly intelligent systems. One simply has to look in nature, starting with, for instance, inanimate crystals, then come amino-acids, then some RNA fragments, then viruses, bacteria, plants, animals, apes, followed by the truly intelligent homo sapiens, and possibly continued by AI systems or ETs. So the best we can expect to find is a partial or total order relation on the set of systems, which orders them w.r.t. their degree of intelligence (like intelligence tests do for human systems, but for a limited class of problems).

Having this order we are, of course, interested in large elements, i.e. highly intelligent systems. If a largest element exists, it would correspond to the most intelligent system which could exist.

Most, if not all known facets of intelligence can be formulated as goal driven or, more precisely, as maximizing some utility function. It is, therefore, sufficient to study goal driven AI. E.g. the (biological) goal of animals and humans is to survive and spread. The goal of AI systems should be to be useful to humans. The problem is that, except for special cases, we know neither the utility function, nor the environment in which the agent will operate, in advance.

Main idea. We propose a theory which formally¹ solves the problem of unknown goal and environment. It might be viewed as a unification of the ideas of universal induction, probabilistic planning and reinforcement learning or as a unification of sequential decision theory with algorithmic information theory. We apply this model to some of the facets of intelligence, including induction, game playing, optimization, reinforcement and supervised learning, and show how it solves these problem classes. This, together with general convergence theorems motivates us to believe that the constructed universal AI system is the best one in a sense to be clarified in the sequel, i.e. that it is the most intelligent environment independent system possible. The intention of this work is to introduce the universal AI model and give an in breadth analysis.

1.2 Simplicity & Uncertainty

This section introduces Occam's razor principle, Kolmogorov complexity, objective/subjective probabilities, to finally arrive at the problem of universal prediction, and its solution due to Solomonoff.

1.2.1 Introduction

An important and highly non-trivial aspect of intelligence is inductive inference. Simply speaking, induction is the process of predicting the future from the past or, more precisely, it is the process of finding rules in (past) data and using these rules to guess future data. Weather or stock-market forecasting, or continuing number series in an IQ test, are non-trivial examples. Making good predictions plays a central role in natural and artificial intelligence in general, and in machine learning in particular. All induction problems can be phrased as sequence prediction tasks. This is, for instance, obvious for time series prediction, but also includes classification tasks. Having observed data x_t at times $t < n$, the task is to predict the n^{th} symbol x_n from sequence $x_1 \dots x_{n-1}$. This *prequential approach* [Daw84] skips over the intermediate step of learning a model based on observed data $x_1 \dots x_{n-1}$ and then using this model

¹With a formal solution we mean a rigorous mathematical definition, uniquely specifying the solution. In the following, a solution is always meant in this formal sense.

to predict x_n . The prequential approach avoids problems of model consistency, how to separate noise from useful data, and many other issues. The goal is to make “good” predictions, where the prediction quality is usually measured by a loss function, which shall be minimized. The key concept to well-define and solve induction problems is *Occam’s razor* (simplicity) principle, which says that “*Entities should not be multiplied beyond necessity*,” which may be interpreted as to keep the simplest theory consistent with the observations $x_1 \dots x_{n-1}$ and to use this theory to predict x_n . Before we can present Solomonoff’s formal solution, we have to quantify Occam’s razor in terms of Kolmogorov complexity, and introduce the notion of subjective/objective probabilities.

1.2.2 Algorithmic Information Theory

Intuitively a string is simple if it can be described in a few words, like “the string of one million ones”, and is complex if there is no such short description, like for a random string whose shortest description is specifying it bit-by-bit. We can restrict the discussion to binary strings, since for other (non-stringy mathematical) objects we may assume some default coding as binary strings. Furthermore, we are only interested in effective descriptions, and hence restrict decoders to be Turing machines. Let us choose some universal (so called prefix) *Turing machine* U with unidirectional binary input and output tapes and a bidirectional work tape. We can then define the (conditional) *prefix Kolmogorov complexity* [Cha75, Gác74, Kol65, Lev74] of a binary string x as the length l of the shortest program p , for which U outputs the binary string x (given y)

$$K(x) := \min_p \{l(p) : U(p) = x\} \quad K(x|y) := \min_p \{l(p) : U(p, y) = x\}.$$

Simple strings like $000\dots 0$ can be generated by short programs, and, hence have low Kolmogorov complexity, but irregular (e.g. random) strings are their own shortest description, and hence have high Kolmogorov complexity. An important property of K is that it is nearly independent of the choice of U . Furthermore it shares many properties with Shannon’s entropy (information measure) S , but K is superior to S in many respects. Figure 2.11 on page 210 contains a schematic graph of K . To be brief, K is an excellent universal complexity measure, suitable for quantifying Occam’s razor. There is (only) one severe disadvantage: K is not finitely computable. More precisely, a function f is said to be *finitely computable* (or *recursive*) if there exists a Turing machine which, given x , computes $f(x)$ and then halts. Some functions are not finitely computable but still *approximable* in the sense that there is a non-halting Turing machine with an infinite output sequence y_1, y_2, y_3, \dots and $\lim_{t \rightarrow \infty} y_t = f(x)$. If additionally the output sequence is monotone increasing/decreasing, then f is said to be *lower/upper semi-computable* (or *enumerable/co-enumerable*). Finally we call f *estimable* if some Turing machine, given x and a precision ε , finitely computes an ε -approximation of x . The major algorithmic property of K is that it is co-enumerable, but not finitely computable.

1.2.3 Uncertainty & Probabilities

²For the *objectivist* probabilities are real aspects of the world. The outcome of an observation or an experiment is not deterministic, but involves physical random processes. Kolmogorov's axioms of probability theory formalize the properties which probabilities should have. In the case of i.i.d. experiments the probabilities assigned to events can be interpreted as limiting frequencies (*frequentist* view), but applications are not limited to this case. Conditionalizing probabilities and the chain rule are the major tools in computing posterior probabilities from prior ones. For instance, given the initial binary sequence $x_1 \dots x_{n-1}$, what is the probability of the next bit being 1? The probability of observing x_n at time n , given past observations $x_1 \dots x_{n-1}$ can be computed with the chain rule³ if the true generating distribution μ of the sequences $x_1 x_2 x_3 \dots$ is known: $\mu(x_n | x_{<n}) = \mu(x_{1:n}) / \mu(x_{<n})$, where we introduced the abbreviations $x_{1:n} \equiv x_1 x_2 \dots x_n$ and $x_{<n} \equiv x_1 x_2 \dots x_{n-1}$. The problem, however, is that one often does not know the true distribution μ (e.g. in the cases of weather and stock-market forecasting).

The *subjectivist* uses probabilities to characterize an agent's degree of belief in (or plausibility of) something, rather than to characterize physical random processes. This is the most relevant interpretation of probabilities in AI. It is somewhat surprising that plausibilities can be shown to also respect Kolmogorov's axioms of probability and the chain rule by assuming only a few plausible qualitative rules they should follow [Cox46]. Hence, if the plausibility of $x_{1:n}$ is $\rho(x_{1:n})$, the degree of belief in x_n given $x_{<n}$ is, again, given by the chain rule: $\rho(x_n | x_{<n}) = \rho(x_{1:n}) / \rho(x_{<n})$.

The chain rule allows computing posterior probabilities/plausibilities from prior ones, but leaves open the question of how to determine the priors themselves. In statistical physics, the principle of indifference (symmetry principle) and the maximum entropy principle can often be exploited to determine prior probabilities, but only Occam's razor is general enough to assign prior probabilities in *every* situation, especially to cope with domains typical for AI.

1.2.4 Algorithmic Probability & Universal Induction

Occam's razor (appropriately interpreted and in compromise with Epicurus' principle of indifference) tells us to assign high/low a priori plausibility to simple/complex strings x . Using K as complexity measure any monotone decreasing function of K , e.g. $\rho(x) = 2^{-K(x)}$ would satisfy this criterion. But ρ also has to satisfy the probability axioms, so we have to be a bit more careful. Solomonoff [Sol64, Sol78] defined the *universal prior* $M(x)$ as the probability that the output of a universal Turing machine U starts with x when provided with fair coin flips on the input tape. Formally,

²Readers not believing in objective and/or subjective probabilities should read the remark at the beginning of Section 2.3.

³Strictly speaking it is just the definition of conditional probabilities.

M can be defined as

$$M(x) := \sum_{p: U(p)=x*} 2^{-l(p)} \geq 2^{-K(x)} \quad (1.1)$$

where the sum is over all (so called minimal) programs p for which U outputs a string starting with x . The inequality follows by dropping all terms in \sum_p except for the shortest p computing x . Strictly speaking M is only a *semimeasure* since it is not normalized to 1, but this is acceptable/correctable. We derive the following bound:

$$\sum_{t=1}^{\infty} (1 - M(x_t | x_{<t}))^2 \leq -\frac{1}{2} \sum_{t=1}^{\infty} \ln M(x_t | x_{<t}) = -\frac{1}{2} \ln M(x_{1:\infty}) \leq \frac{1}{2} \ln 2 \cdot Km(x_{1:\infty})$$

where $Km(x_{1:\infty})$ is the length of the shortest (non-halting) program computing $x_{1:\infty}$. In the first inequality we have used $(1-a)^2 \leq -\frac{1}{2} \ln a$ for $0 \leq a \leq 1$. In the equality we exchanged the sum with the logarithm and eliminated the resulting product by the chain rule. In the last inequality we exploited the inequality (1.1). If $x_{1:\infty}$ is a computable sequence, then $K(x_{1:\infty})$ is finite, which implies $M(x_t | x_{<t}) \rightarrow 1$ ($\sum_{t=1}^{\infty} (1-a_t)^2 < \infty \Rightarrow a_t \rightarrow 1$). This means, that if the environment is a computable sequence (whichsoever, e.g. the digits π or e in binary representation), after having seen the first few digits, M correctly predicts the next digit with high probability, i.e. it recognizes the structure of the sequence.

Assume now that the true sequence is drawn from the distribution μ , i.e. the true (objective) probability of $x_{1:n}$ is $\mu(x_{1:n})$, but μ is unknown. How is the posterior (subjective) belief $M(x_n | x_{<n}) = M(x_n) / M(x_{<n})$ related to the true (objective) posterior probability $\mu(x_n | x_{<n})$? Solomonoff's [Sol78] central result is that the posterior (subjective) beliefs converge to the true (objective) posterior probabilities, if the latter are computable. More precisely he showed that

$$\sum_{t=1}^{\infty} \sum_{x_{<t}} \mu(x_{<t}) \left(M(0 | x_{<t}) - \mu(0 | x_{<t}) \right)^2 \stackrel{+}{\leq} \frac{1}{2} \ln 2 \cdot K(\mu), \quad (1.2)$$

where we abbreviated $f(x) \leq g(x) + O(1)$ by $f(x) \stackrel{+}{\leq} g(x)$. $K(\mu)$ is finite if μ is computable, but the infinite sum on the l.h.s. can only be finite if the difference $M(0 | x_{<t}) - \mu(0 | x_{<t})$ tends to zero for $t \rightarrow \infty$ with μ probability 1 (w. μ .p.1). This shows that using M as an estimate for μ may be a reasonable thing to do.

1.2.5 Generalized Universal (Semi)Measures

One can derive a universal prior in a different way. We (loosely) define $\mathcal{M}_U := \{\nu_1, \nu_2, \dots\}$ to be the class of all (lower semi-)computable semimeasures. Let $\mu \in \mathcal{M}$ and assign (consistent with Occam's razor) a prior plausibility of $2^{-K(\nu_a)}$ to ν_a . Then the prior plausibility of $x_{1:n}$ is (by elementary probability theory)

$$\xi_U(x_{1:n}) := \sum_{\nu \in \mathcal{M}_U} 2^{-K(\nu)} \nu(x_{1:n}). \quad (1.3)$$

One can show that ξ_U coincides with M within an (irrelevant) multiplicative constant, i.e. $M(x) \stackrel{\times}{=} \xi_U(x)$, where $f(x) \stackrel{\times}{\leq} g(x)$ abbreviates $f(x) = O(g(x))$, and $\stackrel{\times}{=}$ denotes $\stackrel{\times}{\leq}$ and $\stackrel{\times}{\geq}$. Both ξ_U and M can be shown to be lower semi-computable. The dominance $M(x) \stackrel{\times}{=} \xi_U(x) \geq 2^{-K(\mu)} \mu(x)$ is the central ingredient in the proof of (1.2). The advantage of ξ_U over M is that the definition immediately generalizes to arbitrary weighted sums of (semi)measures in \mathcal{M} for arbitrary countable \mathcal{M} . Most proofs in this work go through for generic \mathcal{M} and weights.

So what is so special about the class of all enumerable semimeasures \mathcal{M}_U ? The larger we choose \mathcal{M} the less restrictive is the assumption that \mathcal{M} should contain the true distribution μ , which will be essential throughout the thesis. Why not restrict to the still rather general class of estimable or finitely computable (semi)measures? For *every* countable class \mathcal{M} , $\xi(x) := \xi_{\mathcal{M}}(x) := \sum_{\nu \in \mathcal{M}} w_{\nu} \nu(x)$ with $w_{\nu} > 0$, the important dominance $\xi(x) \geq w_{\nu} \nu(x)$ is satisfied. The question is what properties ξ possesses. The distinguishing property of \mathcal{M}_U is that ξ_U is itself an element of \mathcal{M}_U . On the other hand, in this work $\xi_{\mathcal{M}} \in \mathcal{M}$ is not by itself an important property. What matters is whether ξ is computable in one of the senses we defined above. There is an enumerable semimeasure (M) which dominates all enumerable semimeasures in \mathcal{M}_U . As we will see, there is *no* estimable semimeasure which dominates all computable measures, and there is *no* approximable semimeasure which dominates all approximable measures. From this it follows that for a universal (semi)measure which at least satisfies the weakest form of computability, namely being approximable, the largest dominated class among the classes considered in this thesis is the class of enumerable semimeasures, but there are even larger classes [Sch02a]. This is the reason why \mathcal{M}_U and M play a special role in this (and other) works. In practice one has to restrict to a finite subset of finitely computable environments ν to get a finitely computable ξ .

1.3 Universal Sequence Prediction

In the following we more closely investigate sequence prediction (SP) schemes based on Solomonoff's universal prior ξ_U and on more general Bayes-mixtures ξ , mainly from a decision theoretic perspective. In particular we show that they are optimal w.r.t. various optimality criteria.

1.3.1 Setup & Convergence

Let $\mathcal{M} := \{\nu_1, \nu_2, \dots\}$ be a countable set of candidate probability distributions on strings over the finite alphabet \mathcal{X} . We define a weighted average on \mathcal{M} :

$$\xi(x_{1:n}) := \sum_{\nu \in \mathcal{M}} w_{\nu} \cdot \nu(x_{1:n}), \quad \sum_{\nu \in \mathcal{M}} w_{\nu} = 1, \quad w_{\nu} > 0. \quad (1.4)$$

It is easy to see that ξ is a probability distribution as the weights w_{ν} are positive and normalized to 1 and the $\nu \in \mathcal{M}$ are probabilities. We call ξ universal relative

to \mathcal{M} , as it multiplicatively dominates all distributions in \mathcal{M} in the sense that $\xi(x_{1:n}) \geq w_\nu \cdot \nu(x_{1:n})$ for all $\nu \in \mathcal{M}$. In the following, we assume that \mathcal{M} is known and contains the true but unknown distribution μ , i.e. $\mu \in \mathcal{M}$. We abbreviate expectations w.r.t. μ by $\mathbf{E}[\cdot]$, for instance $\mathbf{E}[\cdot] = \sum_{x_{1:n} \in \mathcal{X}^n} \mu(x_{1:n})[\cdot]$ if $[\cdot]$ depends only on $x_{1:n}$, i.e. is independent of $x_{n+1:\infty}$. We use the (total) relative entropy D_n and squared Euclidian distance S_n to measure the distance between μ and ξ :

$$D_n := \mathbf{E} \left[\ln \frac{\mu(x_{1:n})}{\xi(x_{1:n})} \right], \quad S_n := \sum_{t=1}^n \mathbf{E} \left[\sum_{x'_t \in \mathcal{X}} \left(\mu(x'_t | x_{<t}) - \xi(x'_t | x_{<t}) \right)^2 \right]. \quad (1.5)$$

The following sequence of inequalities can be shown, which generalize Solomonoff's result (1.2): $S_n \leq D_n \leq \ln w_\mu^{-1} < \infty$. The finiteness of S_∞ implies $\xi(x'_t | x_{<t}) - \mu(x'_t | x_{<t}) \rightarrow 0$ for $t \rightarrow \infty$ w. μ .p.1 for any x'_t ($\sum_{t=1}^\infty s_t^2 < \infty \Rightarrow s_t \rightarrow 0$). We also show that $\sum_{t=1}^n \mathbf{E}[(\sqrt{\xi(x_t | x_{<t})/\mu(x_t | x_{<t})} - 1)^2] \leq D_n \leq \ln w_\mu^{-1} < \infty$ which implies $\xi(x_t | x_{<t})/\mu(x_t | x_{<t}) \rightarrow 1$ for $t \rightarrow \infty$ w. μ .p.1. This convergence motivates the belief that predictions based on (the known) ξ are asymptotically as good as predictions based on (the unknown) μ , with rapid convergence.

1.3.2 Loss Bounds

Most predictions are eventually used as a basis for some decision (or action), which itself leads to some reward or loss. Let $\ell_{x_t y_t} \in [0, 1] \subset \mathbb{R}$ be the received loss when performing prediction/decision/action $y_t \in \mathcal{Y}$ and $x_t \in \mathcal{X}$ is the t^{th} symbol of the sequence. Let $y_t^\Lambda \in \mathcal{Y}$ be the prediction of a (causal) prediction scheme Λ . The true probability of the next symbol being x_t , given $x_{<t}$, is $\mu(x_t | x_{<t})$. The expected loss when predicting y_t is $\mathbf{E}[\ell_{x_t y_t}]$. The total μ -expected loss suffered by the Λ scheme in the first n predictions is

$$L_n^\Lambda := \sum_{t=1}^n \mathbf{E}[\ell_{x_t y_t^\Lambda}]$$

The goal is to minimize the expected loss. More generally, we define the Λ_ρ sequence prediction scheme (later also called $\text{SP}\rho$) $y_t^{\Lambda_\rho} := \operatorname{argmin}_{y_t \in \mathcal{Y}} \sum_{x_t} \rho(x_t | x_{<t}) \ell_{x_t y_t}$ which minimizes the ρ -expected loss. If μ is known, Λ_μ is obviously the best prediction scheme in the sense of achieving minimal expected loss ($L_n^{\Lambda_\mu} \leq L_n^\Lambda$ for any Λ). We prove the following loss bound for the universal Λ_ξ predictor

$$0 \leq L_n^{\Lambda_\xi} - L_n^{\Lambda_\mu} \leq D_n + \sqrt{4L_n^{\Lambda_\mu} D_n + D_n^2} \leq 2D_n + 2\sqrt{L_n^{\Lambda_\mu} D_n} \quad (1.6)$$

Together with $L_n \leq n$ and $D_\infty \leq \ln w_\mu^{-1} < \infty$ this shows that $\frac{1}{n} L_n^{\Lambda_\xi} - \frac{1}{n} L_n^{\Lambda_\mu} = O(n^{-1/2})$, i.e. asymptotically Λ_ξ achieves the optimal average loss of Λ_μ with rapid convergence. Moreover $L_\infty^{\Lambda_\xi}$ is finite if $L_\infty^{\Lambda_\mu}$ is finite and $L_n^{\Lambda_\xi}/L_n^{\Lambda_\mu} \rightarrow 1$ if $L_\infty^{\Lambda_\mu}$ is not finite. Bound (1.6) also implies $L_n^\Lambda \geq L_n^{\Lambda_\xi} - 2\sqrt{L_n^{\Lambda_\xi} D_n}$, which shows that *no* (causal) predictor Λ whatsoever achieves significantly less (expected) loss than Λ_ξ . Note that for $w_\nu = 2^{-K(\nu)}$, $D_n \leq \ln 2 \cdot K(\mu)$ is of "reasonable" size. Instantaneous loss bounds can also be proven.

1.3.3 Optimality Properties

For any predictor Λ a worst case lower bound which asymptotically matches the upper bound (1.6) can be derived. More precisely let Λ be any deterministic predictor not knowing from which distribution $\mu \in \mathcal{M}$ the observed sequence $x_1 x_2 \dots$ is sampled from. Λ knows (depends on) \mathcal{M} , w_ν , and ℓ , and has at time t access to the previous outcomes $x_{<t}$. Then for every n there is an \mathcal{M} and $\mu \in \mathcal{M}$ and ℓ and weights w_ν such that

$$L_n^\Lambda - L_n^{\Lambda_\mu} = S_n + \sqrt{4L_n^{\Lambda_\mu} S_n + S_n^2} \quad \text{and} \quad D_n/S_n \rightarrow 1 \quad \text{for} \quad n \rightarrow \infty.$$

The equality especially holds for the universal predictor Λ_ξ . This shows that the bound (1.6) is tight in the sense that no other predictor can lead to significantly smaller bounds without making extra assumptions on \mathcal{M} , w_ν , or ℓ . For instance, for logarithmic and quadratic loss functions the regret $L_\infty^{\Lambda_\xi} - L_\infty^{\Lambda_\mu}$ is finite and bounded by $\ln w_\mu^{-1}$.

A different kind of optimality is *Pareto-optimality*. Let $\mathcal{F}(\mu, \rho)$ be any performance measure of ρ relative to μ . The universal prior ξ is called Pareto-optimal w.r.t. \mathcal{F} if there is no ρ with $\mathcal{F}(\nu, \rho) \leq \mathcal{F}(\nu, \xi)$ for all $\nu \in \mathcal{M}$ and strict inequality for at least one ν . We show that the universal prior ξ is Pareto-optimal w.r.t. the squared distance S_n , the relative entropy D_n , and the losses L_n . That is, for all performance measures which are relevant from a decision-theoretic point of view (i.e. for all loss functions ℓ) any improvement achieved by some predictor Λ_ρ over Λ_ξ in some environments ν is balanced by a deterioration in other environments. There are non-decision-theoretic performance measures w.r.t. which ξ is *not* Pareto-optimal. Pareto-optimality is a rather weak notion of optimality, but it emphasizes the distinctiveness of Bayes-mixture strategies.

Pareto-optimality of ξ still leaves open the question of how to choose the class \mathcal{M} and the weights w_ν . We have argued that \mathcal{M}_U is the largest \mathcal{M} suitable from a computational point of view. \mathcal{M}_U is also sufficiently large if we make the mild assumption that strings are sampled from a computable probability distribution. We show that within the class of enumerable weight functions with short program, the universal weights $w_\nu = 2^{-K(\nu)}$ lead to the smallest performance bounds within an additive (to $\ln w_\nu^{-1}$) constant in all enumerable environments. This argument justifies to select Solomonoff's prior (1.3) among all possible Bayes-mixtures⁴.

1.3.4 Miscellaneous

Games of Chance. The general loss bound (1.6) can, for instance, be used to estimate the time needed to reach the winning threshold in a game of chance (defined as a sequence of bets, observations and rewards). At time t we bet, depending on

⁴The reader who smells some free lunch here [WM97] should appease his hunger with Section 3.6.5.

the history $x_{<t}$, a certain amount of money s_t , take some action y_t , observe outcome x_t , and receive reward r_t . Our net profit, which we want to maximize, is $p_t = r_t - s_t \in [p_{\max} - p_{\Delta}, p_{\max}]$. The loss, which we want to minimize, can be identified with the negative (scaled) profit, $l_{x_k y_t} = (p_{\max} - p_t)/p_{\Delta}$. The Λ_{ρ} -system acts as to maximize the ρ -expected profit. Let $\bar{p}_{n\Lambda_{\rho}}$ be the average expected profit of the first n rounds. Bound (1.6) shows that the average profit of the Λ_{ξ} system converges to the best possible average profit $\bar{p}_{n\Lambda_{\mu}}$ achieved by the Λ_{μ} scheme ($\bar{p}_{n\Lambda_{\xi}} - \bar{p}_{n\Lambda_{\mu}} = O(n^{-1/2}) \rightarrow 0$ for $n \rightarrow \infty$). If there is a profitable scheme at all, then asymptotically the universal Λ_{ξ} scheme will also become profitable with the same average profit. We further show using ξ_U that $(\frac{2p_{\Delta}}{\bar{p}_{n\Lambda_{\mu}}})^2 \cdot \ln 2 \cdot K(\mu)$ is an upper bound on the number of bets n needed to reach the winning zone. The bound is proportional to the complexity of the environment μ .

Continuous Probability Classes \mathcal{M} . We have considered thus far countable probability classes \mathcal{M} , which makes sense from a computational point. On the other hand in statistical parameter estimation one often has a continuous hypothesis class (e.g. a Bernoulli(θ) process with unknown $\theta \in [0,1]$). Let $\mathcal{M} := \{\mu_{\theta} : \theta \in \Theta \subseteq \mathbb{R}^d\}$ be a family of probability distributions parameterized by a d -dimensional continuous parameter θ . Let $\mu \equiv \mu_{\theta_0} \in \mathcal{M}$ be the true generating distribution. For a continuous weight density $w(\theta) > 0$ the sums in 1.4 are naturally replaced by integrals: $\xi(x_{1:n}) := \int_{\Theta} w(\theta) \cdot \mu_{\theta}(x_{1:n}) d\theta$ with $\int_{\Theta} w(\theta) d\theta = 1$. The most important property of ξ in the discrete case was the dominance $\xi(x_{1:n}) \geq w_{\nu} \cdot \nu(x_{1:n})$ which has been obtained from (1.4) by dropping the sum over ν . The analogous construction here is to restrict the integral over Θ to a small vicinity N_{δ} of θ . For sufficiently smooth μ_{θ} and $w(\theta)$ we expect $\xi(x_{1:n}) \gtrsim |N_{\delta_n}| \cdot w(\theta) \cdot \mu_{\theta}(x_{1:n})$, where $|N_{\delta_n}|$ is the volume of N_{δ_n} . This in turn leads to $D_n \lesssim \ln w_{\mu}^{-1} + \ln |N_{\delta_n}|^{-1}$, where $w_{\mu} := w(\theta_0)$. N_{δ_n} should be the largest possible region in which $\ln \mu_{\theta}$ is approximately flat on average. More precisely, generalizing [CB90] to the non-i.i.d. case, we show $D_n \leq \ln w_{\mu}^{-1} + \frac{d}{2} \ln \frac{n}{2\pi} + O(1)$, where the $O(1)$ term depends on the smoothness of μ_{θ} , measured by the Fisher information. D_n is no longer bounded by a constant, but still grows only logarithmically with n , the intuitive reason being the necessity to describe θ to an accuracy $O(n^{-1/2})$. So bound (1.6) is also applicable to the case of continuously parameterized probability classes.

1.4 Rational Agents in known Probabilistic Environments

1.4.1 The Agent Model

A very general framework for intelligent systems is that of rational agents [RN95]. In cycle k , an agent performs *action* $y_k \in \mathcal{Y}$ (output) which results in a *perception* or *observation* $x_k \in \mathcal{X}$ (input), followed by cycle $k+1$ and so on. We assume that the action and perception spaces \mathcal{X} and \mathcal{Y} are finite. We write $p(x_{<k}) = y_{1:k}$ to denote the output $y_{1:k}$ of the agent's policy p on input $x_{<k}$ and similarly $q(y_{1:k}) = x_{1:k}$ for

the environment q in the case of deterministic environments. We call policy p and environment q behaving in this way *chronological*. The title page of the thesis depicts this interaction in case p and q are modeled by Turing machines. Note that policy and environment are allowed to depend on the complete history. We do not make any MDP or POMDP assumption here, and we do not talk about states of the environment, only about observations. In the more general case of a *probabilistic environment*, given the history $y_{<k}y_k \equiv y_1 \dots y_{k-1}y_k \equiv y_1x_1 \dots y_{k-1}x_{k-1}y_k$, the probability that the environment leads to perception x_k in cycle k is (by definition) $\mu(y_{<k}\underline{y}_k)$. The underlined argument \underline{y}_k in μ is a random variable and the other non-underlined arguments $y_{<k}y_k$ represent conditions.⁵ We call probability distributions like μ *chronological*. Since value optimizing policies (see below) can always be chosen deterministic, there is no real need to generalize the setting to probabilistic policies.

1.4.2 Value Functions and Optimal Policies

The goal of the agent is to maximize future *rewards*, which are provided by the environment through the inputs x_k . The inputs $x_k \equiv r_k x'_k$ are divided into a regular part x'_k and some (possibly empty or delayed) reward $r_k \in [0, r_{max}]$.⁶ We use the abbreviation

$$\mu(y_{<k}\underline{y}_{k:m}) = \mu(y_{<k}\underline{y}_k) \cdot \mu(y_{1:k}\underline{y}_{k+1}) \cdot \dots \cdot \mu(y_{<m}\underline{y}_m),$$

which is essentially the chain rule, and $\epsilon = y_{<1}$ for the empty string. We define the (total) *value* of policy p in environment μ , or shorter, the μ -value of p , as the μ -expected reward sum

$$V_\mu^p := \sum_{x_{1:m}} (r_1 + \dots + r_m) \mu(\underline{y}_{1:m})_{|y_{1:m}=p(x_{<m})}, \quad (1.7)$$

where m is the *lifespan* or initial *horizon* of the agent. The optimal policy p^μ which maximizes the value V_μ^p is

$$p^\mu := \arg \max_p V_\mu^p, \quad V_\mu^* := V_\mu^{p^\mu} = \max_p V_\mu^p \geq V_\mu^p \forall p$$

The policy p^μ , which we call *AI μ model*, is optimal in the sense that no other policy for an agent leads to higher μ -expected reward. Explicit expressions for the action y_k in cycle k of the μ -optimal policy p^μ and their value V_μ^* are

$$y_k = y_k^\mu := \arg \max_{y_k} \sum_{x_k} \max_{y_{k+1}} \sum_{x_{k+1}} \dots \max_{y_m} \sum_{x_m} (r_k + \dots + r_m) \cdot \mu(y_{<k}\underline{y}_{k:m}), \quad (1.8)$$

⁵The standard notation $\mu(x_k|y_{<k}y_k)$ for conditional probabilities destroys the chronological order and would become quite confusing in later expressions.

⁶In the reinforcement learning literature when dealing with (PO)MDPs the reward is usually considered to be a function of the environmental state. The zero-assumption analogue here is that the reward r_k is some probabilistic function μ' depending on the complete history. It is very convenient to integrate r_k into x_k and μ' into μ .

$$V_\mu^* = \max_{y_1} \sum_{x_1} \max_{y_2} \sum_{x_2} \dots \max_{y_m} \sum_{x_m} (r_1 + \dots + r_m) \cdot \mu(\underline{y}_{1:m}). \quad (1.9)$$

where $\underline{y}_{<k}$ is the actual history. We show that these definitions are consistent and correctly capture our intention. For instance, consider the expectimax expression (1.9): The best expected reward is obtained by averaging over possible perceptions x_i and by maximizing over the possible actions y_i . This has to be done in chronological order $y_1 x_1 \dots y_m x_m$ to correctly incorporate the dependency of x_i and y_i on the history. This is the origin of the alternating *expectimax* sequence, which is similar to the well-known minimax sequence/tree/algorithm in games theory.

1.4.3 Sequential Decision Theory & Reinforcement Learning

One can relate (1.9) to the Bellman equations [Bel57] of sequential decision theory by identifying complete histories $\underline{y}_{<k}$ with states, $\mu(\underline{y}_{<k} \underline{y}_k)$ with the state transition matrix, V_μ^* with the value function, and y_k with the action in cycle k [BT96, RN95]. Due to the use of complete histories as state space, the $\text{AI}\mu$ model neither assumes stationarity, nor the Markov property, nor complete accessibility of the environment. Every state occurs at most once in the lifetime of the system. For this and other reasons the explicit formulation (1.8) is much more useful here than to enforce a pseudo-recursive Bellman equation form.

As we have in mind a universal system with complex interactions, the action and perception spaces \mathcal{Y} and \mathcal{X} are huge (e.g. video images), and every action or perception itself occurs usually only once in the lifespan m of the agent. As there is no (obvious) universal similarity relation on the state space, an effective reduction of its size is impossible, but there is no principle problem in determining y_k from (1.8) as long as μ is known and computable and \mathcal{X} , \mathcal{Y} and m are finite.

Things dramatically change if μ is unknown. Reinforcement learning algorithms [KLM96, SB98, BT96] are commonly used in this case to learn the unknown μ . They succeed if the state space is either small or has effectively been made small by generalization or function approximation techniques. In any case, the solutions are either *ad hoc*, work in restricted domains only, have serious problems with state space exploration versus exploitation, are prone to diverge, or have non-optimal learning rate. There is no universal and optimal solution to this problem so far. The central theme of this thesis is to present a new model and argue that it formally solves all these problems in an optimal way. The true probability distribution μ will not be learned directly, but will be replaced by some universal prior ξ , which converges to μ , similarly to the induction (SP) case.

1.5 The Universal Algorithmic Agent AIXI

1.5.1 The Universal AIXI Model

We have developed enough formalism to suggest our universal AIXI model. All we have to do is to suitably generalize Solomonoff's universal prior M and to replace the true but unknown probability μ in the $\text{AI}\mu$ model by this generalized M . Similarly to (1.1), we define M as the $2^{-l(q)}$ weighted sum over all chronological programs (environments) q which output $x_{1:k}$, but with $y_{1:k}$ provided on the "input" tape. This also generalizes ξ_U (within an irrelevant multiplicative constant):

$$\xi(\underline{y}_{1:k}) = \xi_U(\underline{y}_{1:k}) \stackrel{\times}{=} M(\underline{y}_{1:k}) := \sum_{q: q(y_{1:k})=x_{1:k}} 2^{-l(q)}. \quad (1.10)$$

If not clear from context, we add superscripts SP and AI to ξ , to resolve ambiguities between (1.3) and (1.10). Replacing μ by ξ in (1.8) the *AIXI system* outputs

$$y_k = y_k^\xi := \arg \max_{y_k} \sum_{x_k} \dots \max_{y_m} \sum_{x_m} (r_k + \dots + r_m) \cdot \xi(\underline{y}_{<k} \underline{y}_{k:m}) \quad (1.11)$$

in cycle k given the history $\underline{y}_{<k}$. The ξ -value V_ξ^P and the universal value V_ξ^* are defined as in (1.7) and (1.9), with μ replaced by ξ . The AIXI model and its behavior is completely defined by (1.10) and (1.11). It (slightly) depends on the choice of the universal Turing machine, because $K()$ and $l()$ depend on U and hence are defined only up to terms of order one. The AIXI model also depends on the choice of \mathcal{X} and \mathcal{Y} , but we do not expect any bias when the spaces are chosen sufficiently large and simple, e.g. all strings of length 2^{16} . Choosing \mathcal{N} as the word space would be ideal, but whether the maxima (or suprema) exist in this case, has to be shown beforehand. The only non-trivial dependence is on the horizon m . Ideally we would like to chose $m = \infty$, but there are several subtleties to be unravelled later, which prevent at least a naive limit $m \rightarrow \infty$. So apart from m and unimportant details, *the AIXI system is uniquely defined by (1.10) and (1.11) without adjustable parameters.*

1.5.2 On the Optimality of AIXI

Universality and convergence of ξ . One can show that also ξ defined in (1.10) is universal and rapidly converges to μ analogous to the induction (SP) case. If we take a finite product of conditional ξ -s and use the chain rule, we see that also $\xi(\underline{y}_{<k} \underline{y}_{k:k+h})$ converges to $\mu(\underline{y}_{<k} \underline{y}_{k:k+h})$ for $k \rightarrow \infty$. This gives confidence that the outputs y_k^ξ of the AIXI model (1.11) could converge to the outputs y_k^μ of the $\text{AI}\mu$ model (1.8), at least for a bounded moving horizon h . The problems with a fixed horizon m and especially $m \rightarrow \infty$ will be discussed later.

Universally optimal AI systems. We want to call an AI model *universal*, if it is independent of the true environment μ (unbiased, model-free) and is able to solve any solvable problem and learn any learnable task. Further, we call a universal model, *universally optimal*, if there is no program, which can solve or learn significantly

faster (in terms of interaction cycles). As the AIXI model is parameter-free, ξ converges to μ , the $\text{AI}\mu$ model is itself optimal, and we expect no other model to converge faster to $\text{AI}\mu$ by analogy to the SP case,

we expect AIXI to be universally optimal.

This is our main claim. Further support is given below.

Intelligence order relation. We want to call a policy p *more or equally intelligent* than a policy p' and write $p \succeq p'$ if p yields in every cycle k and for every fixed history $y_{<k}$ higher (future) ξ -expected reward sum than p' . It is a formal exercise to show that $p^\xi \succeq p$ for all p . The AIXI model is, hence, the most intelligent agent w.r.t. \succeq . \succeq is a universal order relation in the sense that it is free of any parameters (except m) or specific assumptions about the environment. A proof that \succeq is a reasonable intelligence order (what we believe to be true), would prove that AIXI is universally optimal.

Value bounds. The values V_ρ^* associated with the $\text{AI}\rho$ systems correspond roughly to the negative total loss $-L_n^{\Lambda_\rho}$ (with $n=m$) of the $\text{SP}\rho$ ($=\Lambda_\rho$) systems. In the SP case we were interested in small bounds for the regret $L_n^{\Lambda_\xi} - L_n^{\Lambda_\mu}$. Unfortunately, simple value bounds for AIXI or any other AI system in terms of V_ν^* analogous to the loss bound (1.6) cannot hold. We even have difficulties in specifying what we can expect to hold for AIXI or any AI system which claims to be universally optimal. In SP, the only important property of μ for proving loss bounds was its complexity $K(\mu)$. In the AI case, there are no useful bounds in terms of $K(\mu)$ only. We either have to study restricted problem or environmental classes or consider bounds depending on other properties of μ , rather than on its complexity only.

1.5.3 Value Related Optimality Results

The mixture distribution ξ . In the following, we consider general Bayes-mixtures ξ over classes \mathcal{M} of chronological probability distributions ν :

$$\xi(y_{1:m}) = \sum_{\nu \in \mathcal{M}} w_\nu \nu(y_{1:m}) \quad \text{with} \quad \sum_{\nu \in \mathcal{M}} w_\nu = 1 \quad \text{and} \quad w_\nu > 0 \quad \forall \nu \in \mathcal{M}.$$

V_ξ^p , p^ξ , and V_ξ^* are defined as in with μ replaced by ξ . Policy p^ξ is called the $\text{AI}\xi$ model. For $\xi = \xi_U$ the $\text{AIXI} \equiv \text{AI}\xi_U$ model is recovered. If μ unknown, but is known to belong to the known class \mathcal{M} it is natural to follow policy p^ξ (which maximizes V_ξ^p). The (true μ -)expected reward when following policy p^ξ is $V_\mu^{p^\xi}$. The optimal (but infeasible) policy p^μ yields reward $V_\mu^{p^\mu} \equiv V_\mu^*$. It is now of interest (a) whether there are policies with uniformly larger value than $V_\mu^{p^\xi}$ and (b) how close $V_\mu^{p^\xi}$ is to V_μ^* .

Linearity and convexity of V_ρ in ρ . The following properties of V_ρ are crucial. V_ρ^p is a linear function in ρ and V_ρ^* is a convex function in ρ in the sense that

$$V_\xi^p = \sum_{\nu \in \mathcal{M}} w_\nu V_\nu^p \quad \text{and} \quad V_\xi^* \leq \sum_{\nu \in \mathcal{M}} w_\nu V_\nu^*.$$

Linearity is obvious from the definition of V_ρ^p and convexity follows easily from the convexity of \max_p and non-negativity of the weights w_ν . One loose interpretation of the convexity is that a mixture can never increase performance.

Pareto-optimality of $\text{AI}\xi$. Similarly to the SP case one can show that p^ξ is *Pareto-optimal* in the sense that there is no other policy p with $V_\nu^p \geq V_\nu^{p^\xi}$ for all $\nu \in \mathcal{M}$ and strict inequality for at least one ν . In particular, AIXI is Pareto-optimal.

Self-optimizing Policy p^ξ w.r.t. Average Value. Since we do not know the true environment μ in advance, we are interested under which circumstances⁷

$$\frac{1}{m} V_\nu^{p^\xi} \rightarrow \frac{1}{m} V_\nu^* \quad \text{for horizon } m \rightarrow \infty \quad \text{for all } \nu \in \mathcal{M}. \quad (1.12)$$

Note that V_ν as well as $p^\xi = p_m^\xi$ depend on m . The least we must demand from \mathcal{M} to have a chance that (1.12) is true is that there exists a policy (sequence) $\tilde{p} = \tilde{p}_m$ at all with this property, i.e.

$$\exists \tilde{p} : \frac{1}{m} V_\nu^{\tilde{p}} \rightarrow \frac{1}{m} V_\nu^* \quad \text{for horizon } m \rightarrow \infty \quad \text{for all } \nu \in \mathcal{M}. \quad (1.13)$$

We show that this necessary condition is also sufficient, i.e. (1.13) implies (1.12). This is another (asymptotic) optimality property of policy p^ξ . If universal convergence in the sense of (1.13) is possible at all in a class of environments \mathcal{M} , then policy p^ξ converges in the sense of (1.12). We call policies \tilde{p} with a property like (1.13) *self-optimizing* [KV86].

Unfortunately the result is not an asymptotic convergence statement of a single policy p^ξ , since p^ξ depends on m . The result merely says that under the stated conditions the average value of p_m^ξ is arbitrarily close to optimum for sufficiently large (pre-chosen) horizon m . This weakness will be resolved in the following.

Discounted Future Value Function. We now shift our focus from the total value to future values (value-to-go). First we have to get rid of the horizon parameter m . We eliminate the horizon by discounting the rewards $r_k \rightsquigarrow \gamma_k r_k$ with $\sum_{i=1}^{\infty} \gamma_i < \infty$ and taking $m \rightarrow \infty$. The analogue of m is now an effective horizon h_k^{eff} which may be defined by $\sum_{i=k}^{k+h_k^{\text{eff}}} \gamma_i \approx \sum_{i=k+h_k^{\text{eff}}}^{\infty} \gamma_i$. Furthermore, we renormalize the value V by $\sum_{i=k}^{\infty} \gamma_i$ and denote it by $V_{k\gamma}$. Finally, we extend the definition to probabilistic policies π (which is not essential). We define the γ -discounted weighted-average future *value*

⁷Here and elsewhere we interpret $a_m \rightarrow b_m$ as an abbreviation for $a_m - b_m \rightarrow 0$. $\lim_{m \rightarrow \infty} b_m$ may not exist.

of (probabilistic) policy π in environment ρ given history $y_{<k}$, or shorter, the ρ -value of π given $y_{<k}$, as

$$V_{k\gamma}^{\pi\rho}(y_{<k}) := \frac{1}{\Gamma_k} \lim_{m \rightarrow \infty} \sum_{y_{k:m}} (\gamma_k r_k + \dots + \gamma_m r_m) \rho(y_{<k} y_{k:m}) \pi(y_{<k} y_{k:m})$$

with $\Gamma_k := \sum_{i=k}^{\infty} \gamma_i$. The policy p^ρ is defined as to maximize the future value $V_{k\gamma}^{\pi\rho}$:

$$p^\rho := \arg \max_{\pi} V_{k\gamma}^{\pi\rho}, \quad V_{k\gamma}^{*\rho} := V_{k\gamma}^{p^\rho\rho} = \max_{\pi} V_{k\gamma}^{\pi\rho} \geq V_{k\gamma}^{\pi\rho} \forall \pi.$$

Setting $\gamma_k = 1$ for $k \leq m$ and $\gamma_k = 0$ for $k > m$ gives back the old undiscounted model with horizon m and $V_{1\gamma}^{p\rho} = \frac{1}{m} V_\rho^p$. Note that $V_{k\gamma}$ depends on the realized history $y_{<k}$. More important, p^ρ can be shown to be independent of k . Similarly to the undiscounted case one can prove that for every k and history $y_{<k}$, $V_{k\gamma}^{\pi\rho}$ is a linear function in ρ , $V_{k\gamma}^{*\rho}$ is a convex function in ρ , and p^ξ is Pareto-optimal in the sense that there is no other policy π with $V_{k\gamma}^{\pi\nu} \geq V_{k\gamma}^{p^\xi\nu}$ for all $\nu \in \mathcal{M}$ and strict inequality for at least one ν . Finally, p^ξ is self-optimizing (w.r.t. discounted value) if \mathcal{M} admits self-optimizing policies:

$$\text{If } \exists \tilde{\pi} \forall \nu : V_{k\gamma}^{\tilde{\pi}\nu} \xrightarrow{k \rightarrow \infty} V_{k\gamma}^{*\nu} \text{ w.}\nu.\text{p.1} \implies V_{k\gamma}^{p^\xi\nu} \xrightarrow{k \rightarrow \infty} V_{k\gamma}^{*\nu} \text{ w.}\mu.\text{p.1}.$$

The probability qualifier refers to the historic perceptions $x_{<k}$. The historic actions $y_{<k}$ are arbitrary. Note that k is a real running value, namely the current cycle number, whereas m was a pre-chosen fixed horizon.

1.5.4 Markov Decision Processes

From all possible environments, Markov (decision) processes are probably the most intensively studied ones. μ is called a (completely observable stationary) *Markov Decision Process* (MDP) if the probability of observing $x_k \in \mathcal{X}$, given history $y_{<k} y_k$ does only depend on the last action $y_k \in \mathcal{Y}$ and the last observation x_{k-1} , i.e. if $\mu(y_{<k} y_k \underline{x}_k) = \mu(x_{k-1} y_k \underline{x}_k)$. In this case x_k is called a *state*, \mathcal{X} the *state space*, and $\mu(x_{k-1} y_k \underline{x}_k)$ the *transition matrix*. An MDP μ is called *ergodic* if there exists a policy under which every state is visited infinitely often with probability 1. If an MDP $\mu(x_{k-1} y_k \underline{x}_k)$ is independent of the action y_k it is a *Markov process*, if it is independent of the last observation x_{k-1} it is an *i.i.d.* process.

Stationary MDPs μ with geometric discounting $\gamma_k = \gamma^k$ have stationary optimal policies p^μ mapping the same state/observation x_t always to the same action y_t . On the other hand a mixture ξ of MDPs is itself not an MDP, i.e. $\xi \notin \mathcal{M}_{MDP}$, which implies that p^ξ is, in general, not a stationary policy.

One can construct self-optimizing policies for the class of ergodic MDPs w.r.t. the average value $\frac{1}{m} V_\rho^p$ and if $\frac{\gamma_{k+1}}{\gamma_k} \rightarrow 1$ also w.r.t. to the discounted future value $V_{k\gamma}^{\pi\rho}$. The necessary condition $\frac{\gamma_{k+1}}{\gamma_k} \rightarrow 1$ ensures unboundedly increasing effective horizon h_k^{eff} . The existence of self-optimizing policies for ergodic MDPs implies that

for a countable class \mathcal{M} of ergodic MDPs, the policies p_m^ξ maximizing V_ξ^p and p^ξ maximizing $V_{k\gamma}^{\pi^\xi}$ are self-optimizing in the sense that

$$\forall \nu \in \mathcal{M} : \frac{1}{m} V_{1m}^{p_m^\xi \nu} \xrightarrow{m \rightarrow \infty} \frac{1}{m} V_{1m}^{* \nu} \quad \text{and} \quad V_{k\gamma}^{p^\xi \nu} \xrightarrow{k \rightarrow \infty} V_{k\gamma}^{* \nu} \quad \text{if} \quad \frac{\gamma_{k+1}}{\gamma_k} \rightarrow 1. \quad (1.14)$$

We also show that if \mathcal{M} is finite, then the speed of the first convergence is at least $O(m^{-1/3})$. The conditions $\Gamma_k < \infty$ and $\frac{\gamma_{k+1}}{\gamma_k} \rightarrow 1$ on the discount sequence are, for instance, satisfied for $\gamma_k = 1/k^2$, but *not* for the popular geometric discount $\gamma_k = \gamma^k$, which has finite effective horizon.

(1.14) shows that p^ξ is self-optimizing for Bandits, i.i.d. processes, and classification tasks, since they are special (degenerate) cases of ergodic MDPs. The existence of self-optimizing policies is not limited to (subclasses of ergodic) MDPs. Certain classes of POMDPs, k^{th} order ergodic MDPs, factorizable environments, repeated games, and prediction problems are not MDPs, but nevertheless admit self-optimizing policies, and hence the corresponding Bayes-optimal mixture policy p^ξ is self-optimizing.

1.5.5 The Choice of the Horizon

The only significant arbitrariness in the AIXI model lies in the choice of the lifespan m or in the discounted case in the discount sequence γ_k . We will not discuss *ad hoc* choices for specific problems. We are interested in universal choices. In many cases the time we are willing to run a system depends on the quality of its actions. Hence, the lifetime, if finite at all, is not known in advance. Geometric discounting $r_k \rightsquigarrow r_k \cdot \gamma^k$ solves the mathematical problem of $m \rightarrow \infty$ but is no real solution, since an effective horizon $h^{eff} \sim \ln \gamma^{-1} < \infty$ has been introduced. The scale invariant discounting $r_k \rightsquigarrow r_k \cdot k^{-\alpha}$ with $\alpha > 1$ has a dynamic horizon $h \sim k$. This choice has some appeal, as it seems that humans of age k years also usually do not plan their lives for more than the next $\sim k$ years. It also satisfies the condition $\frac{\gamma_{k+1}}{\gamma_k} \rightarrow 1$, necessary for AIXI being self-optimizing in ergodic MDPs. The largest lower semi-computable horizon with guaranteed finite reward sum $\Gamma_1 < \infty$ is obtained by the discount $r_k \rightsquigarrow r_k \cdot 2^{-K(k)}$, where $K(k)$ is the Kolmogorov complexity of k . This is maybe the most attractive universal discount. It is similar to a near-harmonic discount $r_k \rightsquigarrow r_k \cdot k^{-(1+\varepsilon)}$, since $2^{-K(k)} \leq 1/k$ for most k and $2^{-K(k)} \geq c/(k \log^2 k)$ for some constant c . We are not sure whether the choice of the horizon is of marginal importance, as long as it is chosen sufficiently large, or whether the choice will turn out to be a central topic for the AIXI model or for the planning aspect of any universal AI system in general. Most if not all problems in agent design of balancing exploration and exploitation vanish by a sufficiently large choice of the (effective) horizon and a sufficiently general prior.

1.6 Important Environmental Classes

In this and the next section we define $\xi = \xi_U \stackrel{\times}{=} M$ be Solomonoff's prior, i.e. $\text{AI}\xi = \text{AIXI}$. Each subsection represents an abstract on what will be done in the corresponding section of Chapter 6.

1.6.1 Introduction

In order to give further support for the universality and optimality of the $\text{AI}\xi$ theory, we apply $\text{AI}\xi$ in this second part to a number of problem classes. They include sequence prediction, strategic games, function minimization and, especially, how $\text{AI}\xi$ learns to learn supervised. For some classes we give concrete examples to illuminate the scope of the problem class. We first formulate each problem class in its natural way (when μ^{problem} is known) and then construct a formulation within the $\text{AI}\mu$ model and prove its equivalence. We then consider the consequences of replacing μ by ξ . The main goal is to understand why and how the problems are solved by $\text{AI}\xi$. We only highlight special aspects of each problem class. The goal is to give a better picture of the flexibility of the $\text{AI}\xi$ model.

1.6.2 Sequence Prediction (SP)

Using the $\text{AI}\mu$ model for sequence prediction (SP) is identical to Bayesian sequence prediction $\text{SP}\mu$. One might expect, when using the $\text{AI}\xi$ model for sequence prediction, one would recover exactly the universal sequence prediction scheme $\text{SP}\xi$, as $\text{AI}\xi$ was a unification of the $\text{AI}\mu$ model and the idea of universal probability ξ . Unfortunately this is not the case. One reason is that ξ is only a probability distribution in the inputs x and not in the outputs y . This is also one of the origins of the difficulty of proving loss/value bounds for $\text{AI}\xi$. Nevertheless, we argue that $\text{AI}\xi$ is equally well suited for sequence prediction as $\text{SP}\xi$ is. In a very limited setting we prove a (weak) error bound for $\text{AI}\xi$ which gives hope that a general proof is attainable.

1.6.3 Strategic Games (SG)

A very important class of problems are strategic games (SG). We restrict ourselves to deterministic strictly competitive strategic games like chess. If the environment is a minimax player, the $\text{AI}\mu$ model itself reduces to a minimax strategy. Repeated games of fixed lengths are a special case for factorizable μ . The consequences of variable game length is sketched. The $\text{AI}\xi$ model has to learn the rules of the game under consideration, as it has no prior information about these rules. We describe how $\text{AI}\xi$ actually learns these rules.

1.6.4 Function Minimization (FM)

There are many problems that fall into the category ‘resource bounded function minimization’ (FM). They include the Traveling Salesman Problem, minimizing production costs, inventing new materials or even producing, e.g. nice paintings, which are (subjectively) judged by a human. The task is to (approximately) minimize some function $f: \mathcal{Y} \rightarrow \mathcal{Z}$ within minimal number of function calls. We will see that a greedy model trying to minimize f in every cycle fails. Although the greedy model has nothing to do with downhill or gradient techniques (there is nothing like a gradient or direction for functions over \mathcal{Y}) which are known to fail, we discover the same difficulties. FM has already nearly the full complexity of general AI. The reason being that FM can actively influence the information gathering process by its trials y_k (whereas SP and CF=Classification cannot). We discuss in detail the optimal FM μ model and its inventiveness in choosing the $y \in \mathcal{Y}$. A discussion of the subtleties when using AI ξ for function minimization, follows.

1.6.5 Supervised Learning from Examples (EX)

Reinforcement learning, as the AI ξ model does, is an important learning technique but not the only one. To improve the speed of learning, supervised learning, i.e. learning by acquiring knowledge, or learning from a constructive teacher is necessary. We show, how AI ξ learns to learn supervised. It actually establishes supervised learning very quickly within $O(1)$ cycles.

1.6.6 Other Aspects of Intelligence

Finally, we give a brief survey of other general aspects, ideas and methods in AI, and their connection to the AI ξ model. Some aspects are directly included in the AI ξ model, others are or should be emergent.

1.7 Computational Aspects

Up to now we have shown the universal character of the AIXI model but have completely ignored computational aspects. We start by developing an algorithm M that is capable of solving any well-defined problem p as quickly as the fastest algorithm computing a solution to p , save for a factor of $1+\varepsilon$ and lower-order additive terms. Based on a similar idea we then construct a computable version of the AIXI model.

1.7.1 The Fastest & Shortest Algorithm for All Well-Defined Problems

Introduction. A wide class of problems can be phrased in the following way. Given a formal specification $f: \mathcal{X} \rightarrow \mathcal{Y}$ of a problem depending on some parameter $x \in \mathcal{X}$, we are interested in a fast algorithm computing solution $y \in \mathcal{Y}$.

Levin search is (within a (large) constant factor) the fastest algorithm to invert a function $g: \mathcal{Y} \rightarrow \mathcal{X}$, if g can be evaluated quickly. [Lev73b, Lev84]. Levin search can also handle time-limited optimization problems [Sol86]. Prime factorization, graph coloring, truth assignments, ... are Problems suitable for Levin search, if we want to find a solution, since verification is quick. Levin search cannot decide the corresponding decision problems. It is also not applicable to e.g. matrix multiplication, and reinforcement learning, since the verification task g is as hard as the computation task. Blum's Speed-up Theorem [Blu67, Blu71] shows that there are types of problems f for which an (incomputable) sequence of speed-improving algorithms (of increasing size) exists, but no fastest algorithm.

In the approach presented here, we consider only those algorithms which *provably* solve a given problem, and have a fast (i.e. quickly computable) time bound. Neither the programs themselves, nor the proofs need to be known in advance. Under these constraints we construct the asymptotically fastest algorithm save a factor of $1 + \varepsilon$ that solves any well-defined problem f .

The Fast Algorithm $M_{p^*}^\varepsilon$. Let p^* be a given algorithm computing $p^*(x)$ from x , or, more generally, a specification of a function f . One ingredient to our fastest algorithm $M_{p^*}^\varepsilon$ to compute $p^*(x)$ is an enumeration of proofs of increasing length in some formal axiomatic system. If a proof actually proves that p and p^* are functionally equivalent and p has time bound t_p , the tuple (p, t_p) is added to a list L . The program p in L with the currently smallest time bound $t_p(x)$ is executed. By construction, the result $p(x)$ is identical to $p^*(x)$. The trick to achieve a small running time is to schedule everything in a proper way, in order not to lose too much performance by computing slow p 's and t_p 's before *the* p has been found.

More formally, we say that a program “ p computes function f ”, when a universal reference Turing machine U on input (p, x) computes $f(x)$ for all x . This is denoted by $U(p, x) = f(x)$. To be able to talk about proofs, we need a formal logic system $(\forall, \lambda, y_i, c_i, f_i, R_i, \rightarrow, \wedge, =, \dots)$, and axioms, and inference rules. A proof is a sequence of formulas, where each formula is either an axiom or inferred from previous formulas in the sequence by applying the inference rules. We only need to know that *provability*, *Turing Machines*, and *computation time* can be formalized, and that the set of (correct) proofs is enumerable. We say that p is provably equivalent to p^* if the formula $[\forall y: U(p, y) = U(p^*, y)]$ can be proven. Let us fix $\varepsilon \in (0, \frac{1}{2})$. $M_{p^*}^\varepsilon$ runs three algorithms A , B , and C in parallel:

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

Initialize the shared variables
 $L := \{\}, \quad t_{fast} := \infty, \quad p_{fast} := p^*.$
 Start algorithms A , B , and C
 in parallel with ε , ε , $1-2\varepsilon$
 computational resources, respectively.

 B

Compute all $t(x)$ in parallel
 for all $(p, t) \in L$ with
 relative computation time $2^{-l(p)-l(t)}.$
if for some t , $t(x) < t_{fast}$,
then $t_{fast} := t(x)$ and $p_{fast} := p.$
continue

 A

Run through all proofs.
if a proof proves for some (p, t) that
 $p(\cdot)$ is equivalent to (computes) $p^*(\cdot)$
 and has time-bound $t(\cdot)$
then add (p, t) to $L.$

 C

run U on $(p_{fast}, x).$
 For each time-step decrease t_{fast} by 1.
if U halts **then** print result $U(p_{fast}, x)$
 and abort computation of A , B and $C.$

Note that A and B only terminate when aborted by C . It is obvious that $M_{p^*}^\varepsilon$ is equivalent to (computes) p^* . We show that the computation time of $M_{p^*}^\varepsilon$ is bounded by

$$time_{M_{p^*}^\varepsilon}(x) \leq (1 + \varepsilon) \cdot t_p(x) + \frac{d_p}{\varepsilon} \cdot time_{t_p}(x) + \frac{c_p}{\varepsilon},$$

$$d_p = 3 \cdot 2^{l(p)+l(t_p)}, c_p = 3 \cdot 2^{l(\text{proof}(p))+1} \cdot O(l(\text{proof}(p))^2),$$

where p is any algorithm, provably computing the same function as p^* with computation time provably bounded by the function $t_p(x)$ for all x . $time_{t_p}(x)$ is the time needed to compute the time bound $t_p(x)$. Known time bounds for practical problems can often be computed quickly, i.e. $time_{t_p}(x)/time_p(x)$ often converges very quickly to zero. Furthermore, from a practical point of view, the provability restrictions are often rather weak. Hence, we have constructed for all those problems a solution, which is asymptotically only a factor $1+\varepsilon$ slower than the (provably) fastest algorithm. On the flip side, for realistically sized problems, the lower order terms usually dominate, which limits the practical use of $M_{p^*}^\varepsilon$.

Algorithmic complexity and the shortest algorithm. A natural definition for the (Kolmogorov) complexity of a function f is the length of the shortest program computing f : $K'(f) := \min_p \{l(p) : U(p, x) = f(x) \forall x\}$. Unfortunately K' suffers from not even being approximable, since functional equality of programs is undecidable. Let p^* be a formal specification or a program for f . Using $K(p^*)$ is also not a suitable alternative, since it essentially depends on the choice of p^* , since, e.g. “dead code” in p^* contributes to $K(p^*)$. A satisfactory solution is to take the length of the shortest program *provably* equivalent to p^* :

$$K''(p^*) := \min_p \{l(p) : \text{a proof of } [\forall y : U(p, y) = U(p^*, y)] \text{ exists}\}$$

K'' (like K) is upper semi-computable. Let p' be some short description of p^* . We are now concerned with the computation time of p' . Could we get slower and slower algorithms by compressing p^* more and more? Interestingly this is not the case.

Inventing complex (long) programs is *not* necessary to construct asymptotically fast algorithms, under the stated provability assumptions, in contrast to Blum's Theorem [Blu67, Blu71]. We show that exists a program \tilde{p} , equivalent to p^* with

$$\begin{aligned} i) \quad l(\tilde{p}) &\leq K''(p^*) + O(1) \\ ii) \quad \text{time}_{\tilde{p}}(x) &\leq (1 + \varepsilon) \cdot t_p(x) + \frac{d_p}{\varepsilon} \cdot \text{time}_{t_p}(x) + \frac{c_p}{\varepsilon} \end{aligned}$$

where p is any program provably equivalent to p^* with computation time provably less than $t_p(x)$. That is, \tilde{p} is simultaneously among the shortest *and* fastest programs.

Generalizations. Algorithm $M_{p^*}^\varepsilon$ can be modified to handle I/O streams, definable by a Turing machine with monotone input and output tapes (and bidirectional work tapes) receiving an input stream and producing an output stream, as is the case in the agent setup.

1.7.2 Time Bounded AIXI Model

The major drawback of the $\text{AI}\xi$ model is that it is uncomputable. To overcome this problem, we construct a modified algorithm $\text{AI}\xi^{tl}$, which is still superior to any other time t and length l bounded agent. The computation time of $\text{AI}\xi^{tl}$ is of the order $t \cdot 2^l$.

Non-effectiveness of $\text{AI}\xi$. ξ^{AI} is not a computable but only an enumerable semimeasure. Hence, the output \dot{y}_k of the $\text{AI}\xi$ model is only asymptotically computable. $\text{AI}\xi$ yields an algorithm that produces a sequence of trial outputs eventually converging to the correct output \dot{y}_k , but one can never be sure whether one has already reached it. Besides this, convergence is extremely slow, so this type of asymptotic computability is of no direct practical use. Furthermore, the replacement of ξ^{AI} by time-limited versions [LV91, LV97], which is suitable for sequence prediction, has been shown to fail for the $\text{AI}\xi$ model. This leads to the issues addressed next.

Time bounds and effectiveness. Let \tilde{p} be a policy which calculates an acceptable output within a reasonable time \tilde{t} per interaction cycle. This sort of computability assumption, namely, that a general purpose computer of sufficient power and appropriate program is able to behave in an intelligent way, is the very basis of AI research. Here it is not necessary to discuss what exactly is meant by 'reasonable time/intelligence' and 'sufficient power'. What we are interested in is whether there is a computable version of the $\text{AI}\xi$ system which is superior or equal to any policy p with computation time per cycle of at most \tilde{t} .

What one can realistically hope to construct is an $\text{AI}\xi^{\tilde{t}\tilde{l}}$ system of computation time $c \cdot \tilde{t}$ per cycle for some constant c . The idea is to run all programs p of length $\leq \tilde{l} := l(\tilde{p})$ and time $\leq \tilde{t}$ per cycle and pick the best output in the sense of maximizing the *universal value* V_ξ^* . The total computation time is $c \cdot \tilde{t}$ with $c \approx 2^{\tilde{l}}$. Unfortunately V_ξ^* cannot be used directly since this measure is itself only semi-computable and the

approximation quality by using computable versions of ξ^{AI} given a time of order $c \cdot \tilde{t}$ is crude [LV97]. On the other hand, we *have* to use a measure which converges to V_ξ^* for $\tilde{t}, \tilde{l} \rightarrow \infty$, since we want the $AI\xi^{\tilde{t}\tilde{l}}$ model to converge to the $AI\xi$ model in that case.

Valid approximations. We suggest the following solution satisfying the above conditions: The main idea is to consider *extended chronological incremental policies* p , which in addition to the regular output y_k^p rate their own output with w_k^p . The $AI\xi^{\tilde{t}\tilde{l}}$ model selects the output $\dot{y}_k = y_k^p$ of the policy p with highest rating w_k^p . p might suggest any output y_k^p but it is not allowed to rate itself with an arbitrarily high w_k^p if one wants w_k^p to be a reliable criterion for selecting the best p . One must demand that no policy p is allowed to claim that it is better than it actually is. We define a logical predicate $VA(p)$, called *valid approximation*, which is true if, and only if, p *always* satisfies $w_k^p \leq V_\xi^p(y_{<k})$, i.e. never overrates itself. $V_\xi^p(y_{<k})$ is the ξ^{AI} expected future reward under policy p . Valid policies p can then be (partially) ordered w.r.t. their rating w_k^p .

The universal time bounded $AI\xi^{\tilde{t}\tilde{l}}$ system. In the following, we describe the algorithm p^* underlying the $AI\xi^{\tilde{t}\tilde{l}}$ system. It is essentially based on the selection of the best algorithms p_k^* out of the time \tilde{t} and length \tilde{l} bounded policies p , for which there exists a proof P of $VA(p)$ with length $\leq l_P$.

1. Create all binary strings of length l_P and interpret each as a coding of a mathematical proof in the same formal logic system in which $VA(\cdot)$ has been formulated. Take those strings which are proofs of $VA(p)$ for some p and keep the corresponding programs p .
2. Eliminate all p of length $> \tilde{l}$.
3. Modify all p in the following way: all output $w_k^p y_k^p$ of p is temporarily written on an auxiliary tape. If p stops in \tilde{t} steps the internal ‘output’ is copied to the output tape. If p does not stop after \tilde{t} steps a stop is forced and $w_k^p := 0$ and some arbitrary y_k^p is written on the output tape. Let \mathcal{P} be the set of all those modified programs.
4. Start first cycle: $k := 1$.
5. Run every $p \in \mathcal{P}$ on extended input $\dot{y}_{<k}$, where all outputs are redirected to some auxiliary tape: $p(\dot{y}_{<k}) = w_1^p y_1^p \dots w_k^p y_k^p$. This step is performed incrementally by adding \dot{y}_{k-1} for $k > 1$ to the input tape and continuing the computation of the previous cycle.
6. Select the program p with highest rating w_k^p : $p_k^* := \operatorname{argmax}_p w_k^p$.
7. Write $\dot{y}_k := y_k^{p_k^*}$ to the output tape.
8. Receive input \dot{x}_k from the environment.
9. Begin next cycle: $k := k + 1$, goto step 5.

Properties of the p^* algorithm. Let p be any extended chronological (incremental) policy of length $l(p) \leq \tilde{l}$ and computation time per cycle $t(p) \leq \tilde{t}$, for which there

exists a proof of $\text{VA}(p)$ of length $\leq l_P$. The algorithm p^* , depending on \tilde{l} , \tilde{t} and l_P but not on p , has always higher rating than any such p . The setup time of p^* is $t_{\text{setup}}(p^*) = O(l_P^2 \cdot 2^{l_P})$ and the computation time per cycle is $t_{\text{cycle}}(p^*) = O(2^{\tilde{l}} \cdot \tilde{t})$. Furthermore, for $\tilde{t}, \tilde{l}, l_P \rightarrow \infty$, policy p^* converges to the behavior of the AI ξ model.

Roughly speaking, this means that if there exists a computable solution to some AI problem at all, then the explicitly constructed algorithm p^* is such a solution. This claim is quite general, but there are some limitations and open questions, regarding the setup time, regarding the necessity that the policies must rate their own output, regarding true but not (efficiently) provable $\text{VA}(p)$, and regarding “inconsistent” policies.

1.8 Discussion

This section contains some discussion and remarks on otherwise unmentioned topics.

Value bounds. Rigorous proofs of value bounds for the AI ξ theory are the major theoretical challenge – general ones as well as tighter bounds for special environments μ^{AI} . Of special importance are suitable (and acceptable) conditions to μ^{AI} , under which \dot{y}_k and finite value bounds exist for infinite \mathcal{Y} , \mathcal{X} and m .

Scaling AI ξ down. One can downscale the AI ξ model by using more restricted forms of ξ^{AI} . This could be done in a similar way as the theory of universal induction has been downscaled with many insights to the Minimum Description Length principle [LV92a, Ris89] or to the domain of finite automata [FMG92]. The AI ξ model might similarly serve as a super model, from which specialized models could be derived.

Applications. We have shown how a number of AI problem classes, including *sequence prediction*, *strategic games*, *function minimization* and *supervised learning* fit into the general AI ξ model. All problems are claimed to be formally solved by the AI ξ model. The solution is, however, only formal, because the AI ξ model is uncomputable or, at best, approximable. First, each problem class is formulated in its natural way (when μ^{problem} is known) and then a formulation within the AI μ model is constructed and their equivalence is proven. Then, the consequences of replacing μ^{AI} by ξ^{AI} are considered. The main goal is to understand how the problems are solved by AI ξ .

Implementation and approximation. The AI $\xi^{\tilde{t}\tilde{l}}$ model suffers from the same large factor $2^{\tilde{l}}$ in computation time as Levin search for inversion problems [Lev73b, Lev84]. Nevertheless, Levin search has been implemented and successfully adapted and applied to a variety of problems [Sch95, WS96, Sch97, SZW97, Sch02b]. Hence, a direct implementation of the AI $\xi^{\tilde{t}\tilde{l}}$ model may also be successful, at least in toy environments, e.g. prisoner problems. The AI $\xi^{\tilde{t}\tilde{l}}$ algorithm should be regarded only as the first step toward a *computable universal AI model*. Elimination of the factor

$2^{\tilde{l}}$ without introducing a large additive constant like in M_p^ε and without giving up universality will probably be a very difficult task. One could try to select programs p and prove $\text{VA}(p)$ in a more clever way than by mere enumeration. All kinds of ideas like, heuristic search, genetic algorithms, advanced theorem provers, and many more could be incorporated. But now we have a problem.

Computability. We seem to have transferred the AI problem just to a different level. This shift has some advantages (and also some disadvantages) but presents, in no way, a solution. Nevertheless, we want to stress that we have reduced the AI problem to (mere) computational questions. Even the most general other systems the author is aware of, depend on some (more than complexity) assumptions about the environment, or it is far from clear whether they are, indeed, universally optimal. Although computational questions are themselves highly complicated, this reduction is a non-trivial result. A formal theory of something, even if not computable, is often a great step toward solving a problem and has also merits of its own (see previous paragraphs).

Elegance. Many researchers in AI believe that intelligence is something complicated and cannot be condensed into a few formulas. They believe it is more a combining of enough *methods* and much explicit *knowledge* in the right way. From a theoretical point of view, we disagree as the AI ξ model is simple and seems to serve all needs. From a practical point of view we agree to the following extent. To reduce the computational burden one should provide special purpose algorithms (*methods*) from the very beginning, probably many of them related to reduce the complexity of the input and output spaces \mathcal{X} and \mathcal{Y} by appropriate pre/post-processing methods.

Extra knowledge. There is no need to incorporate extra *knowledge* from the very beginning. It can be presented in the first few cycles in *any* format. As long as the algorithm that interprets the data is of size $O(1)$, the AI ξ system will “understand” the data after a few cycles. If the environment μ^{AI} is complicated but extra knowledge z makes $K(\mu^{AI}|z)$ small, one can show that a bound for ξ^{AI} similarly to (1.2) reduces roughly to $\frac{1}{2}\ln 2 \cdot K(\mu^{AI}|z)$ when $x_1 \equiv z$, i.e. when z is presented in the first cycle. Special purpose algorithms could also be presented in x_1 , but it would be cheating to say that no special purpose algorithms have been implemented in AI ξ . The boundary between implementation and training is blurred in the AI ξ model.

Training. We have not said much about the training process itself, as it is not specific to the AI ξ model and has been discussed in literature in various forms and disciplines. By a training process we mean a sequence simple-to-complex tasks to solve, with the simpler ones hopefully helping in learning the more complex ones. A serious discussion would be out of place. To repeat a truism, it is, of course, important to present enough knowledge x'_k and evaluate the system output y_k with r_k in a reasonable way. To maximize the information content in the reward, one should start with simple tasks and give positive reward to approximately the better half of the outputs y_k , for instance.

Asymptotically fast programs. Will the ultimate search for asymptotically fastest programs typically lead to fast or slow programs for arguments of practical size? Levin search, matrix multiplication and the algorithm $M_{p^*}^\epsilon$ seem to support the latter, but this might be due to our inability to do better.

The big questions. A discussion of the “big” questions concerning the mere existence of any computable, fast, and elegant universal theory of intelligence, related to Penrose’s non-computable environments [Pen94], and Chaitin’s ‘number of wisdom’ Ω [Cha75, Cha91] will be given later.

1.9 History & References

Introductory textbooks. The book of Hopcroft and Ullman, and in the new revision, co-authored by Motwani [HMU01] is a very readable elementary introduction to automata theory, formal languages, and computation theory. The Artificial Intelligence book [RN95] by Russell and Norvig gives a comprehensive overview over AI approaches in general. For an excellent introduction to Algorithmic Information Theory, Kolmogorov complexity, and Solomonoff induction one should consult the book of Li and Vitányi [LV97]. The Reinforcement Learning book by Sutton and Barto [SB98] requires no background knowledge, describes the key ideas, open problems, and great applications of this field. A tougher and more rigorous book by Bertsekas and Tsitsiklis on sequential decision theory provides all (convergence) proofs [Ber95a, Ber95b].

Algorithmic information theory. Kolmogorov [Kol65] suggested to define the information content of an object as the length of the shortest program computing a representation of it. Solomonoff [Sol64] invented the closely related universal prior probability distribution and used it for binary sequence prediction [Sol64, Sol78] and function inversion and minimization [Sol86]. Together with Chaitin [Cha66, Cha75] this was the invention of what is now called Algorithmic Information theory. For further literature and many applications see [LV97]. Other interesting “applications” can be found in [Cha91, Sch99, VW98]. Related topics are the Weighted Majority Algorithm invented by Littlestone and Warmuth [LW94], universal forecasting by Vovk [Vov92], Levin search [Lev73b], pac-learning introduced by Valiant [Val84] and Minimum Description Length [LV92a, Ris89]. Resource bounded complexity is discussed in [Dal73, Dal77, FMG92, Ko86, PF97], resource bounded universal probability in [LV91, LV97, Sch02c]. Implementations are rare and mainly due to Schmidhuber [Sch95, WS96, Sch97, SZW97, Sch02b, Con97]. Good reviews with a philosophical touch are [LV92b, Sol97]. For an older, but general review of inductive inference see Angluin [AS83].

Sequential Decision Theory. The other ingredient in our AIXI model is sequential decision theory. We do not need much more than the maximum expected utility principle and the expectimax algorithm [Mic66, RN95]. The book of von Neumann

and Morgenstern [NM44] might be seen as the initiation of game theory, which already contains the expectimax algorithm as a special case. If the true environmental μ is unknown, it needs to be learned with, e.g. the help of reinforcement learning algorithms. Existing reinforcement learning algorithms are [Sam59, BSA83, Sut88, Wat89, WD92, MA93, Tes94, BT96, KLM96, KLC98, WS98, KS98], but they are rather limited in view of AIXI. The literature on reinforcement learning and sequential decision theory is so vast that we refer to the textbooks [SB98, BT96, KV86] for further references.

The author's contributions. Most of the issues addressed in this thesis can already be found scattered in various reports and publications by the author: The AIXI model has first been introduced and discussed in March 2000 in [Hut00] in a 62 page long report. More succinct descriptions have been published in [Hut01d, Hut01e]. The AIXI model has been argued to formally solve a number of problem classes, including sequence prediction, strategic games, function minimization, reinforcement and supervised learning [Hut00]. The generalization $\text{AI}\xi$ has recently been shown to be self-optimizing and Pareto-optimal [Hut02c]. The construction of a general fastest algorithm (within a factor 5) for all well-defined problems [Hut02a] arose from the construction of the time-bounded $\text{AIXI}t_l$ model [Hut00, Hut01d]. Tight [Hut02b] error [Hut01c, Hut01a] and loss [Hut01b] bounds for Solomonoff's universal sequence prediction scheme have been proven. Loosely related ideas on a market/economy based reinforcement learner [KHS01b] and gradient based reinforcement planner [KHS01a] have been implemented. These and other papers are available at <http://www.idsia.ch/~marcus/ai>.



William of Ockham
(1285–1349)

“Nulla pluralitas est ponenda nisi per rationem vel experiantiam vel auctoritatem illius, qui non potest falli nec errare, potest convivi.”

“A plurality should only be postulated if there is some good reason, experience, or unfallible authority for it.” (William of Ockham 1285 - 1349)

Chapter 2

Simplicity & Uncertainty

2.1	Introduction	202
2.1.1	Ockham, Epicurus, Hume, Bayes, Solomonoff	202
2.1.2	Problem Setup	203
2.2	Algorithmic Information Theory	204
2.2.1	Definitions and Notation	204
2.2.2	Turing Machines	205
2.2.3	Kolmogorov Complexity	207
2.2.4	Computability Concepts	209
2.3	Uncertainty & Probabilities	211
2.3.1	Frequency Interpretation / Counting	212
2.3.2	Objective Interpretation: Probabilities for Uncertain Events	212
2.3.3	Subjective Interpretation: Probabilities for Degrees of Belief	214
2.3.4	Determining Priors	215
2.4	Algorithmic Probability & Universal Induction	215
2.4.1	The Universal Prior M	216
2.4.2	Universal Sequence Prediction	217
2.4.3	Universal (Semi)Measures	218
2.4.4	Martin-Löf Randomness	224
2.5	History & References	225
2.6	Problems	229

This Chapter deals with the question of how to make predictions in unknown environments. Following a brief description of important philosophical attitudes regarding inductive reasoning and inference, we describe more accurately what we mean by induction, and motivate why we can focus on sequence prediction tasks. The most important concept is Occam’s razor (simplicity) principle. Indeed, one can

show that the best way to make predictions is based on the shortest ($\hat{=}$ simplest) description of the data sequence seen so far. The most general effective descriptions can be obtained with the help of general recursive functions, or equivalently by using programs on Turing machines, especially on the universal Turing machine. The length of the shortest program describing the data is called the Kolmogorov complexity of the data. Unfortunately, the Kolmogorov complexity is not finitely computable, which makes it necessary to introduce several weaker computability concepts. Probability theory is needed to deal with uncertainty. The environment may be a stochastic process (e.g. gambling houses or quantum physics) which can be described by “objective” probabilities. But also uncertain knowledge about the environment, which leads to beliefs about it, can be modeled by “subjective” probabilities. The old question left open by subjectivists of how to choose the a priori probabilities is solved by Solomonoff’s universal prior, which is closely related to Kolmogorov complexity. Solomonoff’s major result is that the universal (subjective) prior converges to the true (objective) environment(al probability) μ . The only assumption on μ is that μ (which needs not be known!) is computable. The problem of the unknown environment μ is hence solved for all problems of inductive type, like sequence prediction and classification. Finally, we show the (non)existence of universal priors for the other introduced computability concepts.

What is new? This chapter is mainly a collection of known concepts which are needed in later chapters, some are presented possibly in a new light. We omit all known proofs and are short in discussion, since many issues reappear in later chapters in more general or related contexts and with proofs provided. The major new result is the classification of generalized-computable universal priors in Theorem 2.26. For a slower and more detailed introduction into Kolmogorov complexity and Solomonoff induction and most proofs one should consult the excellent book of Li and Vitányi [LV97].

2.1 Introduction

2.1.1 Ockham, Epicurus, Hume, Bayes, Solomonoff

One very important and highly non-trivial aspect of intelligence is inductive inference. Simply speaking, induction is the process of predicting the future from the past or, more precisely, it is the process of finding rules in (past) data and using these rules to guess future data. Weather prediction, stock-market forecasting, or continuing number series in an IQ test, are non-trivial examples. Making good predictions plays a central role in natural and artificial intelligence in general, and in machine learning in particular.

On the one hand, induction seems to happen in every day life by finding regularities in past observations and using them to predict the future. On the other hand, this procedure seems to add knowledge about the future from past observa-

tions. But how can we know something about the future? This dilemma and the induction principle in general have a long philosophical history. There are

- *Epicurus' principle of multiple explanations* (342?-270? B.C.)
If more than one theory is consistent with the observations, keep all theories.
- *Occam's razor*¹ (*simplicity principle*) (1290?-1349?)
Entities should not be multiplied beyond necessity – or – keep the simplest theory consistent with the observations.
- *Hume's negation of Induction* (1711-1776) [Hum39]
The belief in the possibility of true induction cannot be justified rationally.
- *Bayes rule for conditional probabilities* (1702-1761) [Bay63]
It tells us how to update our beliefs/probabilities when acquiring new data.

Solomonoff [Sol64] cleverly unified the principles of Epicurus, Occam, and Bayes into one formal universal theory of inductive inference. Among all possible induction schemes it is the optimal method for making predictions.

2.1.2 Problem Setup

Every induction problem can be phrased as a sequence prediction task. This is most clearly illustrated in the domain of time series prediction. Having observed data x_t at times $t < n$, the task is to predict the n^{th} symbol x_n from sequence $x_1 \dots x_{n-1}$. Classification can also be seen as a sequence prediction task. The task of classifying a new instance z_n after having seen (instance,class) pairs $(z_1, c_1), \dots, (z_{n-1}, c_{n-1})$ can be phrased as to predict the continuation of the sequence $z_1 c_1 \dots z_{n-1} c_{n-1} z_n$.² Machine learning is often concerned with finding the *true* or a *predictive* or a *causal model* based on observed data. This step is important for *understanding* the domain under consideration. Understanding is often a goal in itself, but finally the goal is to apply the model to make predictions. In this view, model learning is only an intermediate step. The direct study of predictions based on past observations without discussing models has been coined *prequential approach* by Dawid [Daw84] for sequence predictions and *transductive inference* by Vapnik [Vap99, sec.9.1] for classification and regression. Several difficult issues are avoided by abandoning models. This includes questions about model consistency, i.e. whether the true model can be learned, and how to separate noise from useful data [GTV01, VV02]. One may even go one step further and ask why we want to make predictions. Usually the goal of prediction is to maximize one's profit/value or equivalently to minimize one's loss. In considering only profits or losses one avoids questions on whether prediction algorithms

¹Whereas *William of Ockham* is spelled with *ckh*, for some reason *Occam's razor* is usually spelled with *cc*.

²Sequence prediction may also be phrased as a classification task by adding time tags and if one does not assume a random generation of instances. To predict the next symbol of sequence $x_1 x_2 \dots x_{n-1}$ is the same as to try to find the class label of n after having seen (instance,class) pairs $(1, x_1), \dots, (n-1, x_{n-1})$.

converge to the best possible prediction algorithm (i.e. whether they are *self-tuning* [KV86, p232,p272]). Algorithms for which the loss converges to the minimal possible loss are called *self-optimizing* [KV86, p234]. This is a weaker demand than self-tuningness, but is often all we really care about. Our main purpose in this work is to study algorithms which minimize loss. Convergence of posterior probability distributions or algorithms themselves or models are only considered if this is useful for the ultimate goal of minimizing loss. To summarize our setup:

- Every induction problem can be phrased as a sequence prediction task.
- Classification is a special case of sequence prediction.
(With some tricks the other direction is also true)
- We are interested in maximizing profit or minimizing loss.
We are not primarily interested in finding (true/predictive/causal) models or even in convergence of the predictor itself.
- Separating noise from data is *not* necessary in this setting.

After having clarified the setup we now must delve into math before we can present Solomonoff's induction scheme.

2.2 Algorithmic Information Theory

In this section we give a very brief introduction to Kolmogorov complexity. For a slower, more thorough and comprehensive introduction see [LV97].

2.2.1 Definitions and Notation

We write $\mathbb{N} = \{1, 2, 3, \dots\}$ for the set of natural numbers, \mathbb{B}^* for the set of finite binary strings, and \mathbb{B}^∞ for the set of infinite binary sequences. We use letters i, k, n for natural numbers, x, y, z for finite strings, ϵ for the empty string, 1^n the string of n ones, $l(x)$ for the length of string x , and ω for infinite strings. We write xy for the concatenation of string x with y .

Every countable set may be identified with \mathbb{N} (by means of a bijection). We can interpret a string as a binary representation of a natural number. Unfortunately a naive identification will not be unique, since, for instance string 00101 and 101 represent both the number 5. We get a bijection if we map x to the natural number which has binary representation $1x$ (x prefixed with 1). We subtract 1 from this number, since we need a bijection between \mathbb{B}^* and $\mathbb{N}_0 := \{0, 1, 2, 3, \dots\}$ (see Table 2.2). With this identification $\log_2(x+1) - 1 < l(x) \leq \log_2(x+1)$. String x is called a (proper) prefix of y if there is a $z (\neq \epsilon)$ such that $xz = y$. A set of strings is called prefix-free if no element is a proper prefix of another. A prefix-free set \mathcal{P} is also

Table 2.2 ((Prefix) coding natural numbers and strings) Bijection between natural numbers \mathbb{N} and strings \mathbb{B}^* . Further, the length $l(x)$, and first and second order prefix coding $\bar{x} := 1^{l(x)}0x$ and $x' := \overline{l(x)}x$. For illustrational purpose we separated the first part $\overline{l(x)}$ from the second part x by an additional space. x' is longer than \bar{x} only for $x < 15$, but shorter for all $x > 30$.

$x \in \mathbb{N}_0$	0	1	2	3	4	5	6	7	...
$x \in \mathbb{B}^*$	ϵ	0	1	00	01	10	11	000	...
$l(x)$	0	1	1	2	2	2	2	3	...
\bar{x}	0	100	101	11000	11001	11010	11011	1110000	...
x'	0	100 0	100 1	101 00	101 01	101 10	101 11	11000 000	...

called a prefix-code. Prefix-codes have the important property of satisfying Kraft's inequality

$$\sum_{x \in \mathcal{P}} 2^{-l(x)} \leq 1 \quad (2.1)$$

For $\bar{x} := 1^{l(x)}0x$ the set $\{\bar{x} : x \in \mathbb{B}^*\}$ forms a prefix code with $l(\bar{x}) = 2l(x) + 1$. For $x' := \overline{l(x)}x = 1^{l(l(x))}0l(x)x$ the set $\{x' : x \in \mathbb{B}^*\}$ forms an asymptotically shorter prefix code with $l(x') = l(x) + 2l(l(x)) + 1$ (see Table 2.2). We pair strings x and y (and z) by $\langle x, y \rangle := x'y$ (and $\langle x, y, z \rangle := x'y'z$) which are uniquely decodable, since x' and y' are prefix. Since ' serves as a separator we also write $f(x, y)$ instead of $f(x'y)$ for functions f .

We abbreviate $\lim_{n \rightarrow \infty} [f(n) - g(n)] = 0$ by $f(n) \xrightarrow{n \rightarrow \infty} g(n)$ and say f converges to g , without implying that $\lim_{n \rightarrow \infty} g(n)$ itself exists. We write $f(n) \sim g(n)$ and say that f goes asymptotically to g if $\lim_{n \rightarrow \infty} f(n)/g(n) = 1$. We write $a \lesssim b$ if a is not much larger than b , with precision left unspecified. The big O -notation $f(x) = O(g(x))$ means that there are constants c and $x_0 > 0$ such that $|f(x)| \leq c|g(x)| \forall x > x_0$. The small o -notation $f(x) = o(g(x))$ abbreviates $\lim_{x \rightarrow \infty} f(x)/g(x) = 0$. We write $f(x) \overset{\times}{\leq} g(x)$ for $f(x) = O(g(x))$, $f(x) \overset{+}{\leq} g(x)$ for $f(x) \leq g(x) + O(1)$, and $f(x) \overset{\log}{\leq} g(x)$ for $f(x) \leq g(x) + O(\log g(x))$. Corresponding equalities can be defined similarly. They hold if the corresponding inequalities hold in both directions.

2.2.2 Turing Machines

A Turing machine can be considered as an idealized form of a computer. It consists of tapes (memory), read/write heads, table of rules (program), and an internal state (instruction pointer). A formal definition can be found in any textbook on computability theory, e.g. [HMU01]. The set of partial recursive functions coincides with the set of functions computable with a Turing machine. We say that a set of objects $S = \{o_1, o_2, o_3, \dots\}$ can be (effectively) enumerated if there is a Turing machine mapping i to $\langle o_i \rangle$, where $\langle \rangle$ is some default coding of the elements in S .

The importance of partial recursive functions and Turing machines stems from the following theses:

Thesis 2.3 (Turing) Everything that can be reasonably said to be computable by a human using a fixed procedure can also be computed by a Turing machine.

Thesis 2.4 (Church) The class of algorithmically computable numerical functions (in the intuitive sense) coincides with the class of partial recursive functions.

We need to supplement Turing's and Church's theses in the following way:

Assumption 2.5 (Short compiler) Given two *natural* Turing-equivalent formal systems $F1$ and $F2$, then there always exists a single *short* program on $F2$ which is capable of interpreting all $F1$ -programs.

This means that the difference of the size of the shortest $F1$ -description and the shortest $F2$ -description (of something) is not only bounded by a universal constant, but that this constant is also *reasonably small* for *natural* formal systems. It is easy to formally convert the interpreter into a compiler, by attaching the interpreter to the program to be interpreted and by “selling” the result as a compiled version.

This extends Church's and Turing's Theses in two respects. First it says that the equivalence is effective, i.e. that there exists *one* program (interpreter/compiler) which effectively converts $F1$ -programs to $F2$ -programs. Church's & Turing's thesis only state that the classes of computable functions coincide, leaving open the possibility that there is no effective way of transformation. Second, and more important, our extended Thesis states that the compiler is short if both formal systems are natural.

The above theses are not provable or falsifiable theorems, since *human*, *reasonable*, *intuitive*, and *natural* have not been defined rigorously. One may *define intuitively computable* as Turing computable and a *natural Turing-equivalent system* as one which has a small (say $< 10^5$ bits) interpreter/compiler on a once and for all agreed on fixed reference universal Turing machine. The theses would then be that these definitions are reasonable.

For technical reasons we need the following variant of a Turing machine.

Definition 2.6 (Prefix/Monotone Turing machine) A prefix/monotone Turing machine is defined as a Turing machine with one unidirectional input tape, one unidirectional output tape, and some bidirectional work tapes. Input tapes are read only, output tapes are write only, unidirectional tapes are those where the head can only move from left to right. All tapes are binary (no blank symbol), work tapes initially filled with zeros.

Prefix TM. We say T halts on input p with output x and write $T(p) = x$ if p is to the left of the input head and x is to the left of the output head after T halts. The set of p on which T halts forms a prefix-code. We call such codes p *self-delimiting* programs.

Monotone TM. We say T outputs/computes a string starting with x (or a sequence ω) on input p and write $T(p) = x*$ (or $T(p) = \omega$) if p is to the left of the input head when the last bit of x is output. T may continue operation and need not to halt. For given x , the set of such p forms a prefix-code. We call such codes p *minimal* programs.

The table of rules of a Turing machine T can be encoded in a canonical way as a binary string, which we denote by $\langle T \rangle$. Hence, the set of Turing-machines $\{T_1, T_2, \dots\}$ can be effectively enumerated. There are so-called universal Turing machines which can “simulate” all other Turing-machines. We define a particular one below, which also allows for side information y .

Theorem 2.7 (Universal prefix/monotone Turing machine) There exists a universal prefix/monotone Turing machine U which simulates prefix/monotone Turing machine T_i with input $y'i'q$ if fed with input $y'i'q$, i.e.

$$U(y'i'q) = T_i(y'q) \forall i, q$$

We call this particular U the *reference* universal Turing machine. Note that for p not of the form $y'i'q$, $U(p)$ does not halt. The price we have to pay for the existence of a universal Turing machine is the undecidability of the halting problem [Tur36]. In case of no side information $y = \epsilon$ we suppress in the following the initial $y' = \epsilon' = 0$ in the codes. We also drop the adjunct ‘prefix/monotone’ if clear from the context and identify objects with their coding $\langle \rangle$, i.e. we omit the $\langle \rangle$.

2.2.3 Kolmogorov Complexity

In order to exploit Occam’s razor beyond intuition we need to formalize the concept of simplicity and/or complexity. We first discuss the case of zero background knowledge $y = \epsilon$. Intuitively a string is simple if it can be described in a few words, like “the string of one million ones”, and is complex if there is no such short description, like for a random string whose shortest description is specifying it bit-by-bit. We are only interested in descriptions or *codes* which are effective and hence restrict the

decoders to Turing machines. We say that (program) p is a description of string x relative to the prefix Turing machine T if $T(p) = x$. The length of the shortest description is denoted by $K_T(x) := \min_p \{l(p) : T(p) = x\}$. This complexity measure depends on T and one may ask whether there exists a Turing machine which leads to shortest codes among *all* Turing machines for *all* x . Remarkably, there exists a Turing machine (the universal one) which “nearly” has this property. If p is the shortest description of x under $T = T_i$, then $i'p$ is a description of x under U , hence

$$K_U(x) \leq K_T(x) + c_{TU} \quad (2.8)$$

with $c_{TU} = l(i')$, and similarly for other choices of universal Turing machines. The length of the shortest description of x under U is at most a constant number of bits longer than the shortest description under T . The statement and proof of this invariance theorem in [Sol64, Kol65, Cha69] is often regarded as the birth of algorithmic information theory. Furthermore, for each pair of universal Turing machines U' and U'' satisfying the invariance theorem the complexities coincide up to an additive constant ($|K_{U'}(x) - K_{U''}(x)| \leq c_{U'U''}$).

Since $c_{U'U''}$ is essentially a compiler/interpreter constant we recall Assumption 2.5 and interpret the assumption as $c_{U'U''}$ being small for natural universal Turing machines U' and U'' . Henceforth we write $O(1)$ for terms like $c_{U'U''}$ which only depend on the choice of universal Turing machines, but which are independent of the strings under consideration. We extend the definition of complexity to allow side information y .

Definition 2.9 (Kolmogorov complexity) Let U be the reference universal prefix Turing machine U of Theorem 2.7. The (conditional) prefix Kolmogorov complexity is defined as the shortest program p , for which U outputs x (given y):

$$K(x) := \min_p \{l(p) : U(p) = x\}, \quad K(x|y) := \min_p \{l(p) : U(y'p) = x\}$$

For general (non-string) objects (like computable functions) one can specify some default coding and define $K(\text{object}) := K(\langle \text{object} \rangle)$, especially for numbers and pairs, e.g. we abbreviate $K(x, y) := K(\langle x, y \rangle) = K(x'y)$. The most important information theoretic properties of K are listed below.

Theorem 2.10 (Information properties of Kolmogorov complexity)

- i) $K(x) \stackrel{+}{\leq} l(x) + 2\log_2 l(x), \quad K(n) \stackrel{+}{\leq} \log_2 n + 2\log_2 \log n$
- ii) $\sum_x 2^{-K(x)} \leq 1, \quad K(x) \geq l(x)$ for “most” x , $K(n) \rightarrow \infty$ for $n \rightarrow \infty$.
- iii) $K(x|y) \stackrel{+}{\leq} K(x) \stackrel{+}{\leq} K(x, y)$
- iv) $K(xy) \stackrel{+}{\leq} K(x, y) \stackrel{+}{\leq} K(x) + K(y|x) \stackrel{+}{\leq} K(x) + K(y)$
- v) $K(x|y, K(y)) + K(y) \stackrel{\pm}{=} K(x, y) \stackrel{\pm}{=} K(y, x) \stackrel{\pm}{=} K(y|x, K(x)) + K(x)$
- vi) $K(f(x)) \stackrel{+}{\leq} K(x) + K(f)$ for recursive $f: \mathcal{B}^* \rightarrow \mathcal{B}^*$
- vii) $K(x) \stackrel{+}{\leq} -\log_2 P(x) + K(P)$ if $P: \mathcal{B}^* \rightarrow [0, 1]$ is enumerable and $\sum_x P(x) \leq 1$

All (in)equalities remain valid if K is (further) conditioned under some z , i.e. $K(\dots) \rightsquigarrow K(\dots|z)$ and $K(\dots|y) \rightsquigarrow K(\dots|y, z)$. Those stated are all valid within an additive constant of size $O(1)$, but there are others which are only valid to logarithmic accuracy $\stackrel{\log}{\leq}$. K has many properties in common with Shannon entropy as it should be, since both measure the information content of a string. (i) gives an upper bound on K . (ii) is Kraft’s inequality which implies a lower bound on K valid for “most” n , where “most” means that there are only $o(N)$ exceptions for $n \in \{1, \dots, N\}$ (Figure 2.11). Providing side information y can never increase code-length, requiring extra information y can never decrease code-length (iii). Coding x and y separately never helps (iv), and transforming x does not increase its information content (vi). (vi) also shows that switching from one coding scheme to another by means of a recursive bijection leaves K unchanged within additive $O(1)$ terms. The first non-trivial result is the symmetry of information (v), which is the analogue of the chain rule (see below). (vii) is at the heart of the MDL principle [Ris89], which approximates $K(x)$ by $-\log_2 P(x) + K(P)$.

All upper bounds on $K(z)$ are easily proven by devising *some* (effective) code for z of the length of the right hand side of the inequality and by noting that $K(z)$ is the length of the shortest code among all possible effective codes. For instance if T_{i_0} with $i_0 = O(1)$ is a Turing machine with $T_{i_0}(\epsilon' i_0' x') = x$, then $U(\epsilon' i_0' x') = x$, hence $K(x) \leq l(\epsilon' i_0' x') \stackrel{+}{\leq} l(x') \stackrel{+}{\leq} l(x) + 2\log_2 l(x)$, which proves (i). In (vii) one uses the Shannon-Fano code based on probability distribution P . Lower bounds are usually proven by counting arguments (Easy for (ii) by using (2.1) and harder for (vi)).

2.2.4 Computability Concepts

We need several computability concepts weaker than can be captured by halting Turing machines.

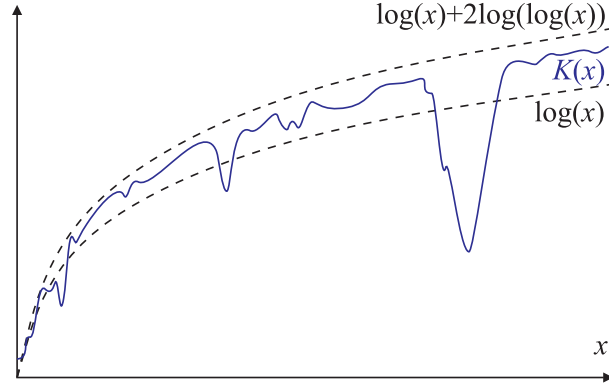


Figure 2.11 (Kolmogorov Complexity) Schematic graph of prefix Kolmogorov complexity $K(x)$ with string x interpreted as integer. $K(x) \geq x$ for “most” x and $K(x) \leq \log_2 x + 2\log_2 \log x + c$ for all x for sufficiently large constant c .

Definition 2.12 (Computable functions) We consider functions $f: \mathbb{N} \rightarrow \mathbb{R}$:

f is *finitely computable* or *recursive* iff there is a Turing machine T with $T(x^d) = n^d$ and $\frac{n}{d} = f(x)$,

f is *approximable* iff there is a Turing machine finitely computing $\phi(\cdot, \cdot)$ such that $\lim_{t \rightarrow \infty} \phi(x, t) = f(x)$.

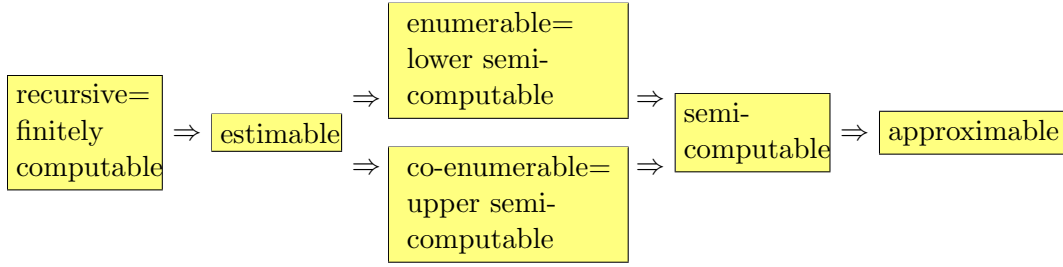
f is *lower semi-computable* or *enumerable* iff additionally $\phi(x, t) \leq \phi(x, t+1)$.

f is *upper semi-computable* or *co-enumerable* iff $[-f]$ is lower semi-computable.

f is *semi-computable* iff f is lower- or upper semi-computable.

f is *estimable* iff f is lower- and upper semi-computable.

If f is estimable we can finitely compute an ε -approximation of f by upper and lower semi-computing f and terminating when differing by less than ε . This means that there is a Turing machine which, given x and ε , finitely computes \hat{y} such that $|\hat{y} - f(x)| < \varepsilon$. Moreover it gives an interval estimate $f(x) \in [\hat{y} - \varepsilon, \hat{y} + \varepsilon]$. An estimable integer-valued function is finitely computable (take any $\varepsilon < 1$). Note that if f is only approximable or semi-computable we can still come arbitrarily close to $f(x)$ but we cannot devise a terminating algorithm which produces an ε -approximation. In the case of lower/upper semi-computability we can at least finitely compute lower/upper bounds to $f(x)$. In case of approximability, the weakest computability form, even this capability is lost. In analogy to lower/upper semi-computability one may think of notions like lower/upper estimability but they are easily shown to coincide with estimability. The following implications are valid:



The major algorithmic property of K is:

Theorem 2.13 ((Non)computability of Kolmogorov complexity) The Kolmogorov complexity $K : B^* \rightarrow \mathbb{N}$ is co-enumerable, but not finitely computable.

In the following we use the term computable synonymous to finitely computable, but sometimes also generically for some of the computability forms of Definition 2.12. What we call *estimable* is often just called *computable*, but it makes sense to separate the concepts of finite computability and estimability here, since the former is conceptually easier and some previous results have only been proved for this case.

2.3 Uncertainty & Probabilities

The aim of probability theory is to describe uncertainty. There are various sources for uncertainty and hence various interpretations of probabilities. There are at least three “schools”:

- the *frequentist*: probabilities are relative frequencies.
- the *objectivist*: probabilities are real aspects of the world.
- the *subjectivist*: probabilities describe an agent’s degree of belief in something.

The following subsections describe these interpretations and discuss approaches to obtain prior probabilities.

Remark. In some communities the domain of applicability and the correct interpretation and form of probability theory is still controversial. For those readers we want to emphasize that probabilities could be completely abandoned from the thesis without trivializing its goals/results. The terminology of subjective probabilities is used in this thesis for motivational and illustrational purposes only. We do not rely on Cox’s justification (see below), but give decision-theoretic justifications. Even the notion of objective probabilities may be abandoned by assuming deterministic environments. Some results in the thesis simplify in this case, but they keep their significance. So readers not believing in objective and/or subjective probabilities can still find the thesis interesting.

2.3.1 Frequency Interpretation / Counting

The *frequentist* interprets probabilities as relative frequencies. If in a sequence of n independent identically distributed (i.i.d.) experiments (trials) an event occurs $k(n)$ times the relative frequency of the event is $k(n)/n$. The limit $\lim_{n \rightarrow \infty} k(n)/n$ is *defined* as the probability of the event. This was the earliest mathematical definition of probabilities by Bernoulli, published in 1713 [Ber13]. For instance, the probability of the event *head* in a sequence of repeatedly tossing a fair coin is $\frac{1}{2}$. The frequentist position is the most easy to grasp, but it has several shortcomings:

- The frequentist obtains probabilities from physical processes as described above. To scientifically reason about probabilities one needs a mathematical theory. The problem is how to define random sequences. This is much more intricate than one might think, and has only been solved in the 1960s by Kolmogorov and Martin-Löf [ML66].
- Philosophically and also often in real experiments it is hard to justify the choice of the, so called, reference class. For instance, a doctor who wants to determine the chances that a patient has a particular disease by counting the frequency of the disease in “similar” patients. But if the doctor considered everything he knows about the patient (symptoms, weight, age, ancestry, ...) there would be no other comparable patients left.
- The frequency approach is limited to a (sufficiently large) sample of i.i.d. data.

2.3.2 Objective Interpretation: Probabilities to Describe Uncertain Events

For the *objectivist* probabilities are real aspects of the world. The outcome of an observation or an experiment is not deterministic, but involves physical random processes. The set S of all possible outcomes is called the *sample space*. It is said that an event $E \subset S$ occurred if the outcome is in E . In the case of i.i.d. experiments the probabilities assigned to events should be interpretable as limiting frequencies, but the application is not limited to this case. The Kolmogorov axioms formalize the properties which probabilities should have.

Axioms 2.14 (Kolmogorov's axioms of probability theory) Let S be the sample space. Events are subsets of S .

- If A and B are events, then also the intersection $A \cap B$, the union $A \cup B$, and the difference $A \setminus B$ are events.
- The sample space S and the empty set $\{\}$ are events.
- There is a function p which assigns non-negative reals, called probabilities, to each event.
- $p(S) = 1$, $p(\{\}) = 0$.
- $p(A \cup B) = p(A) + p(B) - p(A \cap B)$.
- For a decreasing sequence $A_1 \supset A_2 \supset A_3 \dots$ of events with $\bigcap_n A_n = \{\}$ we have $\lim_{n \rightarrow \infty} p(A_n) = 0$.

The function p is called a *probability mass function*, or, probability measure, or, more loosely *probability distribution (function)*. Conditional probabilities are defined in the following way

Definition 2.15 (Conditional probability) If A and B are events with $p(A) > 0$, then the probability that event B will occur under the condition that event A has occurred is defined as

$$p(B|A) := \frac{p(A \cap B)}{p(A)}$$

It is easy to see that $p(\cdot|A)$ (as a function of the first argument) is also a probability measure, if $p(\cdot)$ satisfies the Kolmogorov axioms. One can “verify the correctness” of the Kolmogorov axioms and the definition of conditional probabilities in the case where probabilities are identified with limiting frequencies. But the idea is to take the axioms as a starting point to avoid the frequentist's problems. $p(A \cap B) = p(B|A) \cdot p(A)$ is called the multiplication rule (of conditional probabilities), which is (a special case of) the chain rule.

Theorem 2.16 (Bayes rule) If A and B are events with $p(A) > 0$ and $p(B) > 0$, then

$$p(B|A) = \frac{p(A|B)p(B)}{p(A)}$$

Bayes' theorem is easily proven by applying Definition 2.15 twice.

2.3.3 Subjective Interpretation: Probabilities to Describe Degrees of Belief

The *subjectivist* uses probabilities to characterize an agent's degree of belief in something, rather than to characterize physical random processes. This is the most relevant interpretation of probabilities in AI. We define the plausibility of an event as the degree of belief in the event, or the subjective probability of the event. The problem with the subjective view is that it is much more arguable how to define plausibilities, as compared to objective probabilities. The objectivist can motivate Kolmogorov's axioms by a frequency analysis. There is no frequency interpretation for plausibilities. If an agent believes in extraterrestrials and assigns a plausibility of 0.9 to their existence, it does not make much sense to interpret this as 'in 90 out of 100 parallel universes there are extraterrestrials' or '90 out of 100 similar agents believe in extraterrestrials'. This problem has led to many different systems dealing with uncertain reasoning (see the reference section). They all have their own problems. The most consistent and successful system is, again, based on Kolmogorov's axioms (although not all would agree with this statement). It is surprising that plausibilities follow the same rules as limiting frequencies. It is possible to derive Kolmogorov's axioms from a few plausible qualitative rules they should follow. It is natural to assume that plausibilities can be represented by real numbers, that the rules qualitatively correspond to common sense, and that the rules are mathematically consistent. Cox [Cox46] starts with the following (natural) assumptions on beliefs.

Axioms 2.17 (Cox's axioms for beliefs)

- The degree of belief in event B (plausibility of event B), given that event A occurred can be characterized by a real-valued function $\text{Bel}(B|A)$.
- $\text{Bel}(S \setminus B|A)$ is a twice differentiable function of $\text{Bel}(B|A)$ for $A \neq \{\}$.
- $\text{Bel}(B \cap C|A)$ is a twice continuously differentiable function of $\text{Bel}(C|B \cap A)$ and $\text{Bel}(B|A)$ for $B \cap A \neq \{\}$.

Cox [Cox46] shows that every function $\text{Bel}(\cdot|\cdot)$ satisfying these axioms is isomorphic to a (conditional) probability function. One can motivate the functional relationship in Cox's axioms by analyzing all other possibilities and showing that they violate common sense [Tri69]. The somewhat strong differentiability assumptions can be weakened to more natural continuity and monotonicity assumptions [Ačz66]. Only recently, a loophole in Cox's and other's derivations have been exhibited [Par95]. Several fixes have been suggested by making additional assumptions. Most of them require the range of Bel , and hence the set of events, to be rich enough. We paraphrase these as "additional denseness conditions."

Theorem 2.18 (Cox’s theorem) Under Axioms 2.17 and some additional denseness conditions, $\text{Bel}(\cdot|A)$ is isomorphic to a probability function in the sense that there is a continuous one-to-one onto function $g: \mathbb{R} \rightarrow [0,1]$ such that $p := g \circ \text{Bel}$ satisfies Kolmogorov’s Axioms 2.14 and is consistent with Definition 2.15.

The result of Cox has attracted a great deal of interest, particularly in the maximum entropy and AI community. The qualitative motivation of Cox’s axioms and the derivation of Cox’s theorem from them is the major theoretical justification that subjective ‘degrees of belief’ should satisfy the same Kolmogorov axioms as limiting frequencies. Other approaches to beliefs are missing this strong theoretical foundation and consistency.

2.3.4 Determining Priors

The Kolmogorov axioms of probability (2.14) allow relating probabilities and plausibilities of different events, but they do not uniquely fix a numerical value for each event, except for the sure event S and the empty event $\{\}$. We need new principles for determining values for at least some basis events from which others can then be computed by (2.14). There seem to be only 3 general principles:

- The principle of indifference — the symmetry principle,
- The maximum entropy principle,
- Occam’s razor — the simplicity principle.

Whereas the first two principles are based on the foundations of statistical physics, we will see that only Occam’s razor (keep only the simplest consistent hypothesis) in combination with Epicurus’ principle of multiple explanations (keep all consistent hypotheses) is general enough to assign prior probabilities in *every* situation, especially in the case of induction and other domains typical for AI. The idea is to assign high (subjective) probability to simple events, and low probability to complex events: Simple events (strings) are more plausible a priori than complex ones. This gives (approximately) justice to both Occam’s razor and Epicurus’ principle³. Using K for measuring simplicity/complexity leads to Solomonoff’s universal prior M . In the next section we pursue this approach.

2.4 Algorithmic Probability & Universal Induction

In addition to the notation introduced in Section 2.2.1 we denote binary strings of length n by $x = x_1x_2\dots x_n$ with $x_t \in \mathbb{B}$ and further abbreviate $x_{1:n} := x_1x_2\dots x_{n-1}x_n$ and $x_{<n} := x_1\dots x_{n-1}$.

³In the following we also refer to this general idea as Occam’s razor.

2.4.1 The Universal Prior M

The prefix Kolmogorov complexity $K(x)$ has been defined as the shortest program p , for which the universal prefix Turing machine U outputs string x , and similarly $K(x|y)$ in case of side information y (Definition 2.9). Solomonoff [Sol64, Sol78] defined a closely related quantity, the universal prior $M(x)$.

It is defined as the probability that the output of a universal monotone Turing machine starts with x when provided with fair coin flips on the input tape. Formally, M can be defined as

$$M(x) := \sum_{p: U(p)=x*} 2^{-l(p)} \quad (2.19)$$

where the sum is over minimal programs p for which U outputs a string starting with x (see Definition 2.6). Since the shortest programs p dominate the sum, $M(x)$ is roughly $2^{-K(x)}$ ($M(x) = 2^{-K(x)+O(K(l(x)))}$).

Before we can discuss the stochastic properties of M we need the concept of (semi)measures for strings.

Definition 2.20 ((Semi)measures) $\mu(x)$ denotes the probability that a binary sequence starts with string x . We call $\mu \geq 0$ a semimeasure if $\mu(\varepsilon) \leq 1$ and $\mu(x) \geq \mu(x0) + \mu(x1)$, and a probability measure if equality holds.

The reason for calling μ with the above property a probability measure is that it satisfies Kolmogorov's axioms of probability (Definition 2.14) in the following sense: The sample space is \mathcal{B}^∞ with elements $\omega = \omega_1\omega_2\omega_3\ldots \in \mathcal{B}^\infty$ being infinite binary sequences. The set of events (the σ -algebra) is defined as the set generated from the cylinder sets $\Gamma_{x_{1:n}} := \{\omega : \omega_{1:n} = x_{1:n}\}$ by countable union and complement. A probability measure μ is uniquely defined by giving its values $\mu(\Gamma_{x_{1:n}})$ on the cylinder sets, which we abbreviate by $\mu(x_{1:n})$. We will also call μ a measure, or even more loose a probability distribution.

The reason for extending the definition to semimeasures is that M itself is unfortunately *not* a probability measure. We have $M(x0) + M(x1) < M(x)$ because there are programs p , which output x , neither followed by 0 nor 1. They just stop after printing x or continue forever without any further output. Since $M(\varepsilon) = 1$, M is at least a semimeasure. We can now state the fundamental property of M .

Theorem 2.21 (Universality of M) The universal prior $M(x) := \sum_{p: U(p)=x*} 2^{-l(p)}$ is an enumerable semimeasure which multiplicatively dominates all enumerable semimeasures in the sense that $M(x) \geq \sum_{\rho} 2^{-K(\rho)} \cdot \rho(x)$ if ρ is an enumerable semimeasure. M is enumerable, but not estimable or finitely computable.

The Kolmogorov complexity of a function like ρ is defined as the length of the shortest self-delimiting code of a Turing machine computing this function in the sense

of Definition 2.12. Up to a multiplicative constant, M assigns higher probability to all x than any other computable probability distribution.

It is possible to normalize M to a true probability measure M_{norm} [Sol78, LV97] with dominance still being true, but at the expense of giving up enumerability (M_{norm} is still approximable). We will see that M is more convenient when studying algorithmic questions, but a true probability measure like M_{norm} is more convenient when studying stochastic questions.

2.4.2 Universal Sequence Prediction

In which sense does M incorporate Occam's razor and Epicurus' principle of multiple explanations? From $M(x) \approx 2^{-K(x)}$ we see that M assigns high probability to simple strings (Occam). More useful is to think of x as being the observed history. We see from Definition 2.19 that every program p consistent with history x is allowed to contribute to M (Epicurus). On the other hand shorter programs give significantly larger contribution (Occam). How does all this affect prediction? If $M(x)$ correctly describes our (subjective) prior belief in x , then $M(y|x) := M(xy)/M(x)$ must be our posterior belief in y . From the symmetry of algorithmic information $K(x,y) \pm K(y|x, K(x)) + K(x)$ (Theorem 2.10(v)), and assuming $K(x,y) \approx K(xy)$, and approximating $K(y|x, K(x)) \approx K(y|x)$, $M(x) \approx 2^{-K(x)}$, and $M(xy) \approx 2^{-K(xy)}$ we get $M(y|x) \approx 2^{-K(y|x)}$. This tells us that M predicts y with high probability iff y has an easy explanation, given x (Occam & Epicurus).

The above qualitative discussion should not create the impression that $M(x)$ and $2^{-K(x)}$ always lead to predictors of comparable quality. Indeed in the on-line/incremental setting studied in this work, $K(y) = O(1)$ invalidates the consideration above. The validity of (2.23) below for instance depends on M being a semimeasure and the chain rule being exactly true, neither of them is satisfied by $2^{-K(x)}$.

(Binary) sequence prediction algorithms try to predict the continuation $x_n \in \mathcal{B}$ of a given sequence $x_1 \dots x_{n-1}$. We derive the following bound:

$$\sum_{t=1}^{\infty} (1 - M(x_t | x_{<t}))^2 \leq -\frac{1}{2} \sum_{t=1}^{\infty} \ln M(x_t | x_{<t}) = -\frac{1}{2} \ln M(x_{1:\infty}) \leq \frac{1}{2} \ln 2 \cdot Km(x_{1:\infty}), \quad (2.22)$$

where the monotone complexity $Km(x_{1:\infty})$ is defined as the length of the shortest (non-halting) program computing $x_{1:\infty}$ [ZL70]. In the first inequality we have used $(1-a)^2 \leq -\frac{1}{2} \ln a$ for $0 \leq a \leq 1$. In the equality we exchanged the sum with the logarithm and eliminated the resulting product by the chain rule. In the last inequality we used $M(x) \geq 2^{-Km(x)}$, which follows from definition (2.19) by dropping all terms in \sum_p except for the shortest p computing x . If $x_{1:\infty}$ is a computable sequence, then $Km(x_{1:\infty})$ is finite, which implies $M(x_t | x_{<t}) \rightarrow 1$ ($\sum_{t=1}^{\infty} (1-a_t)^2 < \infty \Rightarrow a_t \rightarrow 1$). This means, that if the environment is a computable sequence (whichsoever, e.g. the digits π or e in binary representation), after having seen the first few digits, M

correctly predicts the next digit with high probability, i.e. it recognizes the structure of the sequence.

Assume now that the true sequence is drawn from a computable probability distribution μ , i.e. the true (objective) probability of $x_{1:n}$ is $\mu(x_{1:n})$. The probability of x_n given $x_{<n}$ hence is $\mu(x_n|x_{<n}) = \mu(x_{1:n})/\mu(x_{<n})$. Solomonoff's [Sol78] central result is that M converges to μ . More precisely he showed that

$$\sum_{t=1}^{\infty} \sum_{x_{<t} \in \mathcal{B}^{t-1}} \mu(x_{<t}) \left(M(0|x_{<t}) - \mu(0|x_{<t}) \right)^2 \stackrel{+}{\leq} \frac{1}{2} \ln 2 \cdot K(\mu) < \infty \quad (2.23)$$

The infinite sum can only be finite if the difference $M(0|x_{<t}) - \mu(0|x_{<t})$ tends to zero for $t \rightarrow \infty$ with μ probability 1 (see Definition 3.8(i)). This holds for *any* computable probability distribution μ . The reason for the astonishing property of a single (universal) function to converge to *any* computable probability distribution lies in the fact that the set of μ -random sequences differ for different μ . Past data $x_{<t}$ are exploited to get a (with $t \rightarrow \infty$) improving estimate $M(x_t|x_{<t})$ of $\mu(x_t|x_{<t})$.

The universality property (Theorem 2.21) is the central ingredient in the proof of (2.23). The proof of Theorem 2.21 involves the construction of a semimeasure ξ whose dominance is obvious. The hard part is to show its enumerability and equivalence to M . Let \mathcal{M} be the (countable) set of all enumerable semimeasures and define

$$\xi(x) := \sum_{\nu \in \mathcal{M}} 2^{-K(\nu)} \nu(x). \quad (2.24)$$

Then dominance

$$\xi(x) \geq 2^{-K(\nu)} \nu(x) \quad \forall \nu \in \mathcal{M} \quad (2.25)$$

is obvious (without $O(1)$ fudge). Is ξ lower semi-computable? To answer this question we have to be more precise. Levin [ZL70] has shown that there is a Turing machine such that for every lower semi-computable semimeasure ν there is an i such that $T(i \cdot x)$ lower semi-computes $\nu_i \equiv \nu$, i.e. T enumerates *all* lower semi-computable semimeasures, possibly with repetition. For the (ordered multi) set $\mathcal{M} = \mathcal{M}_U := \{\nu_1, \nu_2, \nu_3, \dots\}$ and $K(\nu_i) := K(i)$ one can easily see that ξ is lower semi-computable. Finally proving $M(x) \stackrel{\times}{=} \xi(x)$ also establishes universality of M .

The advantage of ξ over M is that it immediately generalizes to arbitrary weighted sums of (semi)measures in \mathcal{M} for arbitrary countable \mathcal{M} . Most proofs in this work go through for generic \mathcal{M} and weights.

2.4.3 Universal (Semi)Measures

What is so special about the set of all enumerable semimeasures \mathcal{M}_U ? The larger we choose \mathcal{M} the less restrictive is the assumption that \mathcal{M} should contain the true distribution μ , which will be essential throughout the thesis. Why do not restrict to the still rather general class of estimable or finitely computable (semi)measures? It is clear that for every countable set \mathcal{M} , $\xi(x) := \xi_{\mathcal{M}}(x) := \sum_{\nu \in \mathcal{M}} w_{\nu} \nu(x)$ with $\sum_{\nu} w_{\nu} \leq 1$

and $w_\nu > 0$ dominates all $\nu \in \mathcal{M}$. This dominance is necessary for the desired convergence $\xi \rightarrow \mu$ similarly to (2.23). The question is what properties ξ possesses. The distinguishing property of \mathcal{M}_U is that $\xi = \xi_U \equiv \xi_{\mathcal{M}_U} \stackrel{\times}{=} M$ is itself an element of \mathcal{M}_U . In this work $\xi_{\mathcal{M}} \in \mathcal{M}$ is not by itself an important property, but whether ξ is computable in one of the senses of Definition 2.12. We define

$$\begin{aligned} \mathcal{M}_1 \stackrel{\times}{\geq} \mathcal{M}_2 &: \Leftrightarrow \text{there is an element of } \mathcal{M}_1 \text{ which dominates all elements of } \mathcal{M}_2 \\ &: \Leftrightarrow \exists \rho \in \mathcal{M}_1 \forall \nu \in \mathcal{M}_2 \exists w_\nu > 0 \forall x : \rho(x) \geq w_\nu \nu(x). \end{aligned}$$

$\stackrel{\times}{\geq}$ is transitive (but not necessarily reflexive) in the sense that $\mathcal{M}_1 \stackrel{\times}{\geq} \mathcal{M}_2 \stackrel{\times}{\geq} \mathcal{M}_3$ implies $\mathcal{M}_1 \stackrel{\times}{\geq} \mathcal{M}_3$ and $\mathcal{M}_0 \supseteq \mathcal{M}_1 \stackrel{\times}{\geq} \mathcal{M}_2 \supseteq \mathcal{M}_3$ implies $\mathcal{M}_0 \stackrel{\times}{\geq} \mathcal{M}_3$. For the computability concepts introduced in Section 2.2.4 we have the following proper set inclusions

$$\begin{array}{ccccccc} \mathcal{M}_{comp}^{msr} & \subset & \mathcal{M}_{est}^{msr} & \equiv & \mathcal{M}_{enum}^{msr} & \subset & \mathcal{M}_{appr}^{msr} \\ \cap & & \cap & & \cap & & \cap \\ \mathcal{M}_{comp}^{semi} & \subset & \mathcal{M}_{est}^{semi} & \subset & \mathcal{M}_{enum}^{semi} & \subset & \mathcal{M}_{appr}^{semi} \end{array}$$

where \mathcal{M}_c^{msr} stands for the set of all probability measures of appropriate computability type $c \in \{\text{comp}=\text{finitely computable}, \text{est}=\text{estimable}, \text{enum}=\text{enumerable}, \text{appr}=\text{approximable}\}$, and similarly for semimeasures \mathcal{M}_c^{semi} . Other classes are briefly discussed in Section 3.2.9. From an enumeration of a measure ρ one can construct a co-enumeration by exploiting $\rho(x_{1:n}) = 1 - \sum_{y_{1:n} \neq x_{1:n}} \rho(y_{1:n})$. This shows that every enumerable measure is also co-enumerable, hence estimable, which proves the identity \equiv above.

With this notation, [ZL70, Th.3.3] reads $\mathcal{M}_{enum}^{semi} \stackrel{\times}{\geq} \mathcal{M}_{enum}^{semi}$. Transitivity allows to conclude, for instance, that $\mathcal{M}_{appr}^{semi} \stackrel{\times}{\geq} \mathcal{M}_{comp}^{msr}$, i.e. that there is an approximable semimeasure which dominates all computable measures.

The standard “diagonalization” way of proving $\mathcal{M}_1 \not\stackrel{\times}{\geq} \mathcal{M}_2$ is to take an arbitrary $\mu \in \mathcal{M}_1$ and “increase” it to ρ such that $\mu \not\stackrel{\times}{\geq} \rho$ and show that $\rho \in \mathcal{M}_2$. There are 7×7 combinations of (semi)measures \mathcal{M}_1 with \mathcal{M}_2 for which $\mathcal{M}_1 \stackrel{\times}{\geq} \mathcal{M}_2$ could be true or false. There are four basic cases, explicated in the following theorem, from which the other 49 combinations displayed in Table 2.27 follow by transitivity.

Theorem 2.26 (Universal (semi)measures) A semimeasure ρ is said to be universal for \mathcal{M} if it multiplicatively dominates all elements of \mathcal{M} in the sense $\forall \nu \exists w_\nu > 0 : \rho(x) \geq w_\nu \nu(x) \forall x$. The following holds true:

- o)* $\exists \rho : \{\rho\} \stackrel{\times}{\geq} \mathcal{M}$: For every countable set of (semi)measures \mathcal{M} , there is a (semi)measure which dominates all elements of \mathcal{M} .
- i)* $\mathcal{M}_{enum}^{semi} \stackrel{\times}{\geq} \mathcal{M}_{enum}^{semi}$: The class of enumerable semimeasures *contains* a universal element.
- ii)* $\mathcal{M}_{appr}^{msr} \stackrel{\times}{\geq} \mathcal{M}_{enum}^{semi}$: There *is* an approximable measure which dominates all enumerable semimeasures.
- iii)* $\mathcal{M}_{est}^{semi} \not\stackrel{\times}{\geq} \mathcal{M}_{comp}^{msr}$: There is *no* estimable semimeasure which dominates all computable measures.
- iv)* $\mathcal{M}_{appr}^{semi} \not\stackrel{\times}{\geq} \mathcal{M}_{appr}^{msr}$: There is *no* approximable semimeasure which dominates all approximable measures.

Table 2.27 (Existence of universal (semi)measures) The entry in row r and column c indicates whether there is a r -able (semi)measure ρ for the set \mathcal{M} which contains all c -able (semi)measures, where $r, c \in \{\text{comput}, \text{estimat}, \text{enumer}, \text{approxim}\}$. Enumerable measures are estimable. This is the reason why the enum. row and column in case of measures is missing. The superscript indicates from which part of Theorem 2.26 the answer follows. For the bold face entries directly, for the others using transitivity of $\stackrel{\times}{\geq}$.

\nwarrow	\mathcal{M}	semimeasure				measure		
ρ	\searrow	comp.	est.	enum.	appr.	comp.	est.	appr.
s	comp.	no ⁱⁱⁱ	no ⁱⁱⁱ	no ⁱⁱⁱ	no ^{iv}	no ⁱⁱⁱ	no ⁱⁱⁱ	no ^{iv}
e	est.	no ⁱⁱⁱ	no ⁱⁱⁱ	no ⁱⁱⁱ	no ^{iv}	noⁱⁱⁱ	no ⁱⁱⁱ	no ^{iv}
m	enum.	yes ⁱ	yes ⁱ	yesⁱ	no ^{iv}	yes ⁱ	yes ⁱ	no ^{iv}
i	appr.	yes ⁱ	yes ⁱ	yes ⁱ	no ^{iv}	yes ⁱ	yes ⁱ	no^{iv}
m	comp.	no ⁱⁱⁱ	no ⁱⁱⁱ	no ⁱⁱⁱ	no ^{iv}	no ⁱⁱⁱ	no ⁱⁱⁱ	no ^{iv}
s	est.	no ⁱⁱⁱ	no ⁱⁱⁱ	no ⁱⁱⁱ	no ^{iv}	no ⁱⁱⁱ	no ⁱⁱⁱ	no ^{iv}
r	appr.	yes ⁱⁱ	yes ⁱⁱ	yesⁱⁱ	no ^{iv}	yes ⁱⁱ	yes ⁱⁱ	no ^{iv}

If we ask for a universal (semi)measure which at least satisfies the weakest form of computability, namely being approximable, we see that the largest dominated set among the 7 sets defined above is the set of enumerable semimeasures. This is the reason why $\mathcal{M}_{enum}^{semi}$ plays a special role in this (and other) works. On the other hand, $\mathcal{M}_{enum}^{semi}$ is not the largest set dominated by an approximable semimeasure, and indeed no such largest set exists. One may, hence, ask for “natural” larger

sets \mathcal{M} . One such set, namely the set of cumulatively enumerable semimeasures \mathcal{M}_{CEM} , has recently been discovered by Schmidhuber [Sch00, Sch02a], for which even $\xi_{CEM} \in \mathcal{M}_{CEM}$ holds.

Theorem 2.26 also holds for *discrete (semi)measures* P defined as

$$P : \mathbb{N} \rightarrow [0, 1] \quad \text{with} \quad \sum_{x \in \mathbb{N}} P(x) \stackrel{(\leq)}{=} 1.$$

We first prove the theorem for this discrete case, since it contains the essential ideas in a cleaner form. We then present the proof for “continuous” (semi)measures μ (Definition 2.20). The proofs naturally generalize from binary to arbitrary finite alphabet. $\text{argmin}_x f(x)$ is the x that minimizes $f(x)$. Ties are broken in an arbitrary but computable way (e.g. by taking the smallest x).

Proof (discrete case).

(o) $Q(x) := \sum_{P \in \mathcal{M}} w_P P(x)$ with $w_P > 0$ obviously dominates all $P \in \mathcal{M}$ (with constant w_P). With $\sum_P w_P = 1$ and all P being discrete (semi)measures also Q is a discrete (semi)measure.

(i) See [LV97, Th.4.3.1].

(ii) Let P be the universal element in $\mathcal{M}_{enum}^{semi}$ and $\alpha := \sum_x P(x)$. We normalize P by $Q(x) := \frac{1}{\alpha} P(x)$. Since $\alpha \leq 1$ we have $Q(x) \geq P(x)$. Hence $Q \geq P \stackrel{\times}{\geq} \mathcal{M}_{enum}^{semi}$. As a ratio between two enumerable functions, Q is still approximable, hence $\mathcal{M}_{appr}^{msr} \stackrel{\times}{\geq} \mathcal{M}_{enum}^{semi}$.

(iii) ⁴Let $P \in \mathcal{M}_{comp}^{semi}$. We partition \mathbb{N} into chunks $I_n := \{2^{n-1}, \dots, 2^n - 1\}$ ($n \geq 1$) of increasing size. With $x_n := \text{argmin}_{x \in I_n} P(x)$ we define $Q(x_n) := \frac{1}{n(n+1)} \forall n$ and $Q(x) := 0$ for all other x . Exploiting that a minimum is smaller than an average we get

$$P(x_n) = \min_{x \in I_n} P(x) \leq \frac{1}{|I_n|} \sum_{x \in I_n} P(x) \leq \frac{1}{|I_n|} = \frac{1}{2^{n-1}} = \frac{n(n+1)}{2^{n-1}} Q(x_n)$$

Since $\frac{n(n+1)}{2^{n-1}} \rightarrow 0$ for $n \rightarrow \infty$, P cannot dominate Q ($P \not\stackrel{\times}{\geq} Q$). With P also Q is computable. Since P was an arbitrary computable semimeasure and Q is a computable measure ($\sum Q(x) = \sum [\frac{1}{n(n+1)}] = \sum [\frac{1}{n} - \frac{1}{n+1}] = 1$) this implies $\mathcal{M}_{comp}^{semi} \stackrel{\times}{\geq} \mathcal{M}_{comp}^{msr}$.

Assume now that there is an estimable semimeasure $S \stackrel{\times}{\geq} \mathcal{M}_{comp}^{msr}$. We construct a finitely computable semimeasure $P \stackrel{\times}{\geq} S$ as follows. Choose an initial $\varepsilon > 0$ and finitely compute an ε -approximation \hat{S} of $S(x)$. If $\hat{S} > 2\varepsilon$ define $P(x) := \frac{1}{2}\hat{S}$, else halve ε and repeat the process. Since $S(x) > 0$ (otherwise it could not dominate, e.g. $T(x) := \frac{1}{x(x+1)} \in \mathcal{M}_{comp}^{msr}$) the loop terminates after finite time. So P

⁴The proof of $\mathcal{M}_{comp}^{semi} \stackrel{\times}{\geq} \mathcal{M}_{comp}^{semi}$ in [LV97, p249] contains minor errors and is not extensible to $\mathcal{M}_{est}^{semi} \stackrel{\times}{\geq} \mathcal{M}_{comp}^{msr}$.

is finitely computable. Inserting $\hat{S} = 2P(x)$ and $\varepsilon < \frac{1}{2}\hat{S} = P(x)$ into $|S(x) - \hat{S}| < \varepsilon$ we get $|S(x) - 2P(x)| < P(x)$, which implies $S(x) \geq P(x)$ and $S(x) \leq 3P(x)$. The former implies $\sum_x P(x) \leq \sum_x S(x) \leq 1$, i.e. P is a semimeasure. The latter implies $P \geq \frac{1}{3}S \geq \mathcal{M}_{comp}^{msr}$. Hence P is a computable semimeasure dominating all computable measures, which contradicts what we have proven in the first half of (iii). Hence the assumption on S was wrong which establishes $\mathcal{M}_{est}^{semi} \not\geq^{\times} \mathcal{M}_{comp}^{msr}$.

(iv) Assume $P \in \mathcal{M}_{app}^{semi}$. We construct an approximable measure Q which is not dominated by P , thus contradicting the assumption. Let P_1, P_2, \dots be a sequence of recursive functions converging to P . We construct x_1, x_2, \dots such that $P(x_n) \not\geq c \cdot Q(x_n) \forall n$ for any constant c . For this we recursively define sequences x_n^1, x_n^2, \dots converging to x_n and from them Q_1, Q_2, \dots converging to Q . Let $I_n := \{2^{n-1}, \dots, 2^n - 1\}$ and $x_n^1 = 2^{n-1} \forall n$. If $P_t(x_n^{t-1}) > n^{-3}$ then $x_n^t := \operatorname{argmin}_{x \in I_n} P_t(x)$ else $x_n^t := x_n^{t-1}$. We show that x_n^t converges for $t \rightarrow \infty$ by assuming the contrary and showing a contradiction. Since $x_n^t \in I_n$ some value, say x_n^* , is assumed infinitely often. Non-convergence implies that the sequence leaves and returns to x_n^* infinitely often. x_n^* is only left ($x_n^{t-1} = x_n^* \neq x_n^t$) if $P_t(x_n^*) > n^{-3}$. On the other hand at the time where x_n^t returns to x_n^* ($x_n^{t-1} \neq x_n^* = x_n^t$) we have $P_t(x_n^*) = P_t(x_n^t) = \min_{x \in I_n} P_t(x) \leq |I_n|^{-1} = 2^{-n+1}$. Hence $P_t(x_n^*)$ oscillates infinitely often between $\leq 2^{-n+1}$ and $\geq n^{-3}$ which contradicts the assumption that P_t converges. Hence the assumption of a non-convergent x_n^t was wrong. x_n^t converges to x_n^* and $P_t(x_n^*)$ to a value $\leq n^{-3}$. With x_n^t also the measure $Q_t(x_n^t) := \frac{1}{n(n+1)}$ (and $Q_t(x) = 0$ for all other x) converges. Since $P(x_n^*) \leq n^{-3}$ does not dominate $Q(x_n^*)$, we have $P \not\geq^{\times} Q$. Since $P \in \mathcal{M}_{app}^{semi}$ was arbitrary and Q is an approximable measure we get $\mathcal{M}_{appr}^{semi} \not\geq^{\times} \mathcal{M}_{appr}^{msr}$. \square

Proof (continuous case).

The major difference to the discrete case is that one also has to take care that $\rho(x) \geq \rho(x_0) + \rho(x_1)$, $x \in \mathbb{B}^*$, is respected. On the other hand the chunking $I_n := \mathbb{B}^n$ is more natural here.

(o) $\rho(x) := \sum_{\nu \in \mathcal{M}} w_\nu \nu(x)$ with $w_\nu > 0$ obviously dominates all $\nu \in \mathcal{M}$ (with domination constant w_ν). With $\sum_\nu w_\nu = 1$ and all ν being (semi)measures also ρ is a (semi)measure.

(i) See [LV97, Th4.5.1].

(ii) Let ξ be a universal element in $\mathcal{M}_{enum}^{semi}$. We define [Sol78]

$$\xi_{norm}(x_{1:n}) := \prod_{t=1}^n \frac{\xi(x_{1:t})}{\xi(x_{<t}0) + \xi(x_{<t}1)}.$$

By induction one can show that ξ_{norm} is a measure and that $\xi_{norm}(x) \geq \xi(x) \forall x$, hence $\xi_{norm} \geq \xi \geq^{\times} \mathcal{M}_{enum}^{semi}$. As a ratio of enumerable functions, ξ_{norm} is still approximable, hence $\mathcal{M}_{appr}^{msr} \geq^{\times} \mathcal{M}_{enum}^{semi}$.

(iii) ⁵Let $\mu \in \mathcal{M}_{comp}^{semi}$. We recursively define the sequence $x_{1:\infty}^*$ by $x_k^* := \operatorname{argmin}_{x_k} \mu(x_{<k}^* x_k)$ and the measure ρ by $\rho(x_{1:k}^*) = 1 \forall k$ and $\rho(x) = 0$ for all x which are not prefixes of $x_{1:\infty}^*$. Exploiting the fact that a minimum is smaller than an average we get

$$\mu(x_{1:k}^*) = \min_{x_k} \mu(x_{<k}^* x_k) \leq \frac{1}{2} [\mu(x_{<k}^* 0) + \mu(x_{<k}^* 1)] \leq \frac{1}{2} \mu(x_{<k}^*).$$

Hence $\mu(x_{1:n}^*) \leq (\frac{1}{2})^n = (\frac{1}{2})^n \rho(x_{1:n}^*)$ which demonstrates that μ does not dominate ρ . Since $\mu \in \mathcal{M}_{comp}^{semi}$ was arbitrary and ρ is a computable measure this implies

$$\mathcal{M}_{comp}^{semi} \not\stackrel{\times}{\geq} \mathcal{M}_{comp}^{msr}.$$

Assume now that there is an estimable semimeasure $\sigma \stackrel{\times}{\geq} \mathcal{M}_{comp}^{msr}$. We construct a finitely computable function $\mu \stackrel{\times}{\geq} \sigma$ as follows. Choose an initial $\varepsilon > 0$ and finitely compute an ε -approximation $\hat{\sigma}$ of $\sigma(x)$. If $\hat{\sigma} > 4\varepsilon$ define $\mu(x) := \hat{\sigma}$, else halve ε and repeat the process. Since $\sigma(x) > 0$ (otherwise it could not dominate, e.g. $2^{-l(x)}$) the loop terminates after finite time. So μ is finitely computable. Inserting $\hat{\sigma} = \mu(x)$ and $\varepsilon < \frac{1}{4}\hat{\sigma} = \frac{1}{4}\mu(x)$ into $|\sigma(x) - \hat{\sigma}| < \varepsilon$ we get $|\sigma(x) - \mu(x)| < \frac{1}{4}\mu(x)$, which implies $\frac{3}{4}\mu(x) \leq \sigma(x) \leq \frac{5}{4}\mu(x)$. Unfortunately μ is not a semimeasure, but it still satisfies the weaker inequality $\mu(x0) + \mu(x1) \leq \frac{4}{3}[\sigma(x0) + \sigma(x1)] \leq \frac{4}{3}\sigma(x) \leq \frac{4}{3} \cdot \frac{5}{4}\mu(x) = \frac{5}{3}\mu(x)$. This is sufficient for the first half of the proof of (iii) to go through with $\frac{1}{2}$ replaced by $\frac{1}{2} \cdot \frac{5}{3} = \frac{5}{6} < 1$, which shows that $\mu \not\stackrel{\times}{\geq} \mathcal{M}_{comp}^{msr}$. But this contradicts $\mu \geq \frac{4}{5}\sigma \stackrel{\times}{\geq} \mathcal{M}_{comp}^{msr}$ showing that our assumed estimable semimeasure σ does not exist, i.e. $\mathcal{M}_{est}^{semi} \not\stackrel{\times}{\geq} \mathcal{M}_{comp}^{msr}$.

(iv) Assume $\mu \in \mathcal{M}_{app}^{semi}$. We construct an approximable measure ρ which is not dominated by μ , thus contradicting the assumption. Let μ_1, μ_2, \dots be a sequence of recursive functions converging to μ . We recursively (in t and n) define sequences y_n^1, y_n^2, \dots converging to y_n and from them ρ_1, ρ_2, \dots converging to ρ . Let $y_n^1 = 0 \forall n$. If $\mu_t(y_{<n}^t y_n^{t-1}) > \frac{2}{3}\mu_t(y_{<n}^t)$ then $y_n^t := \operatorname{argmin}_{x_n} \mu_t(y_{<n}^t x_n)$ else $y_n^t := y_n^{t-1}$. We show that y_n^t converges for $t \rightarrow \infty$ by assuming the contrary and showing a contradiction. Assume that k is the smallest n for which $y_n^t \not\rightarrow y_n$. Since $y_n^t \rightarrow y_n$ for all $n < k$ and $y_n^t \in \mathcal{B}$ is discrete there is a t_0 such that $y_{<k}^t = y_{<k} \forall t > t_0$. Assume $t > t_0$ in the following. Since $y_k^t \in \mathcal{B}$, some value, say \tilde{y}_k , is assumed infinitely often. Non-convergence implies that the sequence leaves and enters to \tilde{y}_k infinitely often. If \tilde{y}_k is left ($y_k^{t-1} = \tilde{y}_k \neq y_k^t$) we have

$$\mu_t(y_{<k} \tilde{y}_k) = \mu_t(y_{<k}^t y_k^{t-1}) > \frac{2}{3}\mu_t(y_{<k}^t) = \frac{2}{3}\mu_t(y_{<k}) \xrightarrow{t \rightarrow \infty} \frac{2}{3}\mu(y_{<k}).$$

If \tilde{y}_k is entered ($y_k^{t-1} \neq \tilde{y}_k = y_k^t$) we have

$$\mu_t(y_{<k} \tilde{y}_k) = \mu_t(y_{<k}^t y_k^t) = \min_{x_k} \mu_t(y_{<k}^t x_k) \leq \frac{1}{2} [\mu_t(y_{<k}^t 0) + \mu_t(y_{<k}^t 1)] \leq$$

⁵The proof in [LV97, p276] only applies to infinite alphabet and not to the binary/finite case considered here.

$$\leq \frac{1}{2}\mu_t(y_{<k}^t) = \frac{1}{2}\mu_t(y_{<k}) \xrightarrow{t \rightarrow \infty} \frac{1}{2}\mu(y_{<k}).$$

Hence $\mu_t(y_{<k}\tilde{y}_k)$ oscillates infinitely often between $> \frac{2}{3}\mu(y_{<k})$ and $\leq \frac{1}{2}\mu(y_{<k})$ which contradicts the assumption that μ_t converges. Hence the assumption of a non-convergent y_k^t was wrong. With y_k^t also the measure $\rho_t(y_{1:n}^t) := 1$ (and $\rho_t(x) = 0$ for all other x which are not prefixes of $y_{1:\infty}^t$) converges. For all sufficiently large t we have $y_{1:n} = y_{1:n}^t$, hence $\mu_t(y_{1:n}) = \mu_t(y_{1:n}^t) \leq \frac{2}{3}\mu_t(y_{<n}^t) \leq \dots \leq (\frac{2}{3})^n$. Since $\mu(y_{1:n}) \leq (\frac{2}{3})^n$ does not dominate $\rho(y_{1:n}) = 1$ ($\forall t > t_0$), we have $\mu \not\stackrel{\times}{\geq} \rho$. Since $\mu \in \mathcal{M}_{app}^{semi}$ was arbitrary and ρ is an approximable measure we get $\mathcal{M}_{app}^{semi} \not\stackrel{\times}{\supseteq} \mathcal{M}_{app}^{msr}$. \square

2.4.4 Martin-Löf Randomness

Martin-Löf randomness is a very important concept of randomness of individual sequences which is closely related to Kolmogorov complexity and Solomonoff's universal prior. Since we refer to this concept only in Section 3.2.7 we will be very brief here. We give a characterization equivalent to Martin-Löf's original definition, in order to bypass the necessity of giving a formal definition of 'effective randomness tests' [Lev73a]:

Theorem 2.28 (Martin-Löf random sequences) A sequence $x_{1:\infty}$ is called μ -Martin-Löf random (μ .M.L.) iff there is a constant c such that $M(x_{1:n}) \leq c \cdot \mu(x_{1:n})$ for all n .

An equivalent formulation for computable μ is:

$$x_{1:\infty} \text{ is } \mu\text{-M.L. random} \quad \Leftrightarrow \quad Km(x_{1:n}) \stackrel{\pm}{=} -\log_2 \mu(x_{1:n}) \quad \forall n, \quad (2.29)$$

where $Km(x_{1:n})$ is the length of the shortest (possibly non-halting) program computing a string starting with $x_{1:n}$. Theorem 2.28 follows from (2.29) by exponentiation, "using $2^{-Km} \approx M$ " and noting that $M \stackrel{\times}{\geq} \mu$ follows from universality of M . Consider the special case of μ being a fair coin, i.e. $\mu(x_{1:n}) = 2^{-n}$, then $x_{1:\infty}$ is M.L. random iff $Km(x_{1:n}) \stackrel{\pm}{=} n$, i.e. if $x_{1:n}$ is incompressible. For general μ , $-\log_2 \mu(x_{1:n})$ is the length of the Shannon-Fano code of $x_{1:n}$, hence $x_{1:\infty}$ is μ .M.L. random iff the Shannon-Fano code is optimal.

One can show that a μ .M.L. random sequence $x_{1:\infty}$ passes *all* thinkable effective randomness tests, e.g. the law of large numbers, the law of the iterated logarithm, etc. In particular, the set of all μ .M.L. random sequences has μ -measure 1. The following generalization is natural when considering general Bayes-mixtures ξ :

Definition 2.30 (μ/ξ -random sequences) A sequence $x_{1:\infty}$ is called μ/ξ -random (μ . ξ .r.) iff there is a constant c such that $\xi(x_{1:n}) \leq c \cdot \mu(x_{1:n})$ for all n .

Typically, ξ is a mixture over some \mathcal{M} as defined in (2.24), in which case the reverse inequality $\xi(x) \stackrel{\times}{\geq} \mu(x)$ is also true (for all x). For finite \mathcal{M} or if $\xi \in \mathcal{M}$, the definition of μ/ξ -randomness depends only on \mathcal{M} , and not on the specific weights used in ξ . For $\mathcal{M} = \mathcal{M}_U$, μ/ξ -randomness is just μ M.L. randomness. The larger \mathcal{M} , the more patterns are recognized as non-random. Roughly speaking, those regularities characterized by some $\nu \in \mathcal{M}$ are recognized by μ/ξ -randomness, i.e. for $\mathcal{M} \subset \mathcal{M}_U$ some μ/ξ -random strings may not be M.L. random. Other randomness concepts, e.g. those by Schnorr, Ko, van Lambalgen, Lutz, Kurtz, von Mises, Wald, and Church (see [Wan96, Lam87, Sch71]), could possibly also be characterized in terms of μ/ξ -randomness for particular choices of \mathcal{M} .

2.5 History & References

Most notation is taken over from [LV97]. The general theory of coding and prefix codes can be found in [Gal68], the important Kraft inequality is due to Kraft [Kra49].

Algorithmic Information Theory. Turing introduced the concept of a Turing machine and demonstrated that the halting problem is undecidable in [Tur36]. Turing machines are formally equivalent to partial recursive functions (see [Rog67, Odi89, Odi99] for an introduction). The halting problem corresponds to Gödel's incompleteness theorem [Göd31, Sho67] whose proof is based on a diagonal argument invented by Cantor [Can1874, Dau90]. The explicit statement of the short compiler Assumption 2.5 is not based on any source. The works [Göd31, Kle36, Tur36, Pos44, ZL70, Sch02a] show the importance of the various computability concepts defined in (2.12). The consideration of (and naming for) estimable functions in the context of universal priors is new.

A coarse picture of the early history of algorithmic information theory could be drawn as follows: Kolmogorov [Kol65] and Chaitin [Cha66, Cha69], suggested defining the information content of an object as the length of the shortest program computing a representation of it. Solomonoff [Sol64] independently invented the closely related universal prior probability distribution and used it for binary sequence prediction [Sol64, Sol78]. Levin worked out most of the mathematical details [ZL70, Lev74] and invented the fastest algorithm for function inversion and optimization, save for a (huge) constant factor [Lev73b]. These papers may be regarded as the invention of what is now called Algorithmic Information theory. The invariance (2.8) is due to [Sol64, Kol65, Cha69], Theorem 2.10(vii) is due to [Lev74], the symmetry of information (vi) due to [ZL70, Gác74, Kol83], (ii) is due to [Lev74], the other parts are elementary.

The short introduction we gave in this chapter necessarily described only the key ideas, ignoring many related and especially newer developments. Some references are given in the following.

There are many variants of “Kolmogorov” complexity. The prefix Kolmogorov complexity K we defined here [Lev74, Gác74, Cha75], the earliest form, “plain”

Kolmogorov complexity C [Kol65], process complexity [Sch73], monotone complexity Km [Lev73a], and uniform complexity [Lov69b, Lov69a], Solomonoff's universal prior $M = 2^{-KM}$ [Sol64, Sol78], Chaitin's complexity Kc [Cha75], extension semimeasure Mc [Cov74], and some others. They often differ from K only by $O(\log K)$, but have otherwise similar properties. For an introduction to Shannon's information theory [Sha48] and its relation to Kolmogorov complexity, see [Kol65, Kol83, ZL70, CT91].

The main drawback of (all these variants of) Kolmogorov complexity is that they are not finitely computable [Kol65, Sol64]. They may be approximated from above [Kol65, Sol64], but no accuracy guarantee can be given, and what is worse, the best upper bound for the running time until one has reasonable accuracy for $K(n)$ grows faster than any computable function in n . This led to the development of time bounded complexity/probability which is finitely computable, or more general resource bounded complexity/probability (e.g. space) [Dal73, Dal77, FMG92, Ko86, PF97, Sch02c].

For an excellent introduction to algorithmic information theory, and a much more accurate treatment of its history (more than 500 references), and many applications one should consult the authoritative book of Li and Vitányi [LV97].

Foundations of probability theory. Although games of chance date back at least to around 300 B.C., the first mathematical analysis of probabilities appears to be much later. Important breakthroughs have been achieved (in chronological order and with significant simplification) by Cardano [Car1565/1663], a systematic way of calculating probabilities by Pascal (in correspondence with Fermat) and conditional probability [Pas1654], Bayes rule [Bay1763], the distinction between subjective and objective interpretation of probabilities and the weak law of large numbers by Bernoulli [Ber1713], equi-probability due to symmetry and other things by Laplace [Lap1814], the principle of indifference by Keynes [Key21], Kolmogorov's axioms of probability theory [Kol33], early attempts to define the notion of randomness of individual objects/sequences by von Mises, Wald and Church [Mis19, Wal37, Chu40], finally successful by Martin-Löf [ML66], the notion of a universal a priori probability by Solomonoff [Sol64], and its mathematical investigation by Levin [ZL70, Lev74].

There is an ongoing debate between objective and subjective probability, which became sharper in the 20th century (not only in AI). Prominent advocates of the relative frequency or objective interpretation were Kolmogorov [Kol63], Fisher [Fis22], and von Mises [Mis28]. There are many advocates of probabilities as degrees of belief [Pop34, Ram31, Fin37, Cox46, Sav54, Jef83]. Carnap [Car48, Car50] tried to supplement logic with probability theory to, so called, inductive logic. This works fine for propositional logic [Jay96], but not for predicate logic [Put63]. The closely related reference class problem is addressed in [Rei49, Kyb77, Kyb83, BGHK92].

There are many books on probability theory with different focus. For a thorough treatment of the early history of the concept of probability the reader is referred to the books by Hacking [Hac75] and Hald [Hal90], and for the foundations developed in the 20th century to the book by Li & Vitányi [LV97]. A good standard textbook

is by Feller [Fel68]. A pleasant to read book with a philosophical touch is by Jaynes [Jay96]. It treats probability theory as a natural extension to (Boolean) logical reasoning, emphasizes the “full” Bayesian approach with priors determined by the maximum entropy principle, and discusses various historical paradoxes, and how these pitfalls could have been avoided by not becoming addicted to measure theory, but by sticking to elementary discrete math. The historic battle between different schools is treated at (over)length in a rather polemic way. Gelman [GCSR95] is a modern and more practical book on Bayesian data analysis.

Alternatives to probability theory. Given the success story of Bayesian probability theory it is somewhat surprising that so many alternatives have been considered in AI. Many reasons why probability theory is unsuitable for AI have been stated: strict numerical values are not appropriate for a qualitative reasoning system, probability theory cannot deal with impreciseness, or vagueness, or subjective beliefs, or is just impractical. Setbacks due to naive and/or inconsistent application are also responsible for Bayesian probabilistic reasoning falling out of favor in the 1970s. Default reasoning [Rei80], nonmonotonic logic [MD80], and circumscription [McC80] treat conclusions or events not as “believed to a certain degree” but as “believed by default until a better reason is found to believe something else.” (see the anthology [Gin87]). Certainty factors (“fudge factors”) have been introduced into classical rule-based expert systems to accommodate uncertainty [Sho76, BS84]. Dempster-Shafer theory uses probability-intervals for probability values if they themselves are not perfectly known, usually because they have been estimated from a finite amount of data [Dem68, Sha76]. More generally, this approach goes under the name imprecise probabilities [Wal91]. In the full Bayesian treatment one defines a (second order) probability distribution over probability values to deal with this kind of ignorance or beliefs. Fuzzy logic deals with vaguely defined events (Fuzzy sets) which are only “sort of” true, like the “Eiffel tower is high” [Zad65, Zim91]. Possibility theory has been introduced to handle uncertainty in fuzzy systems [Zad78]. See [Fin73, Wal91] or [RN95, ch.15.6] for a more detailed account of various uncertain reasoning systems. Finally, quantum systems must be described with complex valued probability amplitudes resulting in strange interference effects. There may come a time (e.g. for nanobots) when quantum logic [Hug89] will be needed in AI.

All these alternate approaches have their problems: Either they have unclear semantics, or they are not selfconsistent, or they don’t scale up, or have other problems. It is not that Bayesian probability theory leaves no wishes open, but it is the most consistent system developed so far. Imprecise probability theory if interpreted as a probabilistic worst-case-reasoning approach is quite consistent and a possibly useful extension of probability theory in game theory and certain safety critical areas. The other approaches may survive as useful (efficient) approximations to a full Bayesian treatment. Although probability theory slowly (re)covers AI, the debate still goes on [Che85, Che88].

Cox's axioms and theorem. In [Cox46] Cox shows that every function $\text{Bel}(\cdot|\cdot)$ satisfying his Axioms 2.17 is isomorphic to a (conditional) probability function. This (with considerable delay) gave a significant boost in using standard probability theory for dealing with subjective beliefs and uncertainty. Cheeseman [Che88] has called Cox's derivation the "strongest argument for use of standard (Bayesian) probability theory". Similar sentiments are expressed by Jaynes [Jay78, p24]; indeed, Cox's Theorem is one of the cornerstones of Jaynes' recent book [Jay96]. Horvitz, Heckerman, and Langlotz [HHL86] used it as a basis for comparison of probability and other nonprobabilistic approaches to reasoning about uncertainty. Heckerman [Hec88] used it as a basis for providing an axiomatization for belief updates. Various variants of Cox's axioms have been considered in the literature [Rei49, Ačz66, Hec88, Jay96, Fin73, Tri69], which simplify the derivation, or weaken, replace or better motivate the assumption. A loophole in all these derivations have only recently been discovered [Par95]. They are all related to the following unwaveringness. The function F of Cox's axioms mapping $\text{Bel}(C|B \cap A)$ and $\text{Bel}(B|A)$ to $\text{Bel}(B \cap C|A)$ is proven to be associative, i.e. $F(x, F(y, z)) = F(F(x, y), z)$, but actually associativity is only proven for (x, y, z) of the form $x = \text{Bel}(D|C \cap B \cap A)$, $y = \text{Bel}(C|B \cap A)$, and $z = \text{Bel}(B|A)$ for some events A , B , C , and D . If the set of such triples (x, y, z) is dense in $[0, 1]^3$, then by continuity, F is associative. Paris provides a rigorous proof of Cox's result, assuming that the range of Bel is contained in $[0, 1]$ and using assumptions similar to [HHL86]. However, he and all others who tried to fix the proof needed to make additional assumptions that are not very appealing. Usually they demand or imply that the belief values are dense in a certain subset of \mathbb{R} , which excludes systems with a finite number of events. It remains an open question whether there is an appropriate strengthening of the assumptions that lead to Cox's result in finite settings. See [Hal99] and references therein for details.

Algorithmic probability & universal induction. The notion of (universal, enumerable) semimeasures has been introduced in [Sol64, ZL70, LV77]. Levin [ZL70] defined universal a priori probability as one dominating all enumerable semimeasures. The dominance (2.21) and the equivalence $M \stackrel{\times}{\approx} \xi_U$ is due to Levin [ZL70]. Convergence of M to μ in the conditional mean squared sense (2.23) is due to Solomonoff [Sol78] (who insists on normalizing M by giving up enumerability). The elementary proof of $M(x_t|x_{<t}) \xrightarrow{t \rightarrow \infty} 1$ for computable $x_{1:\infty}$ is not based on any source. The direct study of predictions based on past observations without discussing models has been coined *prequential approach* by Dawid [Daw84]. Good reviews on universal induction with a philosophical touch are [LV92b, Sol97]. For an older, but general review of inductive inference see Angluin [AS83].

Schmidhuber [Sch00, Sch02a] constructed a natural hierarchy of generalizations of algorithmic probability and complexity and introduced more general, approximable and universal cumulatively enumerable semimeasures. The restriction to time-bounded universal probability is treated in [LV91, LV97, Sch00, Sch02c, Sch02b] and is closely related to resource bounded complexity and universal Levin search.

Other topics related to universal induction are the Weighted Majority algorithm by Littlestone and Warmuth [LW94], universal forecasting by Vovk [Vov92], Levin search [Lev73b], pac-learning introduced by Valiant [Val84], the Minimum Message/Description Length principle [WB68, Ris78, Ris89, LV92a, Grü98, VL00], and Occam's razor, learnability and VC dimension [BEHW87, BEHW89].

Randomness of individual objects in terms of randomness tests has been defined by Martin-Löf [ML66, ML69] and is closely related to Kolmogorov complexity and algorithmic probability. Another interesting randomness criterion for individual sequences by Vovk in terms of the Hellinger distance can be found in [Vov87]. Randomness of individual sequences in a wider context are exhaustively analyzed in the survey papers [KU87, USS90].

Applications of Kolmogorov complexity and Levin's Universal Search.

Schmidhuber [Sch00, Sch02c] defines the *speed prior*, closely related to Levin search, and derives a *computable* strategy for optimal inductive reasoning. He analyzed consequences for computable universes sampled from such priors. Good numerical approximations to Kolmogorov complexity are computationally expensive. But the ongoing decrease of processing costs has permitted the first successful implementations and applications [Sch97, SZW97]. A derivate of Levin's universal search algorithm was used in [Sch97] to discover neural nets with low Levin complexity, low Kolmogorov complexity, and high generalization capability. Adaptive Levin Search (ALS) and the optimal ordered problem solver (OOPS) extends Levin search by making its underlying probability distribution on program space adaptive and by improving it according to experience [SZW97, Sch02b]. This can significantly speed up the discovery of algorithmic solutions.

There are numerous applications of MDL, which can be viewed as an applied form of Kolmogorov complexity [LV97]. Apart from that there are little "direct" approximations, implementations, or practical applications. Conte [Con97] evolves short Lisp programs to estimate Kolmogorov complexity. Chaitin [Cha91] speculates on the computational power of the evolutionary information gathering process and its relation to algorithmic information. Schmidt [Sch99] argues that (time-bounded) Kolmogorov complexity helps and not prevents the search for Extra Terrestrial Intelligence (SETI). Vovk [VW98] describes universal portfolio selection schemes.

2.6 Problems

2.1 ((Un)natural Turing machines) [C10] Show that for every string x there exists a universal Turing machine U' such that $K_{U'}(x) = 1$. Arguments of this sort are often used to demonstrate the arbitrariness/non-absolute character of algorithmic information. Argue that U' is not a natural Turing machine if x is complex. Elaborate on the difficulties in rigorously proving such statements.

2.2 (Exact ξ correspondence) [C20] We have shown that $M(x) := \sum_{p:U(p)=x^*} 2^{-l(p)}$ equals $\xi_U(x) := \sum_{\nu} 2^{-K(\nu)} \nu(x)$ within a multiplicative constant, i.e. $M \stackrel{\times}{\approx} \xi_U$. Improve this result to an exact equality in the sense that $M(x) = \xi_w(x) := \sum_{\nu \in \mathcal{M}_U} w_{\nu} \nu(x)$ for *some* weights with $w_{\nu} \geq 2^{-K(\nu)-O(1)}$ (solution due to Paul Vitányi, private communication). $M \stackrel{\times}{\approx} \xi_w$ is true for *any* choice of the weights $w_{\nu} > 0 \forall \nu \in \mathcal{M}_U$. Show that equality (within a constant) does no longer hold for a similarly generalized M with $2^{-l(p)}$ replaced by arbitrary $w_p > 0$.

2.3 (Martin-Löf random sequences) [C45oi] Show that a theorem true for all μ -random sequences (see Theorem 2.28) is also true with μ probability 1. Under what conditions is the reverse direction true? Especially is $\sum_{t=1}^{\infty} (\mu(x_t|x_{<t}) - M(x_t|x_{<t}))^2 < \infty$ true for every individual μ -random sequence? (cf. Theorem 3.19(i)). It has been shown that $M(x_t|x_{<t})/\mu(x_t|x_{<t}) \xrightarrow{t \rightarrow \infty} 1$ w. μ .p.1 (see [LV97, Th.5.2.2] and Theorem 3.19(v) and Problem 3.10). Does the stronger statement of convergence individually for all Martin-Löf μ -random sequences hold? The argument given in [LV97, Th.5.2.2] and [VL00, Th.10] is incomplete.⁶ The implication “ $M(x_{1:n}) \leq c \cdot \mu(x_{1:n}) \forall n \Rightarrow \lim_{n \rightarrow \infty} M(x_{1:n})/\mu(x_{1:n})$ exists” has been used, but not proven, and may indeed be wrong.

2.4 (Oracle properties of Kolmogorov complexity) [C20s] A function or problem A is said to be Turing-reducible to B if there exists a Turing machine (finitely) computing or solving A provided B is given as an oracle [HMu01]. Let $K: \mathbb{B}^* \rightarrow \mathbb{N}$ be the Kolmogorov complexity, $H: \mathbb{B}^* \rightarrow \mathbb{B}$ be the halting sequence ($H(p) = 1 \Leftrightarrow U(p)$ halts), and $\Omega := \sum_{p:U(p) \text{ halts}} 2^{-l(p)} \in \mathbb{R} \cong [\mathbb{N} \rightarrow \mathbb{B}]$ be the halting probability. Show that K , H , and Ω are Turing-reducible to each other (cf. [LV97, p175]).

2.5 (Weakly Forgetful environments) [C15u] Consider two sequences $x_{1:\infty}^1$ and $x_{1:\infty}^2$, “typical” in the sense that both are μ .M.L.random. Assume a different early history ($x_{<k}^1 \neq x_{<k}^2$, k fixed), continued by the same observations ($x_{k:n-1}^1 = x_{k:n-1}^2 = x_{k:n-1}$) for a long time n . Show that for computable μ the future is not affected by the far back history $x_{<k}^i$ in the sense that $\mu(x_n|x_{<k}^1 x_{k:n-1}) - \mu(x_n|x_{<k}^2 x_{k:n-1}) \rightarrow 0$ for $n \rightarrow \infty$. Hint: show $\xi(x_n|x_{k:n-1}) \rightarrow \mu(x_n|x_{<k}^i x_{k:n-1})$ for $i=1$ and $i=2$. This property of μ for “typical” sequences may be considered as a weak form of forgetfulness. Argue that it is more appropriate to define forgetfulness as asymptotic independence of $x_{<k}^1$ for *all* environments (cf. definition in Section 5.3.6). Suggestion: compare ergodic

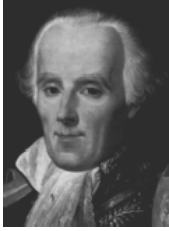
⁶The formulation of their Theorem is quite misleading in general: “Let μ be a positive recursive measure. If the length of y is fixed and the length of x grows to infinity, then $M(y|x)/\mu(y|x) \rightarrow 1$ with μ -probability one. The infinite sequences ω with prefixes x satisfying the displayed asymptotics are precisely [\Rightarrow] and [\Leftarrow] the μ -random sequences.” First, for off-sequence y convergence w.p.1 does not hold (xy must be demanded to be a prefix of ω). Second, the proof of [\Leftarrow] has loopholes (see main text). Last, [\Rightarrow] is given without proof and is probably wrong. Also the assertion in [LV97, Th.5.2.1] that $S_t := \mathbf{E} \sum_{x'_t} (\mu(x'_t|x_{<t}) - M(x'_t|x_{<t}))^2$ converges to zero faster than $1/t$ cannot be made, since S_t may not decrease monotonically.

Markov processes (see Definition 5.37) with μ defined by $\mu(1|x_{<n}) := i/3$ for $x_1 = i-1$ and $\mathcal{X} = \mathbb{B}$.

2.6 (Complexity increase) [C25u/C45oi] We are interested in good upper bounds on the increase in complexity when elongating a string $y := x_{<t}$ to $yx := x_{<t}x_{t:n} = x_{1:n}$. From Theorem 2.10(iv) we know that $K(yx) - K(y) \leq K(x|y) + O(1)$. Later (cf. Problem 3.13) we need similar bounds with K on the l.h.s. replaced by $KM(x) := -\log_2 M(x)$. Furthermore, let $C(x)$ be the plain Kolmogorov complexity, defined as the length of the shortest plain (as opposed to prefix) program computing x (see [LV97, Ch.2]). We have no particular demands on the r.h.s. of the inequality. So let us consider $\tilde{K}(x|y)$ defined as the length of the shortest plain *or* prefix, halting *or* non-halting program, computing x *or* a string starting with x , given y *or* a string starting with y . The only important property of $\tilde{K}(x|y)$ is that it corresponds to the length of a shortest program computing x from y . We don't want $\tilde{K}(x|y)$ to be *defined* as a difference $\tilde{K}(yx) - \tilde{K}(y)$. Prove the following inequalities:

- i) $C(yx) - C(y) \leq \tilde{K}(x|y) + O(?)$
- ii) $KM(yx) - KM(y) \leq \tilde{K}(x|y) + O(?)$
- iii) $KM(yx) - KM(y) \leq \tilde{K}(\mu|y) - \log_2 \mu(x|y) + O(?)$

Since C , K , \tilde{K} , and KM coincide within additive logarithmic terms, all inequalities follow from Theorem 2.10(iv) to logarithmic accuracy $O(\log l(xy))$. Improve the bounds to $O(?)^\pm \{K(C(y)), K(l(y)), K(l(y))\}$ respectively independent of x for suitable \tilde{K} . It is an open question whether the bounds hold within an additive constant independent of x *and* y for any of the \tilde{K} .



Pierre Laplace
(1749 – 1827)

“Probability theory is nothing but common sense reduced to calculation.” (Pierre Laplace, 1819)

Chapter 3

Universal Sequence Prediction

3.1	Introduction	303
3.1.1	Induction	303
3.1.2	Universal Sequence Prediction	304
3.2	Setup and Convergence	304
3.2.1	Random sequences	304
3.2.2	Universal Prior Probability Distribution	305
3.2.3	Universal Posterior Probability Distribution	306
3.2.4	Convergence of Random Sequences	307
3.2.5	Distance Measures between Probability Distributions	308
3.2.6	Convergence of ξ to μ	310
3.2.7	Convergence in Martin-Löf Sense	312
3.2.8	The case where $\mu \notin \mathcal{M}$	316
3.2.9	Probability Classes \mathcal{M}	317
3.3	Error Bounds	318
3.3.1	Bayes-Optimal Predictors	318
3.3.2	Total Expected Numbers of Errors	318
3.3.3	Proof of Theorem 3.36	320
3.4	Loss Bounds	321
3.4.1	Unit Loss Function	321
3.4.2	Loss Bound of Merhav & Feder	323
3.4.3	Example Loss Functions	324
3.4.4	Proof of Theorem 3.48	325
3.4.5	Convergence of Instantaneous Losses	326
3.4.6	General Loss	327

3.5	Application to Games of Chance	328
3.5.1	Introduction	328
3.5.2	Games of Chance	329
3.5.3	Example	330
3.5.4	Information-Theoretic Interpretation	330
3.6	Optimality Properties	331
3.6.1	Lower Error Bound	331
3.6.2	Pareto Optimality of ξ	334
3.6.3	Balanced Pareto Optimality of ξ	336
3.6.4	On the Optimal Choice of Weights	337
3.6.5	Occam's razor versus No Free Lunches	338
3.7	Miscellaneous	338
3.7.1	Multi-Step Predictions	338
3.7.2	Continuous Probability Classes \mathcal{M}	340
3.7.3	Further Applications	342
3.7.4	Prediction with Expert Advice	342
3.7.5	Outlook	344
3.8	Summary	345
3.9	Technical Proofs	345
3.9.1	How to Deal with $\mu=0$	346
3.9.2	Entropy Inequalities (3.11)	347
3.9.3	Error Inequality (3.36)	348
3.9.4	Binary Loss Inequality for $z \leq \frac{1}{2}$ (3.57)	349
3.9.5	Binary Loss Inequality for $z \geq \frac{1}{2}$ (3.58)	350
3.9.6	General Loss Inequality (3.53)	351
3.10	History & References	352
3.11	Problems	353

In this chapter we investigate Solomonoff's universal induction scheme in detail. More generally, we consider a universal (or mixture) distribution ξ , defined as a weighted sum or integral of distributions $\nu \in \mathcal{M}$, where \mathcal{M} is any countable or continuous set of distributions including μ . This is a generalization of Solomonoff induction, in which \mathcal{M} is the set of all enumerable semi-measures. We show for several performance measures that using the universal ξ as a prior is nearly as good as using the unknown true distribution μ . In a sense, this solves the problem of the unknown prior in a universal way. All results are obtained for general finite alphabet. Convergence of ξ to μ in a conditional mean squared sense and of $\xi/\mu \rightarrow 1$ with μ probability 1 is proven. The number of additional errors E_ξ made by the optimal universal prediction scheme based on ξ minus the number of errors E_μ of the optimal informed prediction scheme based on μ is proven to be bounded by $O(\sqrt{E_\mu})$. The prediction framework is generalized to arbitrary loss functions. A system is allowed to take an action y_t , given $x_1 \dots x_{t-1}$ and receives loss $\ell_{x_t y_t}$ if

x_t is the next symbol of the sequence. No assumptions on ℓ are necessary, besides boundedness. Optimal universal Λ_ξ and optimal informed Λ_μ prediction schemes are defined and the total loss of Λ_ξ is bounded in terms of the total loss of Λ_μ , similar to the error bounds. We show that the bounds are tight and that no other predictor can lead to significantly smaller bounds. Furthermore, for various performance measures we show Pareto-optimality of ξ in the sense that there is no other predictor which performs better or equal in all environments $\nu \in \mathcal{M}$ and strictly better in at least one. So, optimal predictors can (w.r.t. to most performance measures in expectation) be based on the mixture ξ . Finally we give an Occam's razor argument that the choice $w_\nu \sim 2^{-K(\nu)}$ for the weights is optimal, where $K(\nu)$ is the length of the shortest program describing ν . Furthermore, games of chance, defined as a sequence of bets, observations, and rewards are studied. The average profit achieved by the Λ_ξ scheme rapidly converges to the best possible profit. The time needed to reach the winning zone is proportional to the relative entropy of μ and ξ . The prediction schemes presented here are compared to predictors based on expert advice. Although the algorithms, the settings, and the proofs are quite different, the bounds of both schemes have a very similar structure. Extensions to infinite alphabets, partial, delayed and probabilistic prediction, classification, and more active systems are briefly discussed.

The **main new contributions** are to

- generalize the convergence [Sol78, LV97] of ξ to μ (Section 3.2),
- derive general error (Section 3.3) and loss (Section 3.4) bounds measuring the performance of ξ relative to μ , improving upon previous results [MF98],
- apply the results to games of chance (Section 3.5),
- show that the error/loss bounds are tight and that Solomonoff's universal prior is optimal (Section 3.6),
- generalize the bound in [CB90] on the relative entropy between ξ and μ for continuous i.i.d. probability classes \mathcal{M} to the non-i.i.d. case (Section 3.7.2),
- compare the universal prediction scheme and its loss bounds to predictors based on expert advice and its loss bounds [CB97] (Section 3.7.4).

3.1 Introduction

3.1.1 Induction

Many problems are of the induction type in which statements about the future have to be made, based on past observations. What is the probability of rain tomorrow, given the weather observations of the last few days? Is the Dow Jones likely to rise tomorrow, given the chart of the last years and possibly additional newspaper information? Can we reasonably doubt that the sun will rise tomorrow? Indeed, one definition of science is to predict the future, where, as an intermediate step, one tries to understand the past by developing theories and, as a consequence of

prediction, one tries to manipulate the future. All induction problems may be studied in the Bayesian framework. The probability of observing x_t at time t , given the observations $x_1 \dots x_{t-1}$ can be computed with the chain rule, if we know the true probability distribution, which generates the observed sequence $x_1 x_2 x_3 \dots$. The problem is that in many cases we do not even have a reasonable guess of the true distribution μ . What is the true probability of weather sequences, stock charts, or sunrises?

3.1.2 Universal Sequence Prediction

In order to overcome the problem of the unknown true distribution, one can define a mixture distribution ξ as a w_ν -weighted sum or integral over distributions $\nu \in \mathcal{M}$, where \mathcal{M} is any discrete or continuous (hypothesis) set including μ . \mathcal{M} is assumed to be known and to contain the true distribution, i.e. $\mu \in \mathcal{M}$. Since the probability ξ can be shown to converge rapidly to the true probability μ in a conditional sense, making decisions based on ξ is often nearly as good as the infeasible optimal decision based on the unknown μ [MF98]. Solomonoff [Sol64] had the idea to define a universal mixture M (see Section 2.4) as a weighted average over deterministic programs. Lower weights were assigned to longer programs. He unified Epicurus' principle of multiple explanations, Occam's razor [simplicity] principle into one formal theory (See [LV97] for this interpretation of [Sol64]). Inspired by Solomonoff's idea, Levin [ZL70] had the idea to define the closely related universal prior ξ_U as a weighted average over *all* semi-computable probability distributions. If the environment possesses some effective structure at all, Solomonoff's posterior [Sol78] "finds" this structure, and allows for a good prediction. In a sense, this solves the induction problem in a universal way, i.e. without making problem specific assumptions.

3.2 Setup and Convergence

3.2.1 Random sequences

We denote strings over a finite alphabet \mathcal{X} by $x_1 x_2 \dots x_n$ with $x_t \in \mathcal{X}$ and $t, n, N \in \mathbb{N}$ and $N = |\mathcal{X}|$. We further use the abbreviations ϵ for the empty string, $x_{t:n} := x_t x_{t+1} \dots x_{n-1} x_n$ for $t \leq n$ and ϵ for $t > n$, and $x_{<t} := x_1 \dots x_{t-1}$, and $\omega = x_{1:\infty}$ for infinite sequences. We use Greek letters for probability distributions (or measures). Let $\rho(x_1 \dots x_n)$ be the probability that an (infinite) sequence starts with $x_1 \dots x_n$:

$$\sum_{x_{1:n} \in \mathcal{X}^n} \rho(x_{1:n}) = 1, \quad \sum_{x_t \in \mathcal{X}} \rho(x_{1:t}) = \rho(x_{<t}), \quad \rho(\epsilon) = 1, \quad (3.1)$$

We also need conditional probabilities. Presuming they exist, we have

$$\rho(x_t | x_{<t}) = \rho(x_{1:t}) / \rho(x_{<t}), \quad (3.2)$$

$$\rho(x_1 \dots x_n) = \rho(x_1) \cdot \rho(x_2 | x_1) \cdot \dots \cdot \rho(x_n | x_1 \dots x_{n-1}), \quad (3.3)$$

called multiplication rule (of conditional probabilities), or chain rule. The first equation states that the probability that a string $x_1 \dots x_{t-1}$ is followed by x_t is equal to the probability that a string starts with $x_1 \dots x_t$ divided by the probability that a string starts with $x_1 \dots x_{t-1}$.

The second equation is the first, applied n times. Whereas ρ might be any probability distribution, μ denotes the true (unknown) generating distribution of the sequences. We denote expectations by \mathbf{E} . The (conditional) expected value of a function $f: \mathcal{X}^t \rightarrow \mathbb{R}$, dependent on $x_{1:t}$, independent of $x_{t+1:\infty}$, (given $x_{<t}$) is

$$\mathbf{E}[f] = \sum'_{x_{1:n} \in \mathcal{X}^n} \mu(x_{1:n}) f(x_{1:t}), \quad \mathbf{E}_t[f] := E[f(x_{1:t}) | x_{<t}] = \sum'_{x_t \in \mathcal{X}} \mu(x_t | x_{<t}) f(x_{1:t}) \quad (3.4)$$

for any choice of $n \geq t$. Expectations \mathbf{E} are *always* w.r.t. the true distribution μ . The prime denotes that the sum is restricted to $x_{1:n}$ with $\mu(x_{1:t}) \neq 0$ ($\mu(x_t | x_{<t}) \neq 0$). If $\mu(x_{<t}) = 0$, then $\mu(x_t | x_{<t})$ and hence \mathbf{E}_t is undefined. Since the sum in \mathbf{E} is restricted to $\mu(x_{1:n}) \neq 0$, $\mathbf{E}[\mathbf{E}_t[\cdot]] = \mathbf{E}[\cdot]$ is valid in any case (by the chain rule).

In a more probabilistic terminology we have a sample space $\Omega = \mathcal{X}^\infty$ with elements $\omega = \omega_1 \omega_2 \omega_3 \dots \in \Omega$ being infinite sequences over the finite alphabet \mathcal{X} . The cylinder sets $\Gamma_{x_{1:n}} := \{\omega : \omega_{1:n} = x_{1:n}\}$ are events. We define the σ -algebra \mathcal{F} as the smallest set containing all cylinder sets and which is closed under complement and countable union. A probability measure μ is uniquely defined by giving its values $\mu(\Gamma_{x_{1:n}})$ on the cylinder sets, which we abbreviate by $\mu(x_{1:n})$. See [LV97, Doo53] or any statistics book for a more thorough treatment.

Similarly, f may be interpreted as a random variable or measurable function. Two functions differing on a set of measure zero have the same expectation. So if we “undefine” f for some $x_{1:t}$ with $\mu(x_{1:t}) = 0$, the expectation should not be affected. Hence, \sum' is the correct definition for partial functions. The prime is ineffective and can be ignored for total functions. Many expressions in this work are undefined on a set of measure zero. Henceforth we will not mention this anymore. See Section 3.9.1 for alternative ways of treating $\mu = 0$.

Finally, the probability of an event $A \subseteq \Omega$ is $\mathbf{P}[A] = \mathbf{E}[\chi_A]$, where χ is the characteristic function of A , i.e. $\chi_A(\omega) = 1$ if $\omega \in A$, and $\chi(\omega) = 0$ otherwise.

3.2.2 Universal Prior Probability Distribution

Every inductive inference problem can be brought into the following form: Given a string $x_{<t}$, take a guess at its continuation x_t . We will assume that the strings which have to be continued are drawn from a probability¹ distribution μ . The maximal prior information a prediction algorithm can possess is the exact knowledge of μ , but in many cases (like for the probability of sun tomorrow) the true generating distribution is not known. Instead, the prediction is based on a guess ρ of μ . We expect that a predictor based on ρ performs well, if ρ is close to μ or converges,

¹This includes deterministic environments, in which case the probability distribution μ is 1 for some sequence $x_{1:\infty}$ and 0 for all others. We call probability distributions of this kind *deterministic*.

in a sense, to μ . Let $\mathcal{M} := \{\nu_1, \nu_2, \dots\}$ be a countable set of candidate probability distributions on strings. Results are generalized to continuous sets \mathcal{M} in Section 3.7.2. We define a weighted average on \mathcal{M}

$$\xi(x_{1:n}) \equiv \xi_{\mathcal{M}}(x_{1:n}) := \sum_{\nu \in \mathcal{M}} w_{\nu} \cdot \nu(x_{1:n}), \quad \sum_{\nu \in \mathcal{M}} w_{\nu} = 1, \quad w_{\nu} > 0. \quad (3.5)$$

It is easy to see that ξ is a probability distribution as the weights w_{ν} are positive and normalized to 1 and the $\nu \in \mathcal{M}$ are probabilities.² For finite \mathcal{M} a possible choice for the w is to give all ν equal weight ($w_{\nu} = \frac{1}{|\mathcal{M}|}$). We call ξ universal relative to \mathcal{M} , as it multiplicatively dominates all distributions in \mathcal{M}

$$\xi(x_{1:n}) \geq w_{\nu} \cdot \nu(x_{1:n}) \quad \text{for all } \nu \in \mathcal{M}. \quad (3.6)$$

In the following, we assume that \mathcal{M} is known and contains the true distribution, i.e. $\mu \in \mathcal{M}$. If \mathcal{M} is chosen sufficiently large, then $\mu \in \mathcal{M}$ is not a serious constraint. Generic classes, especially where \mathcal{M} contains *all* computable probability distributions, are discussed in Subsection 3.2.9. Generalizations to the case where \mathcal{M} does not contain μ are briefly discussed in Subsection 3.2.8. In the next Subsection we motivate and in Subsection 3.2.6 we show the important property of ξ converging to the true distribution $\mu \in \mathcal{M}$ in a sense and, hence, might being a useful substitute for the true, but in general, unknown distribution μ .

3.2.3 Universal Posterior Probability Distribution

All prediction schemes in this work are based on the conditional probabilities $\rho(x_t | x_{<t})$. It is possible to express also the conditional probability $\xi(x_t | x_{<t})$ as a weighted average over the conditional $\nu(x_t | x_{<t})$, but now with time dependent weights:

$$\xi(x_t | x_{<t}) = \sum_{\nu \in \mathcal{M}} w_{\nu}(x_{<t}) \nu(x_t | x_{<t}), \quad w_{\nu}(x_{1:t}) := w_{\nu}(x_{<t}) \frac{\nu(x_t | x_{<t})}{\xi(x_t | x_{<t})}, \quad w_{\nu}(\epsilon) := w_{\nu}. \quad (3.7)$$

The denominator just ensures correct normalization $\sum_{\nu} w_{\nu}(x_{1:t}) = 1$. By induction and the chain rule we see that $w_{\nu}(x_{<t}) = w_{\nu} \nu(x_{<t}) / \xi(x_{<t})$. Inserting this into $\sum_{\nu} w_{\nu}(x_{<t}) \nu(x_t | x_{<t})$ using (3.5) gives $\xi(x_t | x_{<t})$, which proves the equivalence of (3.5) and (3.7). If w_{ν} is interpreted as the prior (subjective) belief in ν , then $w_{\nu}(x_{<t})$ is the posterior belief in ν after having seen $x_{<t}$. The expressions (3.7) can be used to give an intuitive, but non-rigorous, argument why $\xi(x_t | x_{<t})$ converges to $\nu(x_t | x_{<t})$: The weight $w_{\nu}(x_{<t})$ of ν in ξ increases/decreases if ν assigns a high/low probability to the new symbol x_t , given $x_{<t}$. For a μ -random sequence $x_{1:t}$, $\mu(x_{1:t}) \gg \nu(x_{1:t})$ if ν (significantly) differs from μ . We expect the total weight for all ν consistent with

²The weight w_{ν} may be interpreted as the initial degree of belief in ν and $\xi(x_1 \dots x_n)$ as the degree of belief in $x_1 \dots x_n$. If the existence of true randomness is rejected on philosophical grounds one may consider \mathcal{M} containing only deterministic environments. ξ still represents belief probabilities.

μ to converges to 1, and all other weights converge to 0 for $t \rightarrow \infty$. Therefore we expect $\xi(x_t|x_{<t})$ to converge to $\mu(x_t|x_{<t})$ for μ -random strings $x_{1:n}$.

Expressions (3.7) seem to be more suitable than (3.5) for studying convergence and loss bounds of the universal predictor ξ , but it will turn out that (3.6) is all we need, with the sole exception in the proof of Theorem 3.66 and in Section 5.5. Probably (3.7) is useful when one tries to understand the learning aspect in ξ .

3.2.4 Convergence of Random Sequences

A classical (non-random) real-valued sequence a_t is defined to converge to a_* , short $a_t \rightarrow a_*$ if $\forall \varepsilon \exists t_0 \forall t \geq t_0 : |a_t - a_*| < \varepsilon$. We are interested in convergence properties of random sequences $z_t(\omega)$ for $t \rightarrow \infty$ (e.g. $z_t(\omega) = \xi(\omega_t|\omega_{<t}) - \mu(\omega_t|\omega_{<t})$). We define six convergence concepts for random sequences and relate them.

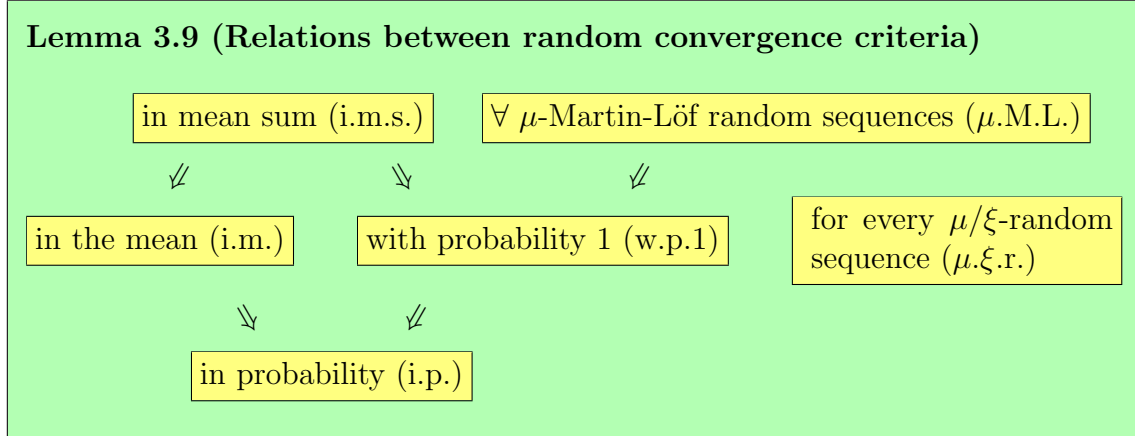
Definition 3.8 (Convergence of random sequences) Let $z_1(\omega), z_2(\omega), \dots$ be a sequence of real-valued random variables (and probability measure be μ). z_t is said to converge for $t \rightarrow \infty$ to (random variable) z_*

- i) with probability 1 (w.p.1) : $\Leftrightarrow \mathbf{P}[\{\omega : z_t(\omega) \rightarrow z_*(\omega)\}] = 1$
 $\Leftrightarrow \forall \varepsilon > 0 : \mathbf{P}[\sup_{s \geq t} |z_s - z_*| \geq \varepsilon] \rightarrow 0$ for $t \rightarrow \infty$,
- ii) in the mean (i.m.) : $\Leftrightarrow \mathbf{E}[(z_t - z_*)^2] \rightarrow 0$ for $t \rightarrow \infty$,
- iii) in mean sum (i.m.s.) : $\Leftrightarrow \sum_{t=1}^{\infty} \mathbf{E}[(z_t - z_*)^2] < \infty$,
- iv) in probability (i.p.) : $\Leftrightarrow \forall \varepsilon > 0 : \mathbf{P}[|z_t - z_*| \geq \varepsilon] \rightarrow 0$ for $t \rightarrow \infty$,
- v) for every μ -Martin-Löf random sequence (μ .M.L.) : \Leftrightarrow
 $\forall \omega : [\exists c \forall n : \xi_U(\omega_{1:n}) \leq c\mu(\omega_{1:n})]$ implies $z_t(\omega) \rightarrow z_*(\omega)$ for $t \rightarrow \infty$,
- vi) for every μ/ξ -random sequence (μ . ξ .r.) : \Leftrightarrow
 $\forall \omega : [\exists c \forall n : \xi(\omega_{1:n}) \leq c\mu(\omega_{1:n})]$ implies $z_t(\omega) \rightarrow z_*(\omega)$ for $t \rightarrow \infty$.

See Section 2.4 for a definition of $\xi_U \equiv \xi_{\mathcal{M}_{\text{enum}}^{\text{semi}}} \stackrel{\times}{=} M$. In statistics, (i) is the “default” characterization of convergence of random sequences. (ii), (iii), and (iv) are also well-known and often more convenient to deal with than (i). Further, convergence i.m.s. is very strong: it provides a “rate” of convergence in the sense that the expected number of times t in which z_t deviates more than ε from z_* is finite and bounded by $\sum_{t=1}^{\infty} \mathbf{E}[(z_t - z_*)^2] / \varepsilon^2$. (v) uses Martin-Löf’s notion of randomness of *individual* sequences to define convergence M.L. Since this Chapter mainly deals with general Bayes-mixtures ξ , we generalized in (vi) the definition of convergence M.L. based on ξ_U to convergence μ . ξ .r. based on ξ in a natural way.³ Convergence

³For finite \mathcal{M} or if $\xi \in \mathcal{M}$, the definition of μ/ξ -randomness depends only on \mathcal{M} , and not on the specific weights used in ξ .

in one sense often implies convergence in another. The following implications for convergence of random sequences are true. Unconnected criteria (in the transitive hull) are incomparable (see also Problem 3.9).



3.2.5 Distance Measures between Probability Distributions

We need several distance measures between vectors $\mathbf{y} = (y_i)$ and $\mathbf{z} = (z_i)$ in general, and probability distributions/vectors for which $y_i \geq 0$, $z_i \geq 0$, and $\sum_i y_i = \sum_i z_i \leq 1$ in particular, $i = 1, \dots, N$, namely the⁴

$$\begin{aligned}
 \text{absolute (or Manhattan) distance} & : a(\mathbf{y}, \mathbf{z}) := \sum_i |y_i - z_i| \\
 \text{quadratic (or squared Euclidian) distance} & : s(\mathbf{y}, \mathbf{z}) := \sum_i (y_i - z_i)^2 \\
 \text{Hellinger distance} & : h(\mathbf{y}, \mathbf{z}) := \sum_i (\sqrt{y_i} - \sqrt{z_i})^2 \quad (3.10) \\
 \text{relative entropy or Kullback-Leibler divergence} & : d(\mathbf{y}, \mathbf{z}) := \sum_i y_i \ln \frac{y_i}{z_i} \\
 \text{absolute divergence} & : b(\mathbf{y}, \mathbf{z}) := \sum_i y_i \left| \ln \frac{y_i}{z_i} \right|
 \end{aligned}$$

The relative entropy is not a metric, but for probability distributions, for which it is defined, it is at least non-negative and zero if and only if $\mathbf{y} = \mathbf{z}$. All bounds we prove in this chapter heavily rely on the following inequalities.

⁴ $0 \ln \frac{0}{z} := 0 \forall z \geq 0$ and $y \ln \frac{y}{0} := \infty \forall y > 0$.

Lemma 3.11 (Entropy Inequalities) Let $\{y_i\}$ and $\{z_i\}$ be two probability distributions, i.e. $y_i \geq 0$, $z_i \geq 0$, and $\sum_i y_i = \sum_i z_i = 1$, and f is a convex and even ($f(x) = f(-x)$) function with $f(0) \leq 0$ then the following inequalities hold:⁵

$$\begin{aligned} \frac{1}{2}\Sigma f &\leq f\sqrt{\frac{1}{2}d} : & \frac{1}{2}\sum_i f(y_i - z_i) &\stackrel{(f)}{\leq} f\left(\sqrt{\frac{1}{2}\sum_i y_i \ln \frac{y_i}{z_i}}\right) \\ s \leq d : & \sum_i (y_i - z_i)^2 &\stackrel{(s)}{\leq} \sum_i y_i \ln \frac{y_i}{z_i} \\ b-d \leq a \leq \sqrt{2d} : & \sum_i y_i \left| \ln \frac{y_i}{z_i} \right| - \sum_i y_i \ln \frac{y_i}{z_i} &\stackrel{(b)}{\leq} \sum_i |y_i - z_i| \stackrel{(a)}{\leq} \sqrt{2\sum_i y_i \ln \frac{y_i}{z_i}} \\ h \leq d : & \sum_i (\sqrt{y_i} - \sqrt{z_i})^2 &\stackrel{(h)}{\leq} \sum_i y_i \ln \frac{y_i}{z_i} \end{aligned}$$

(a), (b), and (h) are known results, (f) and (s) are potentially new. A proof of the lemma is deferred to Section 3.9. Inequality (3.11s) is a generalization of the binary $N=2$ case used in [Sol78, Hut01c, LV97]. If we insert

$$\mathcal{X} = \{1, \dots, N\}, \quad N = |\mathcal{X}|, \quad i = x_t, \quad y_i = \mu(x_t | x_{<t}), \quad z_i = \xi(x_t | x_{<t}) \quad (3.12)$$

into (3.10) we get various *instantaneous distances* (at time t) between μ and ξ . If we take the expectation \mathbf{E} and sum over $\sum_{t=1}^n$ we get various *total distances* between μ and ξ :

$$a_t(x_{<t}) := \sum_{x_t} |\mu(x_t | x_{<t}) - \xi(x_t | x_{<t})|, \quad A_n := \sum_{t=1}^n \mathbf{E}[a_t(x_{<t})] \quad (3.13)$$

$$s_t(x_{<t}) := \sum_{x_t} \left(\mu(x_t | x_{<t}) - \xi(x_t | x_{<t}) \right)^2, \quad S_n := \sum_{t=1}^n \mathbf{E}[s_t(x_{<t})] \quad (3.14)$$

$$h_t(x_{<t}) := \sum_{x_t} \left(\sqrt{\mu(x_t | x_{<t})} - \sqrt{\xi(x_t | x_{<t})} \right)^2, \quad H_n := \sum_{t=1}^n \mathbf{E}[h_t(x_{<t})] \quad (3.15)$$

$$d_t(x_{<t}) := \sum_{x_t} \mu(x_t | x_{<t}) \ln \frac{\mu(x_t | x_{<t})}{\xi(x_t | x_{<t})}, \quad D_n := \sum_{t=1}^n \mathbf{E}[d_t(x_{<t})] \quad (3.16)$$

$$b_t(x_{<t}) := \sum_{x_t} \mu(x_t | x_{<t}) \left| \ln \frac{\mu(x_t | x_{<t})}{\xi(x_t | x_{<t})} \right|, \quad B_n := \sum_{t=1}^n \mathbf{E}[b_t(x_{<t})] \quad (3.17)$$

For D_n the following can be shown [Sol78, LV97]

$$D_n = \sum_{t=1}^n \mathbf{E}[d_t(x_{<t})] = \sum_{t=1}^n \mathbf{E}[\mathbf{E}_t[\ln \frac{\mu(x_t | x_{<t})}{\xi(x_t | x_{<t})}]] = \quad (3.18)$$

⁵In (b) if some $z=0$ we define $y|\ln \frac{y}{0}| - y\ln \frac{y}{0} := 0$.

$$\mathbf{E}[\ln \prod_{t=1}^n \frac{\mu(x_t|x_{<t})}{\xi(x_t|x_{<t})}] = \mathbf{E}[\ln \frac{\mu(x_{1:n})}{\xi(x_{1:n})}] \leq \ln w_\mu^{-1} =: b_\mu$$

In the first line we have inserted (3.16) and used the definition of \mathbf{E}_t . Using $\mathbf{E}[\mathbf{E}_t[.]] = \mathbf{E}[.]$, the t sum can thereafter be exchanged with the expectation \mathbf{E} and transforms to a product inside the logarithm. In the last equality we have used the second form of the chain rule (3.3) for μ and ξ . Using universality (3.6) of ξ , i.e. $\ln[\mu(x_{1:n})/\xi(x_{1:n})] \leq \ln w_\mu^{-1}$ for $\mu \in \mathcal{M}$ yields the final inequality in (3.18).

3.2.6 Convergence of ξ to μ

Theorem 3.19 (Convergence of ξ to μ) Let there be sequences $x_1 x_2 \dots$ over a finite alphabet \mathcal{X} drawn with probability $\mu(x_{1:n})$ for the first n symbols. The universal posterior probability $\xi(x_t|x_{<t})$ of the next symbol x_t given $x_{<t}$ is related to the true posterior probability $\mu(x_t|x_{<t})$ in the following way:

- i) $\sum_{t=1}^n \mathbf{E} \left[\sum_{x'_t} \left(\mu(x'_t|x_{<t}) - \xi(x'_t|x_{<t}) \right)^2 \right] \equiv S_n \leq D_n \leq \ln w_\mu^{-1} < \infty$
- ii) $\sum_{x'_t} \left(\mu(x'_t|x_{<t}) - \xi(x'_t|x_{<t}) \right)^2 \equiv s_t(x_{<t}) \leq d_t(x_{<t}) \rightarrow 0$ for $t \rightarrow \infty$ w. μ .p.1
- iii) $\xi(x'_t|x_{<t}) - \mu(x'_t|x_{<t}) \rightarrow 0$ for $t \rightarrow \infty$ w. μ .p.1 (and i.m.s) for any x'_t
- iv) $\sum_{t=1}^n \mathbf{E} \left[\left(\sqrt{\frac{\xi(x_t|x_{<t})}{\mu(x_t|x_{<t})}} - 1 \right)^2 \right] \leq H_n \leq D_n \leq \ln w_\mu^{-1} < \infty$
- v) $\sqrt{\frac{\xi(x_t|x_{<t})}{\mu(x_t|x_{<t})}} \rightarrow 1$ i.m.s and $\frac{\xi(x_t|x_{<t})}{\mu(x_t|x_{<t})} \rightarrow 1$ w. μ .p.1 for $t \rightarrow \infty$
- vi) $b_t(x_{<t}) - d_t(x_{<t}) \leq a_t(x_{<t}) \leq \sqrt{2d_t(x_{<t})}, \quad B_n - D_n \leq A_n \leq \sqrt{2nD_n},$

where d_t and D_n are the relative entropies (3.16), w_μ is the weight (3.5) of μ in ξ , and $x'_{1:\infty}$ an arbitrary (non-random) sequence.

Proof. Inequality (ii) follows from the definitions (3.14) and (3.16) and from the entropy inequality (3.11s). From the definition and finiteness of D_∞ (3.18), and from $d_t(x_{<t}) \geq 0$ one sees that $\sqrt{d_t(x_{<t})} \rightarrow 0$ for $t \rightarrow \infty$ i.m.s., which implies $d_t(x_{<t}) \rightarrow 0$ w. μ .p.1. The inequality (i) follows from (ii) by taking the \mathbf{E} expectation and the $\sum_{t=1}^n$ sum. (iii) follows from (i) by dropping $\sum_{x'_t}$. The reason for the astonishing property of a single (universal) function ξ to converge to *any* $\mu \in \mathcal{M}$ lies in the fact that the sets of μ -random sequences differ for different μ . (iv) and (v) are related

to (i) and (iii), but are incomparable convergence results. To prove (iv) we use the abbreviations $\mu_t = \mu(x_t|x_{<t})$ and $\xi_t = \xi(x_t|x_{<t})$. For $\mu(x_{<t}) \neq 0$ we have

$$\mathbf{E}_t \left[\left(\sqrt{\frac{\xi_t}{\mu_t}} - 1 \right)^2 \right] = \sum'_{x_t} \mu(x_t|x_{<t}) \left(\sqrt{\frac{\xi_t}{\mu_t}} - 1 \right)^2 = \sum'_{x_t} (\sqrt{\xi_t} - \sqrt{\mu_t})^2 \leq h_t(x_{<t}) \leq d_t(x_{<t}). \quad (3.20)$$

The two inequalities follow from (3.15) and (3.11h). (iv) now follows by taking the \mathbf{E} expectation and the $\sum_{t=1}^n$ sum. (v) follows from (iv) by the Definition 3.8(iii) of convergence i.m.s., which implies convergence w. μ .p.1. The first two inequalities in (vi) immediately follow from inequalities (3.11a,b) and Definitions (3.13), (3.16) and (3.17). The third inequality follows from the first by linearity of \mathbf{E} and \sum . The last inequality $A_n \leq \sqrt{2nD_n}$ follows from

$$\frac{1}{n} A_n \equiv \frac{1}{n} \sum_{t=1}^n \mathbf{E}[a_t] \leq \frac{1}{n} \sum_{t=1}^n \mathbf{E}[\sqrt{2d_t}] \leq \frac{1}{n} \sum_{t=1}^n \sqrt{\mathbf{E}[2d_t]} \leq \sqrt{\frac{1}{n} \sum_{t=1}^n \mathbf{E}[2d_t]} \equiv \sqrt{\frac{2}{n} D_n} \quad (3.21)$$

where we have used Jensen's inequality for exchanging the averages ($\frac{1}{n} \sum_{t=1}^n$ and $\mathbf{E}_{<t}$) with the concave function $\sqrt{\cdot}$. \square

Since the conditional probabilities are the basis of all prediction algorithms considered in this work and ξ converges rapidly to μ (see Problem 3.11), we expect a good prediction performance if we use ξ as a guess of μ . Performance measures are defined in the following sections.

(i)–(iii) generalize Solomonoff's result [Sol78] to arbitrary finite alphabet. Without the use of the Hellinger distance, a somewhat weaker statement than (v) can be derived from (vi):

$$\mathbf{E} \left| \ln \frac{\mu(x_t|x_{<t})}{\xi(x_t|x_{<t})} \right| = \mathbf{E}[b_t] \leq \mathbf{E}[d_t] + \mathbf{E}[\sqrt{2d_t}] \leq \mathbf{E}[d_t] + \sqrt{2\mathbf{E}[d_t]} \xrightarrow{t \rightarrow \infty} 0,$$

since $\mathbf{E}[d_t] \rightarrow 0$. I.e. $|\ln \frac{\mu(x_t|x_{<t})}{\xi(x_t|x_{<t})}|^{1/2} \rightarrow 0$ i.m., which implies $\xi(x_t|x_{<t})/\mu(x_t|x_{<t}) \rightarrow 1$ i.p. The explicit appearance of n in the last expression of (vi) prevents proving stronger convergence w. μ .p.1 from (vi).

The elementary proof for (v) w. μ .p.1 given here does not rely on the semimartingale convergence Theorem [Doo53, pp. 324–325] as the proof of Gács in [LV97, Th.5.2.2]. Furthermore, (iv) gives the “speed” of convergence. Note the subtle difference between (iii) and (v). For *any* sequence $x'_{1:\infty}$ (possibly constant and not necessarily μ -random), $\mu(x'_t|x_{<t}) - \xi(x'_t|x_{<t})$ converges to zero w.p.1 (referring to $x_{1:\infty}$), but no statement is possible for $\xi(x'_t|x_{<t})/\mu(x'_t|x_{<t})$, since $\liminf \mu(x'_t|x_{<t})$ could be zero. On the other hand, if we stay *on* the μ -random sequence ($x'_{1:\infty} = x_{1:\infty}$), (v) shows that $\xi(x_t|x_{<t})/\mu(x_t|x_{<t}) \rightarrow 1$ (whether $\inf \mu(x_t|x_{<t})$ tends to zero or not does not matter). Indeed, it is easy to give an example where $\xi(x'_t|x_{<t})/\mu(x'_t|x_{<t})$ diverges. If we choose

$$\mathcal{M} = \{\mu_1, \mu_2\}, \quad \mu \equiv \mu_1, \quad \mu_1(1|x_{<t}) = \frac{1}{2}t^{-3} \quad \text{and} \quad \mu_2(1|x_{<t}) = \frac{1}{2}t^{-2}$$

the contribution of μ_2 to ξ causes ξ to fall off like $\mu_2 \sim t^{-2}$, much slower than $\mu \sim t^{-3}$ causing the quotient to diverge:

$$\begin{aligned}
\mu_1(0_{1:n}) &= \prod_{t=1}^n (1 - \tfrac{1}{2}t^{-3}) \xrightarrow{n \rightarrow \infty} c_1 = 0.450\dots > 0 \Rightarrow 0_{1:\infty} \text{ is a } \mu\text{-random sequence,} \\
\mu_2(0_{1:n}) &= \prod_{t=1}^n (1 - \tfrac{1}{2}t^{-2}) \xrightarrow{n \rightarrow \infty} c_2 = 0.358\dots > 0 \Rightarrow \xi(0_{1:n}) \rightarrow w_1 c_1 + w_2 c_2 =: c_\xi > 0 \\
\xi(0_{<t}1) &= w_1 \mu_1(1|0_{<t}) \mu_1(0_{<t}) + w_2 \mu_2(1|0_{<t}) \mu_2(0_{<t}) \rightarrow \tfrac{1}{2} w_2 c_2 t^{-2} \\
\Rightarrow \xi(1|0_{<t}) &= \frac{\xi(0_{<t}1)}{\xi(0_{<t})} \rightarrow \frac{w_2 c_2}{2 c_\xi} t^{-2} \Rightarrow \frac{\xi(1|0_{<t})}{\mu(1|0_{<t})} \rightarrow \frac{w_2 c_2}{c_\xi} t \rightarrow \infty \text{ diverges.}
\end{aligned}$$

3.2.7 Convergence in Martin-Löf Sense

An interesting open question is whether ξ converges to μ (in the sense of (iii) or (v)) individually for all Martin-Löf random sequences, short $\xi_U \xrightarrow{\text{M.L.}} \mu$ (see Problem 3.10). Clearly, convergence μ .M.L. implies (iii) and (v) by Lemma 3.9, but the converse may fail on a set of sequences with μ -measure zero. A convergence M.L. result would be particularly interesting and natural for the universal prior ξ_U , since M.L. randomness can be defined in terms of $\xi_U \stackrel{\times}{=} M$. Attempts to convert (ii) or (iv) to effective μ .M.L. randomness tests fail, since $\xi_U(x_t|x_{<t})$, and hence (ii) and (iv) are not enumerable. The argument given in [LV97, Th.5.2.2] and [VL00, Th.10] is incomplete. The implication “ $\xi_U(x_{1:n}) \leq c \cdot \mu(x_{1:n}) \forall n \Rightarrow \lim_{n \rightarrow \infty} \xi_U(x_{1:n})/\mu(x_{1:n})$ exists” has been used, but not proven, and may indeed be wrong (cf. Problem 2.3). Vovk [Vov87] shows that for two finitely computable semi-measures μ and ρ and $x_{1:\infty}$ being μ and ρ M.L. random that

$$\sum_{t=1}^{\infty} \sum_{x'_t} \left(\sqrt{\mu(x'_t|x_{<t})} - \sqrt{\rho(x'_t|x_{<t})} \right)^2 < \infty \quad \text{and} \quad \sum_{t=1}^{\infty} \left(\frac{\rho(x_t|x_{<t})}{\mu(x_t|x_{<t})} - 1 \right)^2 < \infty.$$

If ξ_U were recursive, then $\xi_U \rightarrow \mu$ and $\xi_U/\mu \rightarrow 1$ for every μ .M.L. random sequence $x_{1:\infty}$, since *every* sequence is ξ_U .M.L. random. Since ξ_U is *not* recursive Vovk’s theorem cannot be applied and it is not obvious how to generalize it. So the question of individual convergence remains open. More generally, one may ask whether $\xi_{\mathcal{M}} \rightarrow \mu$ for every \mathcal{M} random sequence. It turns out that this is true for some \mathcal{M} , but false for others.

Theorem 3.22 (μ/ξ -convergence of ξ to μ) Let $\mathcal{X} = \mathbb{B}$ be binary and $\mathcal{M}_\Theta := \{\mu_\theta : \mu_\theta(1|x_{<t}) = \theta \forall t, \theta \in \Theta\}$ be the set of Bernoulli(θ) distributions with parameters $\theta \in \Theta$. Let Θ_D be a countable dense subset of $[0,1]$, e.g. $[0,1] \cap \mathbb{Q}$ and Θ_G be a closed countable subset of $(0,1)$ with a gap, e.g. $\{\frac{1}{4}, \frac{1}{2}\}$ or $\mathbb{Q} \cap ([\frac{1}{5}, \frac{2}{5}] \cup [\frac{3}{5}, \frac{4}{5}])$. Then

- i) If $x_{1:\infty}$ is $\mu/\xi_{\mathcal{M}_{\Theta_D}}$ random with $\mu \in \mathcal{M}_{\Theta_D}$, then $\xi_{\mathcal{M}_{\Theta_D}}(x_t|x_{<t}) \rightarrow \mu(x_t|x_{<t})$,
- ii) There are $\mu \in \mathcal{M}_{\Theta_G}$ and $\mu/\xi_{\mathcal{M}_{\Theta_G}}$ random $x_{1:\infty}$ for which $\xi_{\mathcal{M}_{\Theta_G}}(x_t|x_{<t}) \not\rightarrow \mu(x_t|x_{<t})$.

Our original/main motivation of studying μ/ξ -randomness is the implication of Theorem 3.22 that $\xi_U \xrightarrow{\text{M.L.}} \mu$ cannot be decided from ξ being a mixture distribution (3.5) or from the dominance property (3.6) alone. Further structural properties of \mathcal{M}_U have to be employed. For Bernoulli sequences, convergence $\mu.\xi_{\mathcal{M}_\Theta}.x.$ is related to denseness of \mathcal{M}_Θ . Maybe a denseness characterization of $\mathcal{M}_{\text{enum}}^{\text{semi}}$ can solve the question of convergence M.L. of ξ_U . The property $\xi_U \in \mathcal{M}_U$ is also not sufficient to resolve the question, since there are $\mathcal{M} \ni \xi$ for which $\xi \xrightarrow{\mu.\xi.x} \mu$ and $\mathcal{M} \ni \xi$ for which $\xi \not\xrightarrow{\mu.\xi.x} \mu$ (see also Problem 3.10). Theorem 3.22 can be generalized to i.i.d. sequences over general finite alphabet \mathcal{X} .

The idea to prove (ii) is to construct a sequence $x_{1:\infty}$ which is μ_{θ_0}/ξ -random and μ_{θ_1}/ξ -random for $\theta_0 \neq \theta_1$. This is possible if and only if Θ contains a gap and θ_0 and θ_1 are the boundaries of the gap. Obviously ξ cannot converge to θ_0 and θ_1 , thus proving non-convergence. For no $\theta \in [0,1]$ will this $x_{1:\infty}$ be μ_θ M.L.-random. Finally, the proof of Theorem 3.22 makes essential use of the mixture representation of ξ , as opposed to the proof of Theorem 3.19 which only needs dominance $\xi \stackrel{\times}{\succ} \mathcal{M}$.

An example for (ii) is $\mathcal{M} = \{\mu_0, \mu_1\}$, $\mu_0(1|x_{<t}) = \mu_1(0|x_{<t}) = \frac{1}{4}$, $x_{1:\infty} = (01)^\infty = 01010101\dots \Rightarrow \mu_0(x_{1:2n}) = \mu_1(x_{1:2n}) = \xi(x_{1:2n}) = (\frac{1}{4})^n (\frac{3}{4})^n \Rightarrow x_{1:\infty}$ is μ_0/ξ -random and μ_1/ξ -random, but $\mu_0(x_{2n}|x_{<2n}) = \frac{1}{4}$, $\mu_0(x_{2n+1}|x_{1:2n}) = \frac{3}{4}$, $\mu_1(x_{2n}|x_{<2n}) = \frac{3}{4}$, $\mu_1(x_{2n+1}|x_{1:2n}) = \frac{1}{4}$ and $\xi(x_{2n}|x_{<2n}) = \frac{3}{8}$, $\xi(x_{2n+1}|x_{1:2n}) = \frac{1}{2}$ for $w_0 = w_1 = \frac{1}{2} \Rightarrow \xi(x_n|x_{<n}) \not\rightarrow \mu_{0/1}(x_n|x_{<n})$.

Proof. Let $\mathcal{X} = \mathbb{B}$ and $\mathcal{M} = \{\mu_\theta : \theta \in \Theta\}$ with countable $\Theta \subset [0,1]$ and $\mu_\theta(1|x_{1:n}) = \theta = 1 - \mu_\theta(0|x_{1:n})$, which implies

$$\mu_\theta(x_{1:n}) = \theta^{n_1} (1 - \theta)^{n - n_1}, \quad n_1 := x_1 + \dots + x_n, \quad \hat{\theta} \equiv \hat{\theta}_n := \frac{n_1}{n}$$

$\hat{\theta}$ depends on n ; all other used/defined θ will be independent of n . We assume $\theta. \in \Theta$, where $..$ stands for some (possible empty) index, and $\ddot{\theta} \in [0,1]$ (possibly $\notin \Theta$), where $\ddot{\theta}$ stands for some superscript, i.e. $\mu_{\ddot{\theta}}$ and $w_{\ddot{\theta}}$ make sense, whereas $\mu_{\hat{\theta}}$ and $w_{\hat{\theta}}$ do not. ξ is defined in the standard way as

$$\xi(x_{1:n}) = \sum_{\theta \in \Theta} w_\theta \mu_\theta(x_{1:n}) \quad \Rightarrow \quad \xi(x_{1:n}) \geq w_\theta \mu_\theta(x_{1:n}), \quad (3.23)$$

where $\sum_{\theta} w_{\theta} = 1$ and $w_{\theta} > 0 \forall \theta$. In the following let $\mu = \mu_{\theta_0} \in \mathcal{M}$ be the true environment.

$$\omega = x_{1:\infty} \text{ is } \mu/\xi\text{-random} \Leftrightarrow \exists c_{\omega} : \xi(x_{1:n}) \leq c_{\omega} \cdot \mu_{\theta_0}(x_{1:n}) \forall n \quad (3.24)$$

For binary alphabet it is sufficient to establish whether $\xi(1|x_{1:n}) \xrightarrow{n \rightarrow \infty} \theta_0 \equiv \mu(1|x_{1:n})$ for μ/ξ -random $x_{1:\infty}$ in order to decide $\xi(x_n|x_{<n}) \rightarrow \mu(x_n|x_{<n})$. We need the following posterior representation of ξ :

$$\xi(1|x_{1:n}) = \sum_{\theta \in \Theta} w_n^{\theta} \mu_{\theta}(1|x_{1:n}), \quad w_n^{\theta} := w_{\theta} \frac{\mu_{\theta}(x_{1:n})}{\xi(x_{1:n})} \leq \frac{w_{\theta}}{w_{\theta_0}} \frac{\mu_{\theta}(x_{1:n})}{\mu_{\theta_0}(x_{1:n})}, \quad \sum_{\theta \in \Theta} w_n^{\theta} = 1 \quad (3.25)$$

The ratio $\mu_{\theta}/\mu_{\theta_0}$ can be represented as follows:

$$\frac{\mu_{\theta}(x_{1:n})}{\mu_{\theta_0}(x_{1:n})} = \frac{\theta^{n_1}(1-\theta)^{n-n_1}}{\theta_0^{n_1}(1-\theta_0)^{n-n_1}} = \left[\left(\frac{\theta}{\theta_0} \right)^{\hat{\theta}_n} \left(\frac{1-\theta}{1-\theta_0} \right)^{1-\hat{\theta}_n} \right]^n = e^{n[D(\hat{\theta}_n||\theta_0) - D(\hat{\theta}_n||\theta)]} \quad (3.26)$$

$$\text{where } D(\hat{\theta}||\theta) = \hat{\theta} \ln \frac{\hat{\theta}}{\theta} + (1-\hat{\theta}) \ln \frac{1-\hat{\theta}}{1-\theta}$$

is the relative entropy between $\hat{\theta}$ and θ , which is continuous in $\hat{\theta}$ and θ , and is 0 if and only if $\hat{\theta} = \theta$. We also need the following implication for sets $\Omega \subseteq \Theta$:

$$\text{If } w_n^{\theta} \leq w_{\theta} g_{\theta}(n) \xrightarrow{n \rightarrow \infty} 0 \text{ and } g_{\theta}(n) \leq c \forall \theta \in \Omega, \text{ then } \sum_{\theta \in \Omega} w_n^{\theta} \mu_{\theta}(1|x_{1:n}) \leq \sum_{\theta \in \Omega} w_n^{\theta} \xrightarrow{n \rightarrow \infty} 0, \quad (3.27)$$

which follows from boundedness $\sum_{\theta} w_n^{\theta} \leq 1$ and $\mu_{\theta} \leq 1$. We now prove Theorem 3.22. We leave the special considerations necessary when $0, 1 \in \Theta$ to the reader and assume, henceforth, $0, 1 \notin \Theta$.

(i) Let Θ be a countable dense subset of $(0, 1)$ and $x_{1:\infty}$ be μ/ξ -random. Using (3.23) and (3.24) in (3.26) for $\theta \in \Theta$ to be determined later we can bound

$$e^{n[D(\hat{\theta}_n||\theta_0) - D(\hat{\theta}_n||\theta)]} = \frac{\mu_{\theta}(x_{1:n})}{\mu_{\theta_0}(x_{1:n})} \leq \frac{c_{\omega}}{w_{\theta}} =: c < \infty \quad (3.28)$$

Let us assume that $\hat{\theta} \equiv \hat{\theta}_n \not\rightarrow \theta_0$. This implies that there exists a cluster point $\tilde{\theta} \neq \theta_0$ of sequence $\hat{\theta}_n$, i.e. $\hat{\theta}_n$ is infinitely often in an ε -neighborhood of $\tilde{\theta}$, e.g. $D(\hat{\theta}_n||\tilde{\theta}) \leq \varepsilon$ for infinitely many n . $\tilde{\theta} \in [0, 1]$ may be outside Θ . Since $\tilde{\theta} \neq \theta_0$ this implies that $\hat{\theta}_n$ must be “far” away from θ_0 infinitely often. E.g. for $\varepsilon = \frac{1}{4}(\tilde{\theta} - \theta_0)^2$, using $D(\hat{\theta}||\tilde{\theta}) + D(\hat{\theta}||\theta_0) \geq (\tilde{\theta} - \theta_0)^2$, we get $D(\hat{\theta}||\theta_0) \geq 3\varepsilon$. We now choose $\theta \in \Theta$ so near to $\tilde{\theta}$ such that $|D(\hat{\theta}||\theta) - D(\hat{\theta}||\tilde{\theta})| \leq \varepsilon$ (here we use denseness of Θ). Chaining all inequalities we get $D(\hat{\theta}||\theta_0) - D(\hat{\theta}||\theta) \geq 3\varepsilon - \varepsilon - \varepsilon = \varepsilon > 0$. This, together with (3.28) implies $e^{n\varepsilon} \leq c$ for infinitely many n which is impossible. Hence, the assumption $\hat{\theta}_n \not\rightarrow \theta_0$ was wrong.

Now, $\hat{\theta}_n \rightarrow \theta_0$ implies that for arbitrary $\theta \neq \theta_0$, $\theta \in \Theta$ and for sufficiently large n there exists $\delta_{\theta} > 0$ such that $D(\hat{\theta}_n||\theta) \geq 2\delta_{\theta}$ (since $D(\theta_0||\theta) \neq 0$) and $D(\hat{\theta}_n||\theta_0) \leq \delta_{\theta}$. This implies

$$w_n^{\theta} \leq \frac{w_{\theta}}{w_{\theta_0}} e^{n[D(\hat{\theta}_n||\theta_0) - D(\hat{\theta}_n||\theta)]} \leq \frac{w_{\theta}}{w_{\theta_0}} e^{-n\delta_{\theta}} \xrightarrow{n \rightarrow \infty} 0,$$

where we have used (3.25) and (3.26) in the first inequality and the second inequality holds for sufficiently large n . Hence $\sum_{\theta \neq \theta_0} w_n^\theta \rightarrow 0$ by (3.27) and $w_n^{\theta_0} \rightarrow 1$ by normalization (3.25), which finally gives

$$\xi(1|x_{1:n}) = w_n^{\theta_0} \mu_{\theta_0}(1|x_{1:n}) + \sum_{\theta \neq \theta_0} w_n^\theta \mu_\theta(1|x_{1:n}) \xrightarrow{n \rightarrow \infty} \mu_{\theta_0}(1|x_{1:n}).$$

(ii) We first consider the case $\Theta = \{\theta_0, \theta_1\}$: Let us choose $\bar{\theta} (= \ln(\frac{1-\theta_0}{1-\theta_1}) / \ln(\frac{\theta_1}{\theta_0} \frac{1-\theta_0}{1-\theta_1}))$, potentially $\notin \Theta$ in the (KL) middle of θ_0 and θ_1 such that

$$D(\bar{\theta}||\theta_0) = D(\bar{\theta}||\theta_1), \quad 0 < \theta_0 < \bar{\theta} < \theta_1 < 1, \quad (3.29)$$

and choose $x_{1:\infty}$ such that $\hat{\theta}_n := \frac{n_1}{n}$ satisfies $|\hat{\theta}_n - \bar{\theta}| \leq \frac{1}{n}$ ($\Rightarrow \hat{\theta}_n \xrightarrow{n \rightarrow \infty} \bar{\theta}$)

We will show that $x_{1:\infty}$ is $\mu_{\theta_0} \mathcal{M}$ -random and $\mu_{\theta_1} \mathcal{M}$ -random. Obviously no ξ can converge to θ_0 and θ_1 , thus proving \mathcal{M} -non-convergence. $x_{1:\infty}$ is obviously not $\mu_{\theta_0/1}$ M.L.-random, since the relative frequency $\hat{\theta}_n \not\rightarrow \theta_0/1$. $x_{1:\infty}$ is not even $\mu_{\bar{\theta}}$ M.L.-random, since $\hat{\theta}_n$ converges too fast ($\sim \frac{1}{n}$). $x_{1:\infty}$ is indeed very regular, whereas $\frac{n_1}{n}$ of a truly $\mu_{\bar{\theta}}$ M.L.-random sequence has fluctuations of the order $1/\sqrt{n}$. The fast convergence is necessary for doubly μ/ξ -randomness. The reason that $\hat{\theta}_n$ is μ/ξ -random, but not M.L.-random is that \mathcal{M} randomness is a weaker concept than M.L.-randomness for $\mathcal{M} \subset \mathcal{M}_{enum}^{semi}$. Only regularities characterized by $\nu \in \mathcal{M}$ are recognized by μ/ξ -randomness.

In the following we assume that n is sufficiently large such that $\theta_0 \leq \hat{\theta}_n \leq \theta_1$. We need

$$|D(\hat{\theta}||\theta) - D(\bar{\theta}||\theta)| \leq c|\hat{\theta} - \bar{\theta}| \quad \forall \theta, \hat{\theta}, \bar{\theta} \in [\theta_0, \theta_1] \quad \text{with} \quad c := \ln \frac{\theta_1(1-\theta_0)}{\theta_0(1-\theta_1)} < \infty \quad (3.30)$$

which follows for $\hat{\theta} \geq \bar{\theta}$ (similarly $\hat{\theta} \leq \bar{\theta}$) from

$$D(\hat{\theta}||\theta) - D(\bar{\theta}||\theta) = \int_{\bar{\theta}}^{\hat{\theta}} [\ln \frac{\theta'}{\theta} - \ln \frac{1-\theta'}{1-\theta}] d\theta' \leq \int_{\bar{\theta}}^{\hat{\theta}} [\ln \frac{\theta_1}{\theta_0} - \ln \frac{1-\theta_1}{1-\theta_0}] d\theta' = c \cdot (\hat{\theta} - \bar{\theta})$$

where we have increased θ' to θ_1 and decreased θ to θ_0 in the inequality. Using (3.30) in (3.26) twice we get

$$\frac{\mu_{\theta_1}(x_{1:n})}{\mu_{\theta_0}(x_{1:n})} = e^{n[D(\hat{\theta}_n||\theta_0) - D(\hat{\theta}_n||\theta_1)]} \leq e^{n[D(\bar{\theta}||\theta_0) + c|\hat{\theta}_n - \bar{\theta}| - D(\bar{\theta}||\theta_1) + c|\hat{\theta}_n - \bar{\theta}|]} \leq e^{2c} \quad (3.31)$$

where we have used (3.29) in the last inequality. Now, (3.31) and (3.25) lead to

$$w_n^{\theta_0} = w_{\theta_0} \frac{\mu_{\theta_0}(x_{1:n})}{\xi(x_{1:n})} = [1 + \frac{w_{\theta_1} \mu_{\theta_1}(x_{1:n})}{w_{\theta_0} \mu_{\theta_0}(x_{1:n})}]^{-1} \geq [1 + \frac{w_{\theta_1}}{w_{\theta_0}} e^{2c}]^{-1} =: c_0 > 0, \quad (3.32)$$

which shows that $x_{1:\infty}$ is $\mu_{\theta_0} \mathcal{M}$ -random by (3.24). Exchanging $\theta_0 \leftrightarrow \theta_1$ in (3.31) and (3.32) we similarly get $w_n^{\theta_1} \geq c_1 > 0$, which implies (using $w_n^{\theta_0} + w_n^{\theta_1} = 1$)

$$\xi(1|x_{1:n}) = \sum_{\theta \in \{\theta_0, \theta_1\}} w_n^\theta \mu_\theta(1|x_{1:n}) = w_n^{\theta_0} \cdot \theta_0 + w_n^{\theta_1} \cdot \theta_1 \neq \theta_0 = \mu_{\theta_0}(1|x_{1:n}). \quad (3.33)$$

This shows $\xi(1|x_{1:n}) \xrightarrow{n \rightarrow \infty} \mu(1|x_{1:n})$. One can show that $\xi(1|x_{1:n})$ does not only not converge to θ_0 (and θ_1), but that it does not converge at all. The fast convergence demand $|\hat{\theta}_n - \bar{\theta}| \leq \frac{1}{n}$ on $x_{1:\infty}$ can be weakened to $\hat{\theta}_n \leq \bar{\theta} + O(\frac{1}{n}) \forall n$ and $\hat{\theta}_n \geq \bar{\theta} - O(\frac{1}{n})$ for infinitely many n , then $x_{1:\infty}$ is still $\mu_{\theta_0} \mathcal{M}$ -random, and $w_n^{\theta_1} \geq c'_1 > 0$ for infinitely many n , which is sufficient to prove $\xi \not\rightarrow \mu$.

We now consider general Θ with gap in the sense that there exist $0 < \theta_0 < \theta_1 < 1$ with $[\theta_0, \theta_1] \cap \Theta = \{\theta_0, \theta_1\}$: We show that all $\theta \neq \theta_0, \theta_1$ give asymptotically no contribution to $\xi(1|x_{1:n})$, i.e. (3.33) still applies. Let $\theta \in \Theta \setminus \{\theta_0, \theta_1\}$; all other definitions as before. Then $\delta_\theta := D(\bar{\theta}||\theta) - D(\bar{\theta}||\theta_0) > 0$, since θ is farther than θ_0 away from $\bar{\theta}$ ($|\theta - \bar{\theta}| > |\theta_0 - \bar{\theta}|$). Similarly to (3.31) with θ instead θ_1 we get

$$\frac{\mu_\theta(x_{1:n})}{\mu_{\theta_0}(x_{1:n})} = e^{n[D(\hat{\theta}_n||\theta_0) - D(\hat{\theta}_n||\theta)]} \leq e^{2c} \cdot e^{n[D(\bar{\theta}||\theta_0) - D(\bar{\theta}||\theta)]} = e^{2c} e^{-n\delta_\theta} \xrightarrow{n \rightarrow \infty} 0$$

Hence $w_n^\theta \leq \frac{w_\theta}{w_{\theta_0}} e^{2c} e^{-n\delta_\theta} \rightarrow 0$ from (3.25) and $\varepsilon_n := \sum_{\theta \in \Theta \setminus \{\theta_0, \theta_1\}} w_n^\theta \mu_\theta(1|x_{1:n}) \xrightarrow{n \rightarrow \infty} 0$ from (3.27). Hence $\xi(1|x_{1:n}) = w_n^{\theta_0} \cdot \theta_0 + w_n^{\theta_1} \cdot \theta_1 + \varepsilon_n \neq \theta_0 = \mu_{\theta_0}(1|x_{1:n})$ for sufficiently large n , since $\varepsilon_n \rightarrow 0$, $w_n^{\theta_1} \geq c'_1 > 0$ and $\theta_0 \neq \theta_1$. \square

3.2.8 The case where $\mu \notin \mathcal{M}$

In the following we discuss two cases in which $\mu \notin \mathcal{M}$, but most parts of this work still apply. Actually all theorems remain valid for μ being a finite linear combination $\mu(x_{1:n}) = \sum_{\nu \in \mathcal{L}} v_\nu \nu(x_{1:n})$ of ν 's in $\mathcal{L} \subseteq \mathcal{M}$. Dominance $\xi(x_{1:n}) \geq w_\mu \cdot \mu(x_{1:n})$ is still ensured with $w_\mu := \min_{\nu \in \mathcal{L}} \frac{w_\nu}{v_\nu} \geq \min_{\nu \in \mathcal{L}} w_\nu$. More generally, if μ is an infinite linear combination, dominance is still ensured if w_ν itself dominate v_ν in the sense that $w_\nu \geq \alpha v_\nu$ for some $\alpha > 0$ (then $w_\mu \geq \alpha$).

Another possibly interesting situation is when the true generating distribution $\mu \notin \mathcal{M}$, but a “nearby” distribution $\hat{\mu}$ with weight $w_{\hat{\mu}}$ is in \mathcal{M} . If we measure the distance of $\hat{\mu}$ to μ with the Kullback Leibler divergence $D_n(\mu||\hat{\mu}) := \sum_{x_{1:n}} \mu(x_{1:n}) \ln \frac{\mu(x_{1:n})}{\hat{\mu}(x_{1:n})}$ and assume that it is bounded by a constant c , then

$$D_n = \mathbf{E}[\ln \frac{\mu(x_{1:n})}{\xi(x_{1:n})}] = \mathbf{E}[\ln \frac{\hat{\mu}(x_{1:n})}{\xi(x_{1:n})}] + \mathbf{E}[\ln \frac{\mu(x_{1:n})}{\hat{\mu}(x_{1:n})}] \leq \ln w_{\hat{\mu}}^{-1} + c.$$

So $D_n \leq \ln w_{\hat{\mu}}^{-1}$ remains valid if we define $w_\mu := w_{\hat{\mu}} \cdot e^{-c}$.

3.2.9 Probability Classes \mathcal{M}

In the following we describe some well-known and some less known probability classes \mathcal{M} . This relates our setting to other works in this area, embeds it into the historical context, illustrates the type of classes we have in mind, and discusses computational issues.

We get a rather wide class \mathcal{M} if we include *all* (semi)computable probability distributions in \mathcal{M} . In this case, the assumption $\mu \in \mathcal{M}$ is very weak, as it only

assumes that the strings are drawn from *any (semi)computable* distribution; and all valid physical theories (and, hence, all environments) *are* computable to arbitrary precision (estimable) (in a probabilistic sense).

We will see that it is favorable to assign high weights w_ν to the ν . Simplicity should be favored over complexity, according to Occam's razor. In our context this means that a high weight should be assigned to simple ν . The prefix Kolmogorov complexity $K(\nu)$ is a universal complexity measure [Kol65, Lev74, Gác74, Cha75, LV97]. It is defined as the length of the shortest self-delimiting program (on a universal Turing machine) computing $\nu(x_{1:n})$ given $x_{1:n}$. If we define

$$w_\nu := 2^{-K(\nu)}$$

then, distributions which can be calculated by short programs, have high weights. The relative entropy is bounded by the Kolmogorov complexity of μ in this case ($D_n \leq K(\mu) \cdot \ln 2$). Levin's universal semi-measure ξ_U is obtained if we take $\mathcal{M} = \mathcal{M}_U = \mathcal{M}_{enum}^{semi}$ to be the (multi)set enumerated by a Turing machine which enumerates all enumerable semi-measures [ZL70, LV97]. Recently, \mathcal{M} has been further enlarged to include all cumulatively enumerable semi-measures [Sch02a]. In the enumerable and cumulatively enumerable cases, ξ is not finitely computable, but can still be approximated to arbitrary but not pre-specifiable precision. If we consider *all* approximable (i.e. asymptotically computable) distributions, then the universal distribution ξ , although still well defined, is not even approximable. An interesting and quickly approximable distribution is the Speed prior S defined in [Sch02c]. It is related to Levin complexity and Levin search [Lev73b, Lev84], but it is unclear for now which distributions are dominated by S (see Problem 3.2). If one considers only finite-state automata instead of general Turing machines, ξ is related to the quickly computable, universal finite-state prediction scheme of Feder et al. [FMG92], which itself is related to the famous Lempel-Ziv data compression algorithm. If one has extra knowledge on the source generating the sequence, one might further reduce \mathcal{M} and increase w . A detailed analysis of these and other specific classes \mathcal{M} will be given elsewhere. Note that $\xi \in \mathcal{M}$ in the enumerable and cumulatively enumerable case, but $\xi \notin \mathcal{M}$ in the computable, approximable and finite-state case. If ξ is itself in \mathcal{M} , it is called a universal element of \mathcal{M} [LV97]. As we do not need this property here, \mathcal{M} may be *any* countable set of distributions. In the following we consider generic \mathcal{M} and w . Continuous classes \mathcal{M} are considered in Section 3.7.2.

3.3 Error Bounds

3.3.1 Bayes-Optimal Predictors

We start with a very simple measure: making a wrong prediction counts as one error, making a correct prediction counts as no error. Let Θ_μ be the optimal prediction scheme when the strings are drawn from the probability distribution μ , i.e.

the probability of x_t given $x_{<t}$ is $\mu(x_t|x_{<t})$, and μ is known. Θ_μ predicts (by definition) $x_t^{\Theta_\mu}$ when observing $x_{<t}$. The prediction is erroneous if the true t^{th} symbol is not $x_t^{\Theta_\mu}$. The probability of this event is $1 - \mu(x_t^{\Theta_\mu}|x_{<t})$. It is minimized if $x_t^{\Theta_\mu}$ maximizes $\mu(x_t|x_{<t})$. More generally, let Θ_ρ be a prediction scheme predicting $x_t^{\Theta_\rho} := \operatorname{argmax}_{x_t} \rho(x_t|x_{<t})$ for some distribution ρ . Every deterministic predictor can be interpreted as maximizing some distribution.

3.3.2 Total Expected Numbers of Errors

The μ probability of making a wrong prediction for the t^{th} symbol and the total μ -expected number of errors in the first n predictions of predictor Θ_ρ are

$$e_t^{\Theta_\rho}(x_{<t}) := 1 - \mu(x_t^{\Theta_\rho}|x_{<t}) \quad , \quad E_n^{\Theta_\rho} := \sum_{t=1}^n \mathbf{E}[e_t^{\Theta_\rho}(x_{<t})]. \quad (3.34)$$

If μ is known, Θ_μ is obviously the best prediction scheme in the sense of making the least number of expected errors

$$E_n^{\Theta_\mu} \leq E_n^{\Theta_\rho} \quad \text{for any } \Theta_\rho, \quad (3.35)$$

since

$$e_t^{\Theta_\mu}(x_{<t}) = 1 - \mu(x_t^{\Theta_\mu}|x_{<t}) = \min_{x_t} (1 - \mu(x_t|x_{<t})) \leq 1 - \mu(x_t^{\Theta_\rho}|x_{<t}) = e_t^{\Theta_\rho}(x_{<t})$$

for any ρ . Of special interest is the universal predictor Θ_ξ . As ξ converges to μ the prediction of Θ_ξ might converge to the prediction of the optimal Θ_μ . Hence, Θ_ξ may not make many more errors than Θ_μ and, hence, any other predictor Θ_ρ . Note that $x_t^{\Theta_\rho}$ is a discontinuous function of ρ and $x_t^{\Theta_\xi} \rightarrow x_t^{\Theta_\mu}$ does not follow from $\xi \rightarrow \mu$. Indeed, this problem occurs in related prediction schemes, where the predictor has to be regularized so that it is continuous [FMG92]. Fortunately this is not necessary here. We prove the following error bound.

Theorem 3.36 (Error bound) Let there be sequences $x_1 x_2 \dots$ over a finite alphabet \mathcal{X} drawn with probability $\mu(x_{1:n})$ for the first n symbols. The Θ_ρ -system predicts by definition $x_t^{\Theta_\rho} \in \mathcal{X}$ from $x_{<t}$, where $x_t^{\Theta_\rho}$ maximizes $\rho(x_t|x_{<t})$. Θ_ξ is the universal prediction scheme based on the universal prior ξ . Θ_μ is the optimal informed prediction scheme. The total μ -expected number of prediction errors $E_n^{\Theta_\xi}$ and $E_n^{\Theta_\mu}$ of Θ_ξ and Θ_μ as defined in (3.34) are bounded in the following way

$$0 \leq E_n^{\Theta_\xi} - E_n^{\Theta_\mu} \leq \sqrt{2(E_n^{\Theta_\xi} + E_n^{\Theta_\mu})S_n} \leq S_n + \sqrt{4E_n^{\Theta_\mu}S_n + S_n^2} \leq 2S_n + 2\sqrt{E_n^{\Theta_\mu}S_n}$$

where $S_n \leq D_n \leq \ln w_\mu^{-1}$. S_n is the squared distance (3.14), D_n is the relative entropy (3.18), and w_μ is the weight (3.5) of μ in ξ .

The first bound actually contains $E_n^{\Theta_\xi}$ on the r.h.s., so it is not particularly useful, but this is the major bound we will prove, the others follow easily. Furthermore it has a somewhat nicer structure than the second bound. In Section 3.6 we show that the second bound is optimal. The last bound, which we discuss in the following, has the same asymptotics as the second bound.

First, we observe that the number of errors $E_\infty^{\Theta_\xi}$ of the universal Θ_ξ predictor is finite if the number of errors $E_\infty^{\Theta_\mu}$ of the informed Θ_μ predictor is finite. This is especially the case for deterministic μ , as $E_n^{\Theta_\mu} \equiv 0$ in this case⁶, i.e. Θ_ξ makes only a finite number of errors on deterministic environments. This can be proven by elementary means. Assume $x_1x_2\dots$ is the sequence generated by μ and Θ_ξ makes a wrong prediction $x_t^{\Theta_\xi} \neq x_t$. Since $\xi(x_t^{\Theta_\xi}|x_{<t}) \geq \xi(x_t|x_{<t})$, this implies $\xi(x_t|x_{<t}) \leq \frac{1}{2}$. Hence $e_t^{\Theta_\xi} = 1 \leq -\ln \xi(x_t|x_{<t}) / \ln 2 = d_t / \ln 2$. If Θ_ξ makes a correct prediction $e_t^{\Theta_\xi} = 0 \leq d_t / \ln 2$ is obvious. Using (3.18) proves $E_\infty^{\Theta_\xi} \leq D_\infty / \ln 2 \leq \log_2 w_\mu^{-1}$. A combinatoric argument given in Section 3.6 shows that there are \mathcal{M} and $\mu \in \mathcal{M}$ with $E_\infty^{\Theta_\xi} \geq \log_2 |\mathcal{M}|$. This shows that the upper bound $E_\infty^{\Theta_\xi} \leq \log_2 |\mathcal{M}|$ for uniform w is sharp. From Theorem 3.36 we get the slightly weaker bound $E_\infty^{\Theta_\xi} \leq 2S_\infty \leq 2D_\infty \leq 2\ln w_\mu^{-1}$. For more complicated probabilistic environments, where even the ideal informed system makes an infinite number of errors, the theorem ensures that the error regret $E_n^{\Theta_\xi} - E_n^{\Theta_\mu}$ is only of order $\sqrt{E_n^{\Theta_\mu}}$. The regret is quantified in terms of the information content D_n of μ (relative to ξ), or the weight w_μ of μ in ξ . This ensures that the error densities E_n/n of both systems converge to each other. Actually, the theorem ensures more, namely that the quotient converges to 1, and also gives the speed of convergence $E_n^{\Theta_\xi} / E_n^{\Theta_\mu} = 1 + O((E_n^{\Theta_\mu})^{-1/2}) \rightarrow 1$ for $E_n^{\Theta_\mu} \rightarrow \infty$. Increasing the first occurrence of $E_n^{\Theta_\mu}$ in the theorem to E_n^Θ and the second $E_n^{\Theta_\mu}$ to $E_n^{\Theta_\xi}$ we get the bound $E_n^\Theta \geq E_n^{\Theta_\xi} - 2\sqrt{E_n^{\Theta_\xi} S_n}$, which shows that *no* (causal) predictor Θ whatsoever makes significantly less errors than Θ_ξ . In Section 3.6 we show that the second bound for $E_n^{\Theta_\xi} - E_n^{\Theta_\mu}$ given in Theorem 3.36 can in general not be improved, i.e. for every predictor Θ (and especially Θ_ξ) there exist \mathcal{M} and $\mu \in \mathcal{M}$ such that the upper bound is achieved. See [Hut01c] for some further discussion and bounds for binary alphabet.

3.3.3 Proof of Theorem 3.36

The first inequality in Theorem 3.36 has already been proven (3.35). For the second inequality, let us start more modestly and try to find constants $A > 0$ and $B > 0$ that satisfy the linear inequality

$$E_n^{\Theta_\xi} - E_n^{\Theta_\mu} \leq A(E_n^{\Theta_\xi} + E_n^{\Theta_\mu}) + BS_n. \quad (3.37)$$

⁶Remember that we named a probability distribution *deterministic* if it is 1 for exactly one sequence and 0 for all others.

If we could show

$$e_t^{\Theta_\xi}(x_{<t}) - e_t^{\Theta_\mu}(x_{<t}) \leq A[e_t^{\Theta_\xi}(x_{<t}) + e_t^{\Theta_\mu}(x_{<t})] + Bs_t(x_{<t}) \quad (3.38)$$

for all $t \leq n$ and all $x_{<t}$, (3.37) would follow immediately by summation and the definition of E_n and S_n . With the abbreviations (3.12) and the abbreviations $m = x_t^{\Theta_\mu}$ and $s = x_t^{\Theta_\xi}$ the various error functions can then be expressed by $e_t^{\Theta_\xi} = 1 - y_s$, $e_t^{\Theta_\mu} = 1 - y_m$ and $s_t = \sum_i (y_i - z_i)^2$. Inserting this into (3.38) we get

$$y_m - y_s \leq A[2 - (y_m + y_s)] + B \sum_{i=1}^N (y_i - z_i)^2. \quad (3.39)$$

By definition of $x_t^{\Theta_\mu}$ and $x_t^{\Theta_\xi}$ we have $y_m \geq y_i$ and $z_s \geq z_i$ for all i . We prove a sequence of inequalities which show that

$$B \sum_{i=1}^N (y_i - z_i)^2 + A[2 - (y_m + y_s)] - (y_m - y_s) \geq \dots \quad (3.40)$$

is positive for suitable $A \geq 0$ and $B \geq 0$, which proves (3.39). For $m = s$ (3.40) is obviously positive. So we will assume $m \neq s$ in the following. From the square we keep only contributions from $i = m$ and $i = s$.

$$\dots \geq B[(y_m - z_m)^2 + (y_s - z_s)^2] + A[2 - (y_m + y_s)] - (y_m - y_s) \geq \dots$$

By definition of y , z , \mathcal{M} and s we have the constraints $y_m + y_s \leq 1$, $z_m + z_s \leq 1$, $y_m \geq y_s \geq 0$ and $z_s \geq z_m \geq 0$. From the latter two it is easy to see that the square terms (as a function of z_m and z_s) are minimized by $z_m = z_s = \frac{1}{2}(y_m + y_s)$. Furthermore, we define $x := y_m - y_s$ and increase $(y_m + y_s)$ to 1.

$$\dots \geq \frac{1}{2}Bx^2 + A - x \geq \dots \quad (3.41)$$

(3.41) is quadratic in x and minimized by $x^* = \frac{1}{B}$. Inserting x^* gives

$$\dots \geq A - \frac{1}{2B} \geq 0 \quad \text{for} \quad 2AB \geq 1. \quad (3.42)$$

Inequality (3.37) therefore holds for any $A > 0$, provided we insert $B = \frac{1}{2A}$. Thus we might minimize the r.h.s. of (3.37) w.r.t. A leading to the upper bound

$$E_n^{\Theta_\xi} - E_n^{\Theta_\mu} \leq \sqrt{2(E_n^{\Theta_\xi} + E_n^{\Theta_\mu})S_n} \quad \text{for} \quad A^2 = \frac{S_n}{2(E_n^{\Theta_\xi} + E_n^{\Theta_\mu})} \quad (3.43)$$

which is the first bound in Theorem 3.36. For the second bound we have to prove

$$\sqrt{2(E_n^{\Theta_\xi} + E_n^{\Theta_\mu})S_n} - S_n \leq \sqrt{4E_n^{\Theta_\mu}S_n + S_n^2} \quad (3.44)$$

If we square both sides of this expressions and simplify we just get (3.43). Hence, (3.43) implies (3.44). The last inequality in Theorem 3.36 is a simple triangle inequality. This completes the proof of Theorem 3.36. \square

Note that also the third bound implies the second one:

$$\begin{aligned} E_n^{\Theta_\xi} - E_n^{\Theta_\mu} &\leq \sqrt{2(E_n^{\Theta_\xi} + E_n^{\Theta_\mu})S_n} \Leftrightarrow (E_n^{\Theta_\xi} - E_n^{\Theta_\mu})^2 \leq 2(E_n^{\Theta_\xi} + E_n^{\Theta_\mu})S_n \Leftrightarrow \\ \Leftrightarrow (E_n^{\Theta_\xi} - E_n^{\Theta_\mu} - S_n)^2 &\leq 4E_n^{\Theta_\mu}S_n + S_n^2 \Leftrightarrow E_n^{\Theta_\xi} - E_n^{\Theta_\mu} - S_n \leq \sqrt{4E_n^{\Theta_\mu}S_n + S_n^2} \end{aligned}$$

where we only have used $E_n^{\Theta_\xi} \geq E_n^{\Theta_\mu}$. Nevertheless the bounds are not equal. In Section 3.9 we give an alternative direct proof of the second bound.

3.4 Loss Bounds

3.4.1 Unit Loss Function

We now generalize the error bound derived in the last section to arbitrary bounded loss function. A prediction is very often the basis for some decision. The decision results in an action, which itself leads to some reward or loss. If the action itself can influence the environment we enter the domain of acting agents which will be analyzed in the context of universal probability in later chapters. To stay in the framework of (passive) prediction we have to assume that the action itself does not influence the environment. Let $\ell_{x_t y_t} \in \mathbb{R}$ be the received loss when taking action $y_t \in \mathcal{Y}$ and $x_t \in \mathcal{X}$ is the t^{th} symbol of the sequence. We demand ℓ to be normalized, i.e. $0 \leq \ell_{x_t y_t} \leq 1$. For instance, if we make a sequence of weather forecasts $\mathcal{X} = \{\text{sunny, rainy}\}$ and base our decision, whether to take an umbrella or wear sunglasses $\mathcal{Y} = \{\text{umbrella, sunglasses}\}$ on it, the action of taking the umbrella or wearing sunglasses does not influence the future weather (ignoring the butterfly effect). The losses might be

Loss	sunny	rainy
umbrella	0.3	0.1
sunglasses	0.0	1.0

Note the small loss assignment even when making the right decision to take an umbrella when it rains because sun is still preferable to rain.

In many cases the prediction of x_t can be identified or is already the action y_t . The forecast *sunny* can be identified with the action *wear sunglasses*, and *rainy* with *take umbrella*. $\mathcal{X} \equiv \mathcal{Y}$ in these cases. The error assignment of the previous subsection falls into this class together with a special loss function. It assigns unit loss to an erroneous prediction ($\ell_{x_t y_t} = 1$ for $x_t \neq y_t$) and no loss to a correct prediction ($\ell_{x_t x_t} = 0$).

For convenience we name an action a prediction in the following, even if $\mathcal{X} \neq \mathcal{Y}$. The true probability of the next symbol being x_t , given $x_{<t}$, is $\mu(x_t | x_{<t})$. The

expected loss when predicting y_t is $\mathbf{E}_t[\ell_{x_t y_t}]$. The goal is to minimize the expected loss. More generally we define the Λ_ρ prediction scheme

$$y_t^{\Lambda_\rho} := \arg \min_{y_t \in \mathcal{Y}} \sum_{x_t} \rho(x_t | x_{<t}) \ell_{x_t y_t} \quad (3.45)$$

which minimizes the ρ -expected loss.⁷ As the true distribution is μ , the actual μ -expected loss when Λ_ρ predicts the t^{th} symbol and the total μ -expected loss in the first n predictions are

$$l_t^{\Lambda_\rho}(x_{<t}) := \mathbf{E}_t[\ell_{x_t y_t^{\Lambda_\rho}}], \quad L_n^{\Lambda_\rho} := \sum_{t=1}^n \mathbf{E}[l_t^{\Lambda_\rho}(x_{<t})]. \quad (3.46)$$

Let Λ be *any* (causal) prediction scheme (deterministic or probabilistic) with no constraint at all, predicting *any* $y_t^\Lambda \in \mathcal{Y}$ with losses l_t^Λ and L_n^Λ similarly defined as (3.46). If μ is known, Λ_μ is obviously the best prediction scheme in the sense of achieving minimal expected loss

$$L_n^{\Lambda_\mu} \leq L_n^\Lambda \quad \text{for any } \Lambda \quad (3.47)$$

since

$$l_t^{\Lambda_\mu}(x_{<t}) = \mathbf{E}_t \ell_{x_t y_t^{\Lambda_\mu}} = \min_{y_t} \mathbf{E}_t \ell_{x_t y_t} \leq \mathbf{E}_t \ell_{x_t y_t^\Lambda} = l_t^\Lambda(x_{<t})$$

for any Λ . The predictor Λ_ξ , based on the universal distribution ξ , is, again, of special interest. Theorem 3.36 generalizes to arbitrary loss functions.

Theorem 3.48 (Unit loss bound) Let there be sequences $x_1 x_2 \dots$ over a finite alphabet \mathcal{X} drawn with probability $\mu(x_{1:n})$ for the first n symbols. A system taking action (or predicting) $y_t \in \mathcal{Y}$ given $x_{<t}$ receives loss $\ell_{x_t y_t} \in [0,1]$ if x_t is the true t^{th} symbol of the sequence. The Λ_ρ -system (3.45) acts (or predicts) as to minimize the ρ -expected loss. Λ_ξ is the universal prediction scheme based on the universal prior ξ . Λ_μ is the optimal informed prediction scheme. The total μ -expected losses $L_n^{\Lambda_\xi}$ of Λ_ξ and $L_n^{\Lambda_\mu}$ of Λ_μ as defined in (3.46) are bounded in the following way

$$0 \leq L_n^{\Lambda_\xi} - L_n^{\Lambda_\mu} \leq D_n + \sqrt{4L_n^{\Lambda_\mu} D_n + D_n^2} \leq 2D_n + 2\sqrt{L_n^{\Lambda_\mu} D_n}$$

where $D_n \leq \ln w_\mu^{-1}$ is the relative entropy (3.18), and w_μ is the weight (3.5) of μ in ξ .

The loss bounds have the same form as the error bounds when substituting $S_n \leq D_n$ in Theorem 3.36, so most of the discussion of Theorem 3.36 also applies here. We

⁷ $\arg \min_y(\cdot)$ is defined as the y which minimizes the argument. A tie is broken arbitrarily. If \mathcal{Y} is finite, then $y_t^{\Lambda_\rho}$ always exists. For infinite action space \mathcal{Y} we assume that a minimizing $y_t^{\Lambda_\rho} \in \mathcal{Y}$ exists. This is for instance the case if \mathcal{Y} is compact and ℓ_{xy} is continuous in y , or for $\mathcal{Y} = \mathbb{N}$, if $\lim_{y \rightarrow \infty} \ell_{xy}$ exists for all x and is larger or equal to ℓ_{xy} for most y .

were not able to derive loss bounds in terms of S_n as in the error case, and indeed one can show that substituting S_n for D_n in Theorem 3.48 gives an invalid bound. For convenience we collect the most important consequences of Theorem 3.48 in the following corollary.

Corollary 3.49 (Unit loss bound) Under the same conditions as in Theorem 3.48 the following relations hold.

- i) $L_\infty^{\Lambda_\xi}$ is finite $\iff L_\infty^{\Lambda_\mu}$ is finite,
- ii) $L_\infty^{\Lambda_\xi} \leq 2D_\infty \leq 2 \ln w_\mu^{-1}$ for deterministic μ if $\forall x \exists y \ell_{xy} = 0$,
- iii) $L_n^{\Lambda_\xi} / L_n^{\Lambda_\mu} = 1 + O((L_n^{\Lambda_\mu})^{-1/2}) \rightarrow 1$ for $L_n^{\Lambda_\mu} \rightarrow \infty$,
- iv) $L_n^{\Lambda_\xi} - L_n^{\Lambda_\mu} = O(\sqrt{L_n^{\Lambda_\mu}})$,

Let Λ be *any* prediction scheme.

- v) $L_n^{\Lambda_\mu} \leq L_n^\Lambda$, $l_t^{\Lambda_\mu}(x_{<t}) \leq l_t^\Lambda(x_{<t})$,
- vi) $L_n^\Lambda \geq L_n^{\Lambda_\xi} - 2\sqrt{L_n^{\Lambda_\xi} D_n}$,
- vii) $L_n^{\Lambda_\xi} / L_n^\Lambda \leq 1 + O((L_n^\Lambda)^{-1/2})$.

3.4.2 Loss Bound of Merhav & Feder

The first general loss bound with no structural assumptions on μ and ℓ (except boundedness) has been derived in a survey paper by Merhav and Feder in [MF98, Sec.III.A.2]. They showed that the regret $L_n^{\Lambda_\xi} - L_n^{\Lambda_\mu}$ is bounded by $\ell_{\max} \sqrt{2nD_n}$ for $\ell \in [0, \ell_{\max}]$. Assuming $\ell_{\max} = 1$ (general ℓ_{\max} can be recovered by scaling) their bound reads (in our notation)

$$L_n^{\Lambda_\xi} - L_n^{\Lambda_\mu} \leq A_n \leq \sqrt{2nD_n}. \quad (3.50)$$

Later (Theorem 3.59) we prove

$$l_t^{\Lambda_\xi}(x_{<t}) - l_t^{\Lambda_\mu}(x_{<t}) \leq a_t(x_{<t}) \leq \sqrt{2d_t(x_{<t})}$$

Taking the expectation \mathbf{E} and the average $\frac{1}{n} \sum_{t=1}^n$ and using Jensen's inequality for the concave square root (similarly to (3.21)) or directly Theorem 3.19(vi) shows (3.50).

Bound (3.50) and our bound (Theorem 3.48) are in general incomparable. Since $2D_\infty$ is finite and $L_n^{\Lambda_\mu} \leq n$, bound (3.50) can be at best a factor $\sqrt{2}$ and an additive constant better than our bound. On the other hand, for large n and for $L_n^{\Lambda_\mu} < \frac{n}{2}$ our bound is tighter. The latter condition is satisfied if the best predictor Λ_μ suffers

small instantaneous loss $< \frac{1}{2}$ on average. Significant improvement occurs if $L_n^{\Lambda_\mu}$ does not grow linearly with n , but is for instance finite (see Corollary 3.49, especially (i) and (ii)).

3.4.3 Example Loss Functions

The case $\mathcal{X} \equiv \mathcal{Y}$ with unit error assignment $\ell_{xy} = 1 - \delta_{xy}$ ($\delta_{xy} = 1$ for $x = y$ and $\delta_{xy} = 0$ for $x \neq y$) has already been discussed and proven in Section 3.3.

$$y_t^{\Lambda_\rho} = \arg \min_{y_t} \sum_{x_t} \rho(x_t | x_{<t}) (1 - \delta_{x_t y_t}) = \arg \max_{x_t} \rho(x_t | x_{<t}) = x_t^{\Theta_\rho}$$

In this case $L_n^{\Lambda_\rho} \equiv E_n^{\Theta_\rho}$ is the total expected number of prediction errors. For $\mathcal{X} = \mathcal{Y} = \{0, 1\}$, like in the weather example above, Λ_ρ is a threshold strategy with $y_t^{\Lambda_\rho} = \arg \min_{y \in \{0, 1\}} \{\rho_1 \ell_{1y} + \rho_0 \ell_{0y}\} = 0/1$ for $\rho_1 \gtrless \gamma$, where $\gamma := \frac{\ell_{01} - \ell_{00}}{\ell_{01} - \ell_{00} + \ell_{10} - \ell_{11}}$ and $\rho_i = \rho(i | x_{<t})$. In the special error case $\ell_{xy} = 1 - \delta_{xy}$, the bit with the highest ρ probability is predicted ($\gamma = \frac{1}{2}$). In the following we consider some standard loss functions for binary outcome $\mathcal{X} = \{0, 1\}$ and continuous action y in the unit interval $\mathcal{Y} = [0, 1]$. The *absolute loss* is defined as $\ell_{xy} = |x - y| \in [0, 1]$. The Λ_ρ scheme predicts $y_t^{\Lambda_\rho} = \arg \min_{y \in [0, 1]} \{\rho_1(1 - y) + \rho_0 y\} = 0/1$ for $\rho_0 \gtrless \rho_1$. Since all predictions y lie in the subset $\{0, 1\} \subset [0, 1]$ and $|x - y| = 1 - \delta_{xy}$ for $y \in \{0, 1\}$ this case coincides with the binary error case above. The same holds for the α -loss $|x - y|^\alpha$ with $0 < \alpha \leq 1$. The μ -expected loss is $l_t^{\Lambda_\rho} = \mu(i | x_{<t})$ for the i with $\rho_i > \frac{1}{2}$. For the *quadratic loss* $\ell_{xy} = (x - y)^2 \in [0, 1]$ the action/prediction $y_t^{\Lambda_\rho} = \arg \min_{y \in [0, 1]} \{\rho_1(1 - y)^2 + \rho_0 y^2\} = \rho_1$ is proportional to the ρ probability of $x_t = 1$ and $l_t^{\Lambda_\rho} = \mathbf{E}_t(1 - \rho(x_t | x_{<t}))^2$. For the α -loss $|x - y|^\alpha$ with $\alpha > 1$ we get $y_t^{\Lambda_\rho} = (1 + \alpha^{-1} \sqrt{\rho_0/\rho_1})^{-1}$. For arbitrary finite alphabet \mathcal{X} and vector-valued predictions \mathbf{y} the quadratic loss may be generalized to $\ell_{xy} = \frac{1}{2} \mathbf{y}^T \mathbf{A}_x \mathbf{y} + \mathbf{b}_x^T \mathbf{y} + c_x$. The *Hellinger loss* can be written for binary outcome in the form $\ell_{xy} = 1 - \sqrt{|1 - x - y|} \in [0, 1]$ with $y_t^{\Lambda_\rho} = \rho_1^2 / (\rho_0^2 + \rho_1^2)$ and $l_t^{\Lambda_\rho} = 1 - (\mu_0 \rho_0 + \mu_1 \rho_1) / \sqrt{\rho_0^2 + \rho_1^2}$. The *logarithmic loss* $\ell_{xy} = -\ln |1 - x - y| \in [0, \infty]$ is unbounded. But since the corresponding action is $y_t^{\Lambda_\rho} = \rho_1$ the expected loss is $l_t^{\Lambda_\rho} = -\mathbf{E}_t \ln \rho(x_t | x_{<t})$. Hence $l_t^{\Lambda_\xi} - l_t^{\Lambda_\mu} = d_t$ and the total loss regret $L_n^{\Lambda_\xi} - L_n^{\Lambda_\mu} = D_n \leq \ln w_\mu^{-1}$ is finite anyway and Theorem 3.48 is not needed. Continuous outcome spaces \mathcal{X} are briefly discussed in Section 3.7.5.

3.4.4 Proof of Theorem 3.48

The first inequality in Theorem 3.48 has already been proven (3.47). For the second and last inequality, we start, as in Theorem 3.36, by looking for constants $A > 0$ and $B > 0$, which satisfy the linear inequality

$$L_n^{\Lambda_\xi} \leq (A + 1)L_n^{\Lambda_\mu} + (B + 1)D_n. \quad (3.51)$$

If we could show

$$l_t^{\Lambda_\xi}(x_{<t}) \leq A' l_t^{\Lambda_\mu}(x_{<t}) + B' d_t(x_{<t}), \quad A' := A + 1, \quad B' := B + 1 \quad (3.52)$$

for all $t \leq n$ and all $x_{<t}$, (3.51) would follow immediately by summation and the definition of L_n and D_n . With the abbreviations $m = y_t^{\Lambda^\mu}$ and $s = y_t^{\Lambda^\xi}$ and the abbreviations (3.12) the loss and entropy can then be expressed by $l_t^{\Lambda^\xi} = \sum_i y_i \ell_{is}$, $l_t^{\Lambda^\mu} = \sum_i y_i \ell_{im}$ and $d_t = \sum_i y_i \ln \frac{y_i}{z_i}$. Inserting this into (3.52) we get

$$\sum_{i=1}^N y_i \ell_{is} \leq A' \sum_{i=1}^N y_i \ell_{im} + B' \sum_{i=1}^N y_i \ln \frac{y_i}{z_i} \quad (3.53)$$

By definition (3.45) of $y_t^{\Lambda^\mu}$ and $y_t^{\Lambda^\xi}$ we have

$$\sum_i y_i \ell_{im} \leq \sum_i y_i \ell_{ij} \quad \text{and} \quad \sum_i z_i \ell_{is} \leq \sum_i z_i \ell_{ij} \quad \text{for all } j. \quad (3.54)$$

Actually, we need the first constraint only for $j = s$ and the second for $j = m$. In Section 3.9 we reduce the problem to the binary $N=2$ case, which we will consider in the following. We take $\sum_{i=0}^1$ instead of $\sum_{i=1}^2$ for convenience.

$$B' \sum_{i=0}^1 y_i \ln \frac{y_i}{z_i} + \sum_{i=0}^1 y_i (A' \ell_{im} - \ell_{is}) \stackrel{?}{\geq} 0 \quad (3.55)$$

The cases $\ell_{im} > \ell_{is} \forall i$ and $\ell_{is} > \ell_{im} \forall i$ contradict the first/second inequality (3.54). Hence we can assume $\ell_{0m} \geq \ell_{0s}$ and $\ell_{1m} \leq \ell_{1s}$. The symmetric case $\ell_{0m} \leq \ell_{0s}$ and $\ell_{1m} \geq \ell_{1s}$ is proven analogously or can be reduced to the first case by renumbering the indices ($0 \leftrightarrow 1$). Using the abbreviations $a := \ell_{0m} - \ell_{0s}$, $b := \ell_{1s} - \ell_{1m}$, $c := y_1 \ell_{1m} + y_0 \ell_{0s}$, $y = y_1 = 1 - y_0$ and $z = z_1 = 1 - z_0$ we can write (3.55) as

$$f(y, z) := B' [y \ln \frac{y}{z} + (1-y) \ln \frac{1-y}{1-z}] + A'(1-y)a - yb + Ac \stackrel{?}{\geq} 0 \quad (3.56)$$

for $zb \leq (1-z)a$ and $0 \leq a, b, c, y, z \leq 1$. The constraint (3.54) on y has been dropped since (3.56) will turn out to be true for all y . Furthermore, we can assume that $d := A'(1-y)a - yb \leq 0$ since for $d > 0$, f is trivially positive. Multiplying d with a constant ≥ 1 will decrease f . Let us first consider the case $z \leq \frac{1}{2}$. We multiply the d term by $1/b \geq 1$, i.e. replace it with $A'(1-y)\frac{a}{b} - y$. From the constraint on z we know that $\frac{a}{b} \geq \frac{z}{1-z}$. We can decrease f further by replacing $\frac{a}{b}$ by $\frac{z}{1-z}$ and by dropping Ac . Hence, (3.56) is proven for $z \leq \frac{1}{2}$ if we can prove

$$B' [y \ln \frac{y}{z} + (1-y) \ln \frac{1-y}{1-z}] + A'(1-y)\frac{z}{1-z} - y \stackrel{?}{\geq} 0 \quad \text{for } z \leq \frac{1}{2}. \quad (3.57)$$

Section 3.9 we prove that it holds for $B \geq \frac{1}{A} + 1$. The case $z \geq \frac{1}{2}$ is treated similarly. We scale d with $1/a \geq 1$, i.e. replace it with $A'(1-y) - y\frac{b}{a}$. From the constraint on z we know that $\frac{b}{a} \leq \frac{1-z}{z}$. We decrease f further by replacing $\frac{b}{a}$ by $\frac{1-z}{z}$ and by dropping Ac . Hence (3.56) is proven for $z \geq \frac{1}{2}$ if we can prove

$$B' [y \ln \frac{y}{z} + (1-y) \ln \frac{1-y}{1-z}] + A'(1-y) - y\frac{1-z}{z} \stackrel{?}{\geq} 0 \quad \text{for } z \geq \frac{1}{2}. \quad (3.58)$$

In Section 3.9 we prove that it holds for $B \geq \frac{1}{A} + 1$. So in summary we proved that (3.51) holds for $B \geq \frac{1}{A} + 1$. Inserting $B = \frac{1}{A} + 1$ into (3.51) and minimizing the r.h.s. w.r.t. A leads to the last bound of Theorem 3.48 with $A = \sqrt{D_n/L_n^{\Lambda_\mu}}$. Actually inequalities (3.57) and (3.58) also hold for $B \geq \frac{1}{4}A + \frac{1}{A}$, which, by the same minimization argument, proves the slightly tighter second bound in Theorem 3.48. Unfortunately, the current proof is very long and complex, and involves some numerical or graphical analysis for determining intersection properties of some higher order polynomials. This or a hopefully simplified proof will be postponed. The cautious reader may check the inequalities (3.57) and (3.58) numerically for $B = \frac{1}{4}A + \frac{1}{A}$. \square

3.4.5 Convergence of Instantaneous Losses

Since $L_n^{\Lambda_\xi} - L_n^{\Lambda_\mu}$ is not finitely bounded by Theorem 3.48 it cannot be used directly to conclude $l_t^{\Lambda_\xi} - l_t^{\Lambda_\mu} \rightarrow 0$. It would follow from $\xi \rightarrow \mu$ by continuity if $l_t^{\Lambda_\xi}$ and $l_t^{\Lambda_\mu}$ would be continuous functions of ξ and μ . $l_t^{\Lambda_\mu}$ is a continuous piecewise linear concave function of μ , but $l_t^{\Lambda_\xi}$ is an, in general, discontinuous function of ξ (and μ). Fortunately it is continuous at the one necessary point $\xi = \mu$. This allows to bound $l_t^{\Lambda_\xi} - l_t^{\Lambda_\mu}$ in terms of $\xi(x_t|x_{<t}) - \mu(x_t|x_{<t})$.

Theorem 3.59 (Instantaneous Loss Bound) Under the same conditions as in Theorem 3.48 the following relations hold for the instantaneous losses $l_t^{\Lambda_\mu}(x_{<t})$ and $l_t^{\Lambda_\xi}(x_{<t})$ at time t of the informed and universal prediction schemes Λ_μ and Λ_ξ :

$$\begin{aligned} i) \quad & \sum_{t=1}^n \mathbf{E}[(l_t^{\Lambda_\xi}(x_{<t}) - l_t^{\Lambda_\mu}(x_{<t}))^2] \leq 2D_n \leq 2 \ln w_\mu^{-1} < \infty \\ ii) \quad & 0 \leq l_t^{\Lambda_\xi}(x_{<t}) - l_t^{\Lambda_\mu}(x_{<t}) \leq \sum_{x_t} |\xi(x_t|x_{<t}) - \mu(x_t|x_{<t})| \leq \sqrt{2d_t(x_{<t})} \xrightarrow[t \rightarrow \infty]{w.\mu.p.1} 0. \\ iii) \quad & 0 \leq l_t^{\Lambda_\xi}(x_{<t}) - l_t^{\Lambda_\mu}(x_{<t}) \leq 2d_t(x_{<t}) + 2\sqrt{l_t^{\Lambda_\mu}(x_{<t}) d_t(x_{<t})} \xrightarrow[t \rightarrow \infty]{w.\mu.p.1} 0. \end{aligned}$$

(i) implies that the expected number of times t in which $l_t^{\Lambda_\xi}$ exceeds $l_t^{\Lambda_\mu}$ by more than ε is finite and bounded by $2\varepsilon^{-2} \ln w_\mu^{-1}$.

Proof. (ii) follows from

$$\begin{aligned} l_t^{\Lambda_\xi}(x_{<t}) - l_t^{\Lambda_\mu}(x_{<t}) & \equiv \sum_i y_i \ell_{is} - \sum_i y_i \ell_{im} \leq \sum_i (y_i - z_i)(\ell_{is} - \ell_{im}) \leq \\ & \leq \sum_i |y_i - z_i| \cdot |\ell_{is} - \ell_{im}| \leq \sum_i |y_i - z_i| \leq \sqrt{2 \sum_i y_i \ln \frac{y_i}{z_i}} \equiv \sqrt{2d_t(x_{<t})} \end{aligned}$$

To arrive at the first inequality we added $\sum_i z_i(\ell_{im} - \ell_{is})$ which is positive due to (3.54). $|\ell_{is} - \ell_{im}| \leq 1$ since $\ell \in [0,1]$. The last inequality follows from Lemma 3.11a. $d_t \rightarrow 0$ has been proven in Theorem 3.19(ii). (i) follows by inserting (ii) and using (3.18). (iii) follows from the proof of Theorem 3.48 by inserting $B = \frac{1}{A} + 1 = \sqrt{l_t^{\Lambda_\mu}/d_t} + 1$ into (3.52). Convergence to zero holds for μ random sequences, i.e. with μ probability 1, since $l_t^{\Lambda_\mu} \leq 1$ is bounded. The losses $l_t^{\Lambda_\rho}(x_{<t})$ itself need not to converge. \square

Note, that the inequalities in (ii) and (iii) hold for all individual sequences. The sum/average is only taken over the current outcome x_t , but the history $x_{<t}$ is fixed. Bound (ii) and (iii) are in general incomparable, but for large t and for $l_t^{\Lambda_\mu} < \frac{1}{2}$ (especially if $l_t^{\Lambda_\mu} \rightarrow 0$) bound (iii) is tighter than bound (ii).

3.4.6 General Loss

There are only very few restrictions imposed on the loss $\ell_{x_t y_t}$ in Theorem 3.48, namely that it is static and in the unit interval $[0,1]$. If we look at the proof of Theorem 3.48, we see that the time-independence has not been used at all. The proof is still valid for an individual loss function $\ell_{x_t y_t}^t \in [0,1]$ for each step t . The loss might even depend on the actual history $x_{<t}$. The case of a loss $\ell_{x_t y_t}^t(x_{<t})$ bounded to a general interval $[\ell_{min}, \ell_{max}]$ can be reduced to the unit interval case by rescaling ℓ . We introduce a scaled loss ℓ'

$$0 \leq \ell_{x_t y_t}^t(x_{<t}) := \frac{\ell_{x_t y_t}^t(x_{<t}) - \ell_{min}}{\ell_\Delta} \leq 1, \quad \text{where} \quad \ell_\Delta := \ell_{max} - \ell_{min}.$$

The prediction scheme Λ'_ρ based on ℓ' is identical to the original prediction scheme Λ_ρ based on ℓ , since argmin in (3.45) is not affected by linear transformation of its argument. From $y_t^{\Lambda'_\rho} = y_t^{\Lambda_\rho}$ it follows that $l_t^{\Lambda'_\rho} = (l_t^{\Lambda_\rho} - \ell_{min})/\ell_\Delta$ and $L_n^{\Lambda'_\rho} = (L_n^{\Lambda_\rho} - \ell_{min})/\ell_\Delta$ ($D'_n \equiv D_n$, since ℓ is not involved). Theorem 3.48 is valid for the primed quantities, since $\ell' \in [0,1]$. Inserting $L'_{n\Lambda_\mu/\xi}$ and rearranging terms we get

Theorem 3.60 (General loss bound) Let there be sequences $x_1x_2\dots$ over a finite alphabet \mathcal{X} drawn with probability $\mu(x_{1:n})$ for the first n symbols. A system taking action (or predicting) $y_t \in \mathcal{Y}$ given $x_{<t}$ receives loss $\ell_{x_t y_t}^t(x_{<t}) \in [\ell_{\min}, \ell_{\min} + \ell_{\Delta}]$ if x_t is the true t^{th} symbol of the sequence. The Λ_{ρ} -system (3.45) acts (or predicts) as to minimize the ρ -expected loss. Λ_{ξ} is the universal prediction scheme based on the universal prior ξ . Λ_{μ} is the optimal informed prediction scheme. The total μ -expected losses $L_n^{\Lambda_{\xi}}$ and $L_n^{\Lambda_{\mu}}$ of Λ_{ξ} and Λ_{μ} as defined in (3.46) are bounded in the following way

$$0 \leq L_n^{\Lambda_{\xi}} - L_n^{\Lambda_{\mu}} \leq \ell_{\Delta} D_n + \sqrt{4(L_n^{\Lambda_{\mu}} - n\ell_{\min})\ell_{\Delta} D_n + \ell_{\Delta}^2 D_n^2}$$

where $D_n \leq \ln w_{\mu}^{-1}$ is the relative entropy (3.18), and w_{μ} is the weight (3.5) of μ in ξ .

3.5 Application to Games of Chance

3.5.1 Introduction

Consider investing in the stock market. At time t an amount of money s_t is invested in portfolio y_t , where we have access to past knowledge $x_{<t}$ (e.g. charts). After our choice of investment we receive new information x_t , and the new portfolio value is r_t . The best we can expect is to have a probabilistic model μ of the behavior of the stock-market. The goal is to maximize the net μ -expected profit $p_t = r_t - s_t$. Nobody knows μ , but the assumption of all traders is that there *is* a computable, profitable μ they try to find or approximate. From Theorem 3.19 we know that Solomonoff's universal prior $\xi(x_t|x_{<t})$ converges to any computable $\mu(x_t|x_{<t})$ with probability 1. If there is a computable, asymptotically profitable trading scheme at all, the Λ_{ξ} scheme should also be profitable in the long run. To get a practically useful, computable scheme we have to restrict \mathcal{M} to a finite set of computable distributions, e.g. with bounded Levin complexity Kt [LV97, Ch.7]. Although convergence of ξ to μ is pleasing, what we are really interested in is whether Λ_{ξ} is asymptotically profitable and how long it takes to become profitable. This will be explored in the following.

3.5.2 Games of Chance

We use Theorem 3.60 to estimate the time needed to reach the winning threshold when using Λ_{ξ} in a game of chance. We assume a game (or a sequence of possibly correlated games) which allows a sequence of bets and observations. In step t we bet, depending on the history $x_{<t}$, a certain amount of money s_t , take some action y_t , observe outcome x_t , and receive reward r_t . Our profit, which we want to maximize, is $p_t = r_t - s_t$. The loss, which we want to minimize, can be defined as the negative profit, $\ell_{x_t y_t} = -p_t$. The probability of outcome x_t , possibly depending on the history

$x_{<t}$, is $\mu(x_t|x_{<t})$. The total μ -expected profit when using scheme Λ_ρ is $P_n^{\Lambda_\rho} = -L_n^{\Lambda_\rho}$. If we knew μ , the optimal strategy to maximize our expected profit is just Λ_μ . We assume $P_n^{\Lambda_\mu} > 0$ (otherwise there is no winning strategy at all, since $P_n^{\Lambda_\mu} \geq P_n^{\Lambda_\rho} \forall \rho$). Often we are not in the favorable position of knowing μ , but we know (or assume) that $\mu \in \mathcal{M}$ for some \mathcal{M} , for instance that μ is a computable probability distribution. From Theorem 3.60 we see that the average profit per round $\bar{p}_n^{\Lambda_\xi} := \frac{1}{n} P_n^{\Lambda_\xi}$ of the universal Λ_ξ scheme converges to the average profit per round $\bar{p}_n^{\Lambda_\mu} := \frac{1}{n} P_n^{\Lambda_\mu}$ of the optimal informed scheme, i.e. asymptotically we can make the same money even without knowing μ , by just using the universal Λ_ξ scheme. Theorem 3.60 allows us to lower bound the universal profit $P_n^{\Lambda_\xi}$

$$P_n^{\Lambda_\xi} \geq P_n^{\Lambda_\mu} - p_\Delta D_n - \sqrt{4(np_{max} - P_n^{\Lambda_\mu})p_\Delta D_n + p_\Delta^2 D_n^2} \quad (3.61)$$

where p_{max} is the maximal profit per round and p_Δ the profit range. The time needed for Λ_ξ to perform well can also be estimated. An interesting quantity is the expected number of rounds needed to reach the winning zone. Using $P_n^{\Lambda_\mu} > 0$ one can show that the r.h.s. of (3.61) is positive if, and only if

$$n > \frac{2p_\Delta(2p_{max} - \bar{p}_n^{\Lambda_\mu})}{(\bar{p}_n^{\Lambda_\mu})^2} \cdot D_n. \quad (3.62)$$

Theorem 3.63 (Time to Win) Let there be sequences $x_1 x_2 \dots$ over a finite alphabet \mathcal{X} drawn with probability $\mu(x_{1:n})$ for the first n symbols. In step t we make a bet, depending on the history $x_{<t}$, take some action y_t , and observe outcome x_t . Our net profit is $p_t \in [p_{max} - p_\Delta, p_{max}]$. The Λ_ρ -system (3.45) acts as to maximize the ρ -expected profit. $P_n^{\Lambda_\rho}$ is the total and $\bar{p}_n^{\Lambda_\rho} = \frac{1}{n} P_n^{\Lambda_\rho}$ is the average expected profit of the first n rounds. For the universal Λ_ξ and for the optimal informed Λ_μ prediction scheme the following holds:

$$\begin{aligned} i) \quad \bar{p}_n^{\Lambda_\xi} &= \bar{p}_n^{\Lambda_\mu} - O(n^{-1/2}) \longrightarrow \bar{p}_n^{\Lambda_\mu} \quad \text{for } n \rightarrow \infty \\ ii) \quad n &> \left(\frac{2p_\Delta}{\bar{p}_n^{\Lambda_\mu}} \right)^2 \cdot b_\mu \quad \wedge \quad \bar{p}_n^{\Lambda_\mu} > 0 \quad \implies \quad \bar{p}_n^{\Lambda_\xi} > 0 \end{aligned}$$

where $w_\mu = e^{-b_\mu}$ is the weight (3.5) of μ in ξ .

By dividing (3.61) by n and using $D_n \leq b_\mu$ (3.18) we see that the leading order of $\bar{p}_n^{\Lambda_\xi} - \bar{p}_n^{\Lambda_\mu}$ is bounded by $\sqrt{4p_\Delta p_{max} b_\mu / n}$, which proves (i). The condition in (ii) is actually a weakening of (3.62). $P_n^{\Lambda_\xi}$ is trivially positive for $p_{min} > 0$, since in this wonderful case *all* profits are positive. For negative p_{min} the condition of (ii) implies (3.62), since $p_\Delta > p_{max}$, and (3.62) implies positive (3.61), i.e. $P_n^{\Lambda_\xi} > 0$, which proves (ii).

If a winning strategy Λ with $\bar{p}_n^\Lambda > \varepsilon > 0$ exists, then Λ_ξ is asymptotically also a winning strategy with the same average profit.

3.5.3 Example

Let us consider a game with two dice, one with two black and four white faces, the other with four black and two white faces. The dealer who repeatedly throws the dice uses one or the other die according to some deterministic rule, which correlates the throws (e.g. the first die could be used in round t iff the t^{th} digit of π is 7). We can bet on black or white; the stake s is 3\$ in every round; our return r is 5\$ for every correct prediction.

The profit is $p_t = r\delta_{x_t y_t} - s$. The coloring of the dice and the selection strategy of the dealer unambiguously determine μ . $\mu(x_t | x_{<t})$ is $\frac{1}{3}$ or $\frac{2}{3}$ depending on which die has been chosen. One should bet on the more probable outcome ($\gamma = \frac{1}{2}$ in Section 3.4.3). If we knew μ the expected profit per round would be $\bar{p}_n^{\Lambda_\mu} = p_n^{\Lambda_\mu} = \frac{2}{3}r - s = \frac{1}{3}\$ > 0$. If we don't know μ we should use Solomonoff's universal prior with $D_n \leq b_\mu = K(\mu) \cdot \ln 2$, where $K(\mu)$ is the length of the shortest program coding μ (see Subsection 3.2.9). Then we know that betting on the outcome with higher ξ probability leads asymptotically to the same profit (Theorem 3.63(i)) and Λ_ξ reaches the winning threshold no later than $n_{\text{thresh}} = 900 \ln 2 \cdot K(\mu)$ (Theorem 3.63(ii)) or sharper $n_{\text{thresh}} = 330 \ln 2 \cdot K(\mu)$ from (3.62), where $p_{\max} = r - s = 2\$$ and $p_\Delta = r = 5\$$ have been used.

If the die selection strategy reflected in μ is not too complicated, the Λ_ξ prediction system reaches the winning zone after a few thousand rounds. The number of rounds is not really small because the expected profit per round is one order of magnitude smaller than the return. This leads to a constant of two orders of magnitude size in front of $K(\mu)$. Stated otherwise, it is due to the large stochastic noise, which makes it difficult to extract the signal, i.e. the structure of the rule μ (see next subsection). Furthermore, this is only a bound for the turnaround value of t_{thresh} . The true expected turnaround t might be smaller. However, every game for which there exists a computable winning strategy with $\bar{p}_{n\rho} > \varepsilon > 0$, Λ_ξ is guaranteed to get into the winning zone for some $t \sim K(\mu)$.

3.5.4 Information-Theoretic Interpretation

We try to give an intuitive explanation of Theorem 3.63(ii). We know that $\xi(x_t | x_{<t})$ converges to $\mu(x_t | x_{<t})$ for $t \rightarrow \infty$. In a sense Λ_ξ learns μ from past data $x_{<t}$. The information content in μ relative to ξ is $\ln 2 \cdot D_\infty \leq b_\mu \cdot \ln 2$. One might think of a Shannon-Fano prefix code of $\nu \in \mathcal{M}$ of length $\lceil b_\nu \cdot \ln 2 \rceil$, which exists since the Kraft inequality $\sum_\nu 2^{-\lceil b_\nu \cdot \ln 2 \rceil} \leq \sum_\nu w_\nu \leq 1$ is satisfied. $b_\mu \cdot \ln 2$ bits have to be learned before Λ_ξ can be as good as Λ_μ . In the worst case, the only information conveyed by x_t is in form of the received profit p_t . Remember that we always know the profit p_t before the next cycle starts.

Assume that the distribution of the profits in the interval $[p_{\min}, p_{\max}]$ is mainly due to noise, and there is only a small informative signal of amplitude $\bar{p}_n^{\Lambda_\mu}$. To reliably determine the sign of a signal of amplitude $\bar{p}_n^{\Lambda_\mu}$, disturbed by noise of amplitude

p_Δ , we have to resubmit a bit $O((p_\Delta/\bar{p}_n^{\Lambda_\mu})^2)$ times (this reduces the standard deviation below the signal amplitude $\bar{p}_n^{\Lambda_\mu}$). To learn μ , $b_\mu \ln 2$ bits have to be transmitted, which requires $n \geq O((p_\Delta/\bar{p}_n^{\Lambda_\mu})^2) \cdot b_\mu \ln 2$ cycles. This expression coincides with the condition in (ii). Identifying the signal amplitude with $\bar{p}_n^{\Lambda_\mu}$ is the weakest part of this consideration, as we have no argument why this should be true. It may be interesting to make the analogy more rigorous, which may also lead to a simpler proof of (ii) not based on Theorems 3.48 and 3.60 with their rather complex proofs.

3.6 Optimality Properties

3.6.1 Lower Error Bound

We want to show that there exists a class \mathcal{M} of distributions such that *any* predictor Θ ignorant of the distribution $\mu \in \mathcal{M}$ from which the observed sequence is sampled must make some minimal additional number of errors as compared to the best informed predictor Θ_μ .

For deterministic environments a lower bound can easily be obtained by a combinatoric argument. Consider a class \mathcal{M} containing 2^n binary sequences such that each prefix of length n occurs exactly once. Assume any deterministic predictor Θ (not knowing the sequence in advance), then for every prediction x_t^Θ of Θ at times $t \leq n$ there exists a sequence with opposite symbol $x_t = 1 - x_t^\Theta$. Hence, $E_\infty^\Theta \geq E_n^\Theta = n = \log_2 |\mathcal{M}|$ is a lower worst case bound for every predictor Θ , (this includes Θ_ξ , of course). This shows that the upper bound $E_\infty^{\Theta_\xi} \leq \log_2 |\mathcal{M}|$ for uniform w obtained in the discussion after Theorem 3.36 is sharp. In the general probabilistic case we can show by a similar argument that the upper bound of Theorem 3.36 is sharp for θ_ξ and “static” predictors, and sharp within a factor of 2 for general predictors. We do not know whether the factor two gap can be closed.

Theorem 3.64 (Lower Error Bound) For every n there is an \mathcal{M} and $\mu \in \mathcal{M}$ and weights w_ν such that

$$(i) \quad e_t^{\Theta_\xi} - e_t^{\Theta_\mu} = \sqrt{2s_t} \quad \text{and} \quad E_n^{\Theta_\xi} - E_n^{\Theta_\mu} = S_n + \sqrt{4E_n^{\Theta_\mu} S_n + S_n^2}$$

where $E_n^{\Theta_\xi}$ and $E_n^{\Theta_\mu}$ are the total expected number of errors of Θ_ξ and Θ_μ , and s_t and S_n are defined in (3.14). More generally, the equalities hold for *any* “static” deterministic predictor θ for which y_t^θ is independent of $x_{<t}$. For every n and *arbitrary* deterministic predictor Θ , there exists an \mathcal{M} and $\mu \in \mathcal{M}$ such that

$$(ii) \quad e_t^\Theta - e_t^{\Theta_\mu} \geq \frac{1}{2} \sqrt{2s_t(x_{<t})} \quad \text{and} \quad E_n^\Theta - E_n^{\Theta_\mu} \geq \frac{1}{2} [S_n + \sqrt{4E_n^{\Theta_\mu} S_n + S_n^2}]$$

Proof. (i) The proof parallels and generalizes the deterministic case. Consider a

class \mathcal{M} of 2^n distributions (over binary alphabet) indexed by $a \equiv a_1 \dots a_n \in \{0,1\}^n$. For each t we want a distribution with posterior probability $\frac{1}{2}(1+\varepsilon)$ for $x_t=1$ and one with posterior probability $\frac{1}{2}(1-\varepsilon)$ for $x_t=1$ independent of the past $x_{<t}$ with $0 < \varepsilon \leq \frac{1}{2}$. That is

$$\mu_a(x_1 \dots x_n) = \mu_{a_1}(x_1) \cdot \dots \cdot \mu_{a_n}(x_n), \quad \text{where} \quad \mu_{a_t}(x_t) = \begin{cases} \frac{1}{2}(1+\varepsilon) & \text{for } a_t = x_t \\ \frac{1}{2}(1-\varepsilon) & \text{for } a_t \neq x_t \end{cases}$$

We are not interested in predictions beyond time n but for completeness we may define μ_a to assign probability 1 to $x_t=1$ for all $t > n$. If $\mu = \mu_a$, the informed scheme Θ_μ always predicts the bit which has highest μ probability, i.e. $y_t^{\Theta_\mu} = a_t$

$$\implies e_t^{\Theta_\mu} = 1 - \mu_{a_t}(y_t^{\Theta_\mu}) = \frac{1}{2}(1-\varepsilon) \implies E_n^{\Theta_\mu} = \frac{n}{2}(1-\varepsilon).$$

Since $E_n^{\Theta_\mu}$ is the same for all a we seek to maximize E_n^Θ for a given predictor Θ in the following. Assume Θ predicts y_t^Θ (independent of history $x_{<t}$). Since we want lower bounds we seek for a worst case μ . A success $y_t^\Theta = x_t$ has lowest possible probability $\frac{1}{2}(1-\varepsilon)$ if $a_t = 1 - y_t^\Theta$.

$$\implies e_t^\Theta = 1 - \mu_{a_t}(y_t^\Theta) = \frac{1}{2}(1+\varepsilon) \implies E_n^\Theta = \frac{n}{2}(1+\varepsilon).$$

So we have $e_t^\Theta - e_t^{\Theta_\mu} = \varepsilon$ and $E_n^\Theta - E_n^{\Theta_\mu} = n\varepsilon$ for the regrets. We need to eliminate n and ε in favor of s_t , S_n , and $E_n^{\Theta_\mu}$. If we assume uniform weights $w_{\mu_a} = 2^{-n}$ for all μ_a we get

$$\xi(x_{1:n}) = \sum_a w_{\mu_a} \mu_a(x_{1:n}) = 2^{-n} \prod_{t=1}^n \sum_{a_t \in \{0,1\}} \mu_{a_t}(x_t) = 2^{-n} \prod_{t=1}^n 1 = 2^{-n},$$

i.e. ξ is an unbiased Bernoulli sequence ($\xi(x_t|x_{<t}) = \frac{1}{2}$).

$$\implies s_t(x_{<t}) = \sum_{x_t} \left(\frac{1}{2} - \mu_{a_t}(x_t)\right)^2 = \frac{1}{2}\varepsilon^2 \quad \text{and} \quad S_n = \frac{n}{2}\varepsilon^2.$$

So we have $\varepsilon = \sqrt{2s_t}$ which proves the instantaneous regret formula $e_t^\Theta - e_t^{\Theta_\mu} = \sqrt{2s_t}$ for static Θ . Inserting $\varepsilon = \sqrt{\frac{2}{n}S_n}$ into $E_n^{\Theta_\mu}$ and solving w.r.t. $\sqrt{2n}$ we get $\sqrt{2n} = \sqrt{S_n} + \sqrt{4E_n^{\Theta_\mu} + S_n}$. So we finally get

$$E_n^\Theta - E_n^{\Theta_\mu} = n\varepsilon = \sqrt{S_n}\sqrt{2n} = S_n + \sqrt{4E_n^{\Theta_\mu}S_n + S_n^2}$$

which proves the total regret formula in (ii) for static Θ . We can choose⁸ $y_t^{\Theta_\xi} \equiv 0$ to be a static predictor. Together this shows (i).

(ii) For non-static predictors, $a_t = 1 - y_t^\Theta$ in the proof of (i) depends on $x_{<t}$, which is not allowed. For general, but fixed a_t we have $e_t(x_{<t}) = 1 - \mu_{a_t}(y_t^\Theta)$. This quantity

⁸This choice may be made unique by slightly ($\ll \varepsilon$) non-uniformizing w_{μ_a} .

may assume any value between $\frac{1}{2}(1-\varepsilon)$ and $\frac{1}{2}(1+\varepsilon)$, when averaged over $x_{<t}$, and is, hence of little direct help. But if we average the result also over all environments μ_a , we get

$$\langle E_n^\Theta \rangle_a = \langle \sum_{t=1}^n \mathbf{E}[e_t^\Theta(x_{<t})] \rangle_a = \sum_{t=1}^n \mathbf{E}[\langle e_t^\Theta(x_{<t}) \rangle_a] = \sum_{t=1}^n \mathbf{E}[\frac{1}{2}] = \frac{1}{2}n$$

whatever Θ is chosen: a sort of No-Free-Lunch theorem [WM97], stating that on average all predictors perform equally well/bad. The expectation of E_n^Θ w.r.t. a can only be $\frac{1}{2}n$ if there exist a , such that $E_n^\Theta \geq \frac{1}{2}n$. Fixing such an a and choosing $\mu = \mu_a$ we get $E_n^\Theta - E_n^{\Theta_\mu} \geq \frac{1}{2}n\varepsilon = \dots = \frac{1}{2}[S_n + \sqrt{4E_n^{\Theta_\mu}S_n + S_n^2}]$, and similarly $e_n^\Theta - e_n^{\Theta_\mu} \geq \frac{1}{2}\varepsilon = \dots = \frac{1}{2}\sqrt{2s_t(x_{<t})}$. \square

Since $d_t/s_t = 1 + O(\varepsilon^2)$ we have $D_n/S_n \rightarrow 1$ for $\varepsilon \rightarrow 0$. Hence the error bound of Theorem 3.36 with S_n replaced by D_n is asymptotically tight for $E_n^{\Theta_\mu}/D_n \rightarrow \infty$ (which implies $\varepsilon \rightarrow 0$). This shows that without restrictions on the loss function which exclude the error-loss, the loss bound in Theorem 3.48 can also not be improved. Note that the bounds are tight even when \mathcal{M} is restricted to Markov or i.i.d. environments, since the presented counterexample is i.i.d. Finally, $E_n^\Theta - E_n^{\Theta_\mu} = n\varepsilon = n\sqrt{\frac{2S_n}{n}} \rightarrow \sqrt{2nD_n}$, which shows that the bound (3.50) of Merhav and Feder is also asymptotically tight.

A set \mathcal{M} independent of leading to a good (but not tight) lower bound is $\mathcal{M} = \{\mu_1, \mu_2\}$ with $\mu_{1/2}(1|x_{<t}) = \frac{1}{2} \pm \varepsilon_t$ with $\varepsilon_t = \min\{\frac{1}{2}, \sqrt{\ln w_{\mu_1}^{-1}}/\sqrt{t} \ln t\}$. For $w_{\mu_1} \ll w_{\mu_2}$ and $n \rightarrow \infty$ one can show that $E_n^{\Theta_\xi} - E_n^{\Theta_{\mu_1}} \sim \frac{1}{\ln n} \sqrt{E_n^{\Theta_{\mu_1}} \ln w_{\mu_1}^{-1}}$ (Problem 3.6).

Unfortunately there are many important special cases for which the loss bound (3.48) is not tight. For continuous \mathcal{Y} and logarithmic or quadratic loss function, for instance, we have seen that the regret $L_\infty^{\Lambda_\xi} - L_\infty^{\Lambda_\mu} \leq \ln w_\mu^{-1} < \infty$ is finite. For arbitrary loss function, but μ bounded away from certain critical values, the regret is also finite. For instance, consider the special error-loss, binary alphabet, and $|\mu(x_t|x_{<t}) - \frac{1}{2}| > \varepsilon$ for all t and x . Θ_μ predicts 0 if $\mu(0|x_{<t}) < \frac{1}{2}$. If also $\xi(0|x_{<t}) < \frac{1}{2}$, then Θ_ξ makes the same prediction as Θ_μ , for $\xi(0|x_{<t}) > \frac{1}{2}$ the predictions differ. In the latter case $|\xi(0|x_{<t}) - \mu(0|x_{<t})| > \varepsilon$. Conversely for $\mu(0|x_{<t}) > \frac{1}{2}$. So in any case $e_t^{\Theta_\xi} - e_t^{\Theta_\mu} \leq \frac{1}{\varepsilon^2} [\xi(x_t|x_{<t}) - \mu(x_t|x_{<t})]^2$. Using Definition 3.34 and Theorem 3.19(i) we see that $E_\infty^{\Theta_\xi} - E_\infty^{\Theta_\mu} \leq \frac{1}{\varepsilon^2} \ln w_\mu^{-1} < \infty$ is finite too. Nevertheless, Theorem 3.64 is important as it tells us that bound (3.48) can only be strengthened by making further assumptions on ℓ or \mathcal{M} .

3.6.2 Pareto Optimality of ξ

In this subsection we want to establish a different kind of optimality property of ξ . Let $\mathcal{F}(\mu, \rho)$ be any of the performance measures of ρ relative to μ considered in the previous sections (e.g. s_t , or D_n , or L_n , ...). It is easy to find ρ more tailored towards μ such that $\mathcal{F}(\mu, \rho) < \mathcal{F}(\mu, \xi)$. This improvement may be achieved by increasing w_μ ,

but probably at the expense of increasing \mathcal{F} for other ν , i.e. $\mathcal{F}(\nu, \rho) > \mathcal{F}(\nu, \xi)$ for some $\nu \in \mathcal{M}$. Since we do not know μ in advance we may ask whether there exists a ρ with better or equal performance for *all* $\nu \in \mathcal{M}$ and a strictly better performance for one $\nu \in \mathcal{M}$. This would clearly render ξ suboptimal w.r.t. to \mathcal{F} . We show that there is no such ρ for all performance measures studied in this work.

Definition 3.65 (Pareto Optimality) Let $\mathcal{F}(\mu, \rho)$ be any performance measure of ρ relative to μ . The universal prior ξ is called Pareto-optimal w.r.t. \mathcal{F} if there is no ρ with $\mathcal{F}(\nu, \rho) \leq \mathcal{F}(\nu, \xi)$ for all $\nu \in \mathcal{M}$ and strict inequality for at least one ν .

Theorem 3.66 (Pareto Optimality) The universal prior ξ is Pareto-optimal w.r.t. the instantaneous and total squared distances s_t and S_n (3.14), entropy distances d_t and D_n (3.16), errors e_t and E_n (3.34), and losses l_t and L_n (3.46).

Proof. We first prove Theorem 3.66 for the instantaneous expected loss l_t . We need the more general ρ expected instantaneous losses

$$l_{t\rho}^\Lambda(x_{<t}) := \sum_{x_t} \rho(x_t | x_{<t}) \ell_{x_t y_t}^\Lambda \quad (3.67)$$

for a predictor Λ . We want to arrive at a contradiction by assuming that ξ is not Pareto-optimal, i.e. by assuming the existence of a predictor⁹ Λ with $l_{t\nu}^\Lambda \leq l_{t\nu}^{\Lambda_\xi}$ for all $\nu \in \mathcal{M}$ and strict inequality for some ν . Implicit to this assumption is the assumption that $l_{t\nu}^\Lambda$ and $l_{t\nu}^{\Lambda_\xi}$ exist. $l_{t\nu}^\Lambda$ exists iff $\nu(x_t | x_{<t})$ exists iff $\nu(x_{<t}) > 0$ iff $w_\nu(x_{<t}) > 0$.

$$l_{t\xi}^\Lambda = \sum_{\nu} w_\nu(x_{<t}) l_{t\nu}^\Lambda < \sum_{\nu} w_\nu(x_{<t}) l_{t\nu}^{\Lambda_\xi} = l_{t\xi}^{\Lambda_\xi} \leq l_{t\xi}^\Lambda$$

The two equalities follow from inserting (3.7) into (3.67). The strict inequality follows from the assumption and $w_\nu(x_{<t}) > 0$. The last inequality follows from the fact that Λ_ξ minimizes by definition (3.45) the ξ -expected loss (similarly to (3.47)). The contradiction $l_{t\xi}^\Lambda < l_{t\xi}^{\Lambda_\xi}$ proves Pareto-optimality of ξ w.r.t. l_t .

In the same way we can prove Pareto-optimality of ξ w.r.t. the total loss L_n by defining the ρ expected total losses

$$L_{n\rho}^\Lambda := \sum_{t=1}^n \sum_{x_{<t}} \rho(x_{<t}) l_{t\rho}^\Lambda(x_{<t}) = \sum_{t=1}^n \sum_{x_{1:t}} \rho(x_{1:t}) \ell_{x_t y_t}^\Lambda \quad (3.68)$$

for a predictor Λ , and by assuming $L_{n\nu}^\Lambda \leq L_{n\nu}^{\Lambda_\xi}$ for all ν and strict inequality for some ν , from which we get the contradiction $L_{n\xi}^\Lambda = \sum_{\nu} w_\nu L_{n\nu}^\Lambda < \sum_{\nu} w_\nu L_{n\nu}^{\Lambda_\xi} = L_{n\xi}^{\Lambda_\xi} \leq L_{n\xi}^\Lambda$ with

⁹According to definition 3.65 we should look for a ρ , but for each deterministic predictor Λ there exists a ρ with $\Lambda = \Lambda_\rho$.

the help of (3.5). The instantaneous and total expected errors e_t and E_n can be considered as special loss functions.

Pareto-optimality of ξ w.r.t. s_t (and hence S_n) can be understood from geometrical insight. A formal proof for s_t goes as follows: With the abbreviations $i = x_t$, $y_{\nu i} = \nu(x_t | x_{<t})$, $z_i = \xi(x_t | x_{<t})$, $r_i = \rho(x_t | x_{<t})$, and $w_\nu = w_\nu(x_{<t}) \geq 0$ we ask for a vector \mathbf{r} with $\sum_i (y_{\nu i} - r_i)^2 \leq \sum_i (y_{\nu i} - z_i)^2 \forall \nu$. This implies

$$\begin{aligned} 0 &\geq \sum_\nu w_\nu \left[\sum_i (y_{\nu i} - r_i)^2 - \sum_i (y_{\nu i} - z_i)^2 \right] = \sum_\nu w_\nu \left[\sum_i -2y_{\nu i} r_i + r_i^2 + 2y_{\nu i} z_i - z_i^2 \right] \\ &= \sum_i -2z_i r_i + r_i^2 + 2z_i z_i - z_i^2 = \sum_i (r_i - z_i)^2 \geq 0 \end{aligned}$$

where we have used $\sum_\nu w_\nu = 1$ and $\sum_\nu w_\nu y_{\nu i} = z_i$ (3.7). $0 \geq \sum_i (r_i - z_i)^2 \geq 0$ implies $\mathbf{r} = \mathbf{z}$ proving unique Pareto-optimality of ξ w.r.t. s_t . Similarly for d_t the assumption $\sum_i y_{\nu i} \ln \frac{y_{\nu i}}{r_i} \leq \sum_i y_{\nu i} \ln \frac{y_{\nu i}}{z_i} \forall \nu$ implies

$$0 \geq \sum_\nu w_\nu \left[\sum_i y_{\nu i} \ln \frac{y_{\nu i}}{r_i} - y_{\nu i} \ln \frac{y_{\nu i}}{z_i} \right] = \sum_\nu w_\nu \sum_i y_{\nu i} \ln \frac{z_i}{r_i} = \sum_i z_i \ln \frac{z_i}{r_i} \geq 0$$

which implies $\mathbf{r} = \mathbf{z}$ proving unique Pareto-optimality of ξ w.r.t. d_t . The proofs for S_n and D_n are similar. \square

We have proven that ξ is *uniquely* Pareto-optimal w.r.t. s_t , S_n , d_t and D_n . In the case of e_t , E_n , l_t and L_n there are other $\rho \neq \xi$ with $\mathcal{F}(\nu, \rho) = \mathcal{F}(\nu, \xi) \forall \nu$, but the actions/predictions they invoke are unique ($y_t^{\Lambda_\rho} = y_t^{\Lambda_\xi}$) (if ties in $\arg\max_{y_t}$ are broken in a consistent way), and this is all that counts.

Note that ξ is *not* Pareto-optimal w.r.t. to all thinkable performance measures. Counterexamples can be given for $\mathcal{F}(\nu, \xi) = \sum_{x_t} |\nu(x_t | x_{<t}) - \xi(x_t | x_{<t})|^\alpha$ for $\alpha \neq 2$ (see Problem 3.5). Nevertheless, for all measures which are relevant from a decision theoretic point of view, i.e. for all loss functions l_t and L_n , ξ has the welcomed property of being Pareto-optimal.

3.6.3 Balanced Pareto Optimality of ξ

Pareto-optimality should be regarded as a necessary condition for a prediction scheme aiming to be optimal. From a practical point of view a significant decrease of \mathcal{F} for many ν may be desirable even if this causes a small increase of \mathcal{F} for a few other ν . The impossibility of such a “balanced” improvement is a more demanding condition on ξ than pure Pareto-optimality. The next theorem shows that Λ_ξ is also balanced-Pareto-optimal. We only consider the performance measure L_n and suppress the index n for convenience.

Theorem 3.69 (Balanced Pareto Optimality w.r.t. L)

$$\Delta_\nu := L_\nu^{\tilde{\Lambda}} - L_\nu^{\Lambda_\xi}, \quad \Delta := \sum_{\nu \in \mathcal{M}} w_\nu \Delta_\nu \quad \Rightarrow \quad \Delta \geq 0.$$

This implies the following: Assume $\tilde{\Lambda}$ has larger loss than Λ_ξ on environments \mathcal{L} by a total weighted amount of $\Delta_{\mathcal{L}} := \sum_{\lambda \in \mathcal{L}} w_\lambda \Delta_\lambda$. Then $\tilde{\Lambda}$ can have smaller loss on $\eta \in \mathcal{H} := \mathcal{M} \setminus \mathcal{L}$, but the improvement is bounded by $\Delta_{\mathcal{H}} := |\sum_{\eta \in \mathcal{H}} w_\eta \Delta_\eta| \leq \Delta_{\mathcal{L}}$. Especially $|\Delta_\eta| \leq w_\eta^{-1} \max_{\lambda \in \mathcal{L}} \Delta_\lambda$.

This means that a weighted loss decrease $\Delta_{\mathcal{H}}$ by using $\tilde{\Lambda}$ instead of Λ_ξ is compensated by an at least as large weighted increase $\Delta_{\mathcal{L}}$ on other environments. If the increase is small, the decrease can also only be small. In the special case of only a single environment with increased loss Δ_λ , the decrease is bound by $\Delta_\eta \leq \frac{w_\lambda}{w_\eta} |\Delta_\lambda|$, i.e. an increase by an amount Δ_λ can only cause a decrease by at most the same amount times a factor $\frac{w_\lambda}{w_\eta}$. A increase can only cause a smaller decrease in simpler environments, but a scaled decrease in more complex environments. Finally note that pure Pareto-optimality (3.66) follows from balanced Pareto-optimality in the special case of no increase $\Delta_{\mathcal{L}} \equiv 0$.

Proof. $\Delta \geq 0$ follows from $\Delta = \sum_\nu w_\nu [L_\nu^{\tilde{\Lambda}} - L_\nu^{\Lambda_\xi}] = L_\xi^{\tilde{\Lambda}} - L_\xi^{\Lambda_\xi} \geq 0$, where we have used linearity of L_ρ in ρ and $L_\xi^{\Lambda_\xi} \leq L_\xi^{\tilde{\Lambda}}$. The remainder of Theorem 3.69 is obvious from $0 \leq \Delta = \Delta_{\mathcal{L}} - \Delta_{\mathcal{H}}$ and by bounding the weighted average Δ_η by its maximum. \square

The term *Pareto-optimal* has been taken from the economics literature, but there is the closely related notion of unimprovable strategies [BM98] or admissible estimators [Fer67] in statistics for parameter estimation, for which results similar to Theorem 3.66 exist. Furthermore, it would be interesting to show under which conditions, the class of *all* Bayes-mixtures (i.e. with all possible values for the weights) is complete in the sense that *every* Pareto-optimal strategy can be based on a Bayes-mixtures. Pareto-optimality is sort of a minimal demand on a prediction scheme aiming to be optimal. A scheme which is not even Pareto-optimal cannot be regarded as optimal in any reasonable sense. Pareto-optimality of ξ w.r.t. most performance measures emphasizes the distinctiveness of Bayes-mixture strategies.

3.6.4 On the Optimal Choice of Weights

In the following we indicate the dependency of ξ on w explicitly by writing ξ_w . We have shown that the Λ_{ξ_w} prediction schemes are (balanced) Pareto optimal, i.e. that *no* prediction scheme Λ (whether based on a Bayes mix or not) can be uniformly better. Least assumptions on the environment are made for \mathcal{M} which are as large as possible. In Section 2.4 we have discussed the set \mathcal{M} of all enumerable semimeasures which we regarded as sufficiently large from a computational point

of view (see [Sch02a] for even larger sets, but which are still in the computational realm). Agreeing on this \mathcal{M} still leaves open the question of how to choose the weights (prior beliefs) w_ν , since every ξ_w with $w_\nu > 0 \forall \nu$ is Pareto-optimal and leads asymptotically to optimal predictions.

We have derived bounds for the mean squared sum $S_{n\nu}^{\xi_w} \leq \ln w_\nu^{-1}$ and for the loss regret $L_{n\nu}^{\Lambda_{\xi_w}} - L_{n\nu}^{\Lambda_\nu} \leq 2 \ln w_\nu^{-1} + 2\sqrt{\ln w_\nu^{-1} L_{n\nu}^{\Lambda_\nu}}$. All bounds monotonically decrease with increasing w_ν . So it is desirable to assign high weights to all $\nu \in \mathcal{M}$. Due to the (semi)probability constraint $\sum_\nu w_\nu \leq 1$ one has to find a compromise. In the following we will argue that in the class of enumerable weight functions with short program there is an optimal compromise, namely $w_\nu = 2^{-K(\nu)}$ which gives Solomonoff's prior.

Consider the class of enumerable weight functions with short programs, namely $\mathcal{V} := \{v_{(\cdot)} : \mathcal{M} \rightarrow \mathbb{R}^+ \text{ with } \sum_\nu v_\nu \leq 1 \text{ and } K(v) = O(1)\}$. Let $w_\nu := 2^{-K(\nu)}$ and $v_{(\cdot)} \in \mathcal{V}$. Corollary 4.3.1 of [LV97, p255] says that $K(x) \leq -\log_2 P(x) + K(P) + O(1)$ for all x if P is an enumerable discrete semimeasure. Identifying P with v and x with (the program index describing) ν we get

$$\ln w_\nu^{-1} \leq \ln v_\nu^{-1} + O(1).$$

This means that the bounds for ξ_w depending on $\ln w_\nu^{-1}$ are at most $O(1)$ larger than the bounds for ξ_v depending on $\ln v_\nu^{-1}$. So we lose at most an additive constant of order 1 in the bounds when using ξ_w instead of ξ_v . In using Solomonoff's prior ξ_w we are on the safe side, getting (within $O(1)$) best bounds for *all* environments.

Theorem 3.70 (Optimality of universal weights) Within the set \mathcal{V} of enumerable weight functions with short program, the universal weights $w_\nu = 2^{-K(\nu)}$ lead to the smallest loss bounds within an additive (to $\ln w_\nu^{-1}$) constant in all enumerable environments.

Since the above justifies the use of Solomonoff's prior and Solomonoff's prior assigns high probability to an environment if and only if it has low (Kolmogorov) complexity, one may interpret the result as a justification of Occam's razor¹⁰. But note that this is more of a bootstrap argument, since we used Occam's razor in Section 2.1 to justify the restriction to enumerable semimeasures. We also considered only weight functions v with low complexity $K(v) = O(1)$. What did not enter as an assumption but came out as a result is that the specific universal weights $w_\nu = 2^{-K(\nu)}$ are optimal. See Problem 3.7 for a discussion of the (non)uniqueness of w_ν .

3.6.5 Occam's razor versus No Free Lunches

We do not regard Theorem 3.69 as a "No Free Lunch" (NFL) theorem [WM97]. Since most environments are completely random, a small concession on the loss in each of these completely uninteresting environments provides enough margin $\Delta_{\mathcal{H}}$

¹⁰The *only if* direction can be shown by a more easy and direct argument [Sch02a].

to yield distinguished performance on the few non-random (interesting) environments. Indeed, we would interpret the NFL theorems for optimization and search in [WM97] as balanced Pareto-optimality results. Interestingly, whereas for prediction only Bayes-mixes are Pareto-optimal, for search and optimization every algorithm is Pareto-optimal. There is an ongoing battle between believers in Occam's razor and believers in "no free lunches" that cannot be dealt with here [Sto01, SH02].

3.7 Miscellaneous

3.7.1 Multi-Step Predictions

Introduction. In multi-step prediction we want to predict $x_{t:n}$ from $x_{<t}$. For instance, every day, a weather forecaster in the morning of day t predicts the weather for the next 3 days t , $t+1$, and $t+2=n$. Up to now we have considered prediction problems with a lookahead of one time-step only: Given $x_{<t}$, predict x_t . Greedy minimization of the expected loss $l_t(x_{<t})$ at time t was optimal. Farther lookahead ($>t$) was not necessary, the reason being that the prediction/decision/action y_t has no influence on the environment μ . For acting agents, described in detail in later Chapters, multi-step lookahead is necessary for optimal actions. Another application of multi-step predictions is 'delayed sequence prediction', in which not the next, but next-to-next or h^{th} -next symbol shall be predicted.

Notation and basic relations. We are interested in multi-step posteriors $\rho(x_{t:n}|x_{<t}) = \rho(x_{1:n})/\rho(x_{<t})$, which generalize the one-step posteriors ($n=t$) considered so far. We abbreviate $\rho_{t:n}^{\dots} := \rho(x_{t:n}^{\dots}|x_{<t})$, where \dots are any superscripts (e.g. empty or '). We define the conditional probability vector $\vec{\rho}_{t:n} := \rho_{t:n}(\cdot|x_{<t}) \in \mathbb{R}^N$, where $N = |\mathcal{X}|^{n-t+1}$ and the i^{th} component of vector $\vec{\rho}_{t:n}$ is $\rho_{t:n}(i|x_{<t})$ with identification $\{1, \dots, N\} \ni i \hat{=} x_{t:n} \in \mathcal{X}^{n-t+1}$. Let $f \in \{a, b, d, h, s\}$ be any of the distances defined in (3.10), i.e. $f(\vec{y}, \vec{z}) = \sum_{i=1}^N \hat{f}(y_i, z_i)$ with $\hat{a}(y, z) = |y - z|$, $\hat{b}(y, z) = y |\ln \frac{y}{z}|$, $\hat{d}(y, z) = y \ln \frac{y}{z}$, $\hat{h}(y, z) = (\sqrt{y} - \sqrt{z})^2$, $\hat{s}(y, z) = (y - z)^2$. We define $f_{t:n}(x_{<t}) := f(\vec{\mu}_{t:n}, \vec{\xi}_{t:n})$, generalizing (3.13–3.17). The definitions $f_t(x_{<t}) := f_{t:t}(x_{<t})$ and $F_n := \sum_{t=1}^n \mathbf{E}[f_t]$, $F \in \{A, B, D, H, S\}$ are consistent with (3.13–3.17). Lemma 3.11 ($b-d \leq a \leq \sqrt{2d}$, $h \leq d$, $s \leq d$) implies $b_{t:n} - d_{t:n} \leq a_{t:n} \leq \sqrt{2d_{t:n}}$, $h_{t:n} \leq d_{t:n}$, $s_{t:n} \leq d_{t:n}$. For the relative entropy we have $D_n = d_{1:n}$ and¹¹ $\mathbf{E}_{t:k}[d_{k+1:n}] = d_{t:n} - d_{t:k}$ for $t \leq k < n$, which implies $\mathbf{E}[d_{t:n}] = D_n - D_{t-1} = \sum_{k=t}^n \mathbf{E}[d_k] \geq 0$ and $d_{t:n}$ is monotone increasing in n .

Convergence i.m.s. for bounded horizon. Henceforth, we don't need the Hellinger distance anymore, and we will use h for the horizon. Assume we want to predict the next h symbols, i.e. $n = n_t = t + h - 1$. We want to determine how fast $\xi'_{t:n_t} \equiv \xi(x'_{t:t+h-1}|x_{<t})$ converges to $\mu'_{t:n_t} \equiv \mu(x'_{t:t+h-1}|x_{<t})$. To prove convergence

¹¹ $\mathbf{E}_{t:k}[f(x_{1:k})] := \sum_{x_{t:k}} \mu_{t:k} f(x_{1:k})$, cf. (3.4).

i.m.s. we have to bound the expectation sum of $s_{t:n_t}$ or $a_{t:n_t}^2 \equiv (\sum_{x'_{t:n_t}} |\xi'_{t:n_t} - \mu'_{t:n_t}|)^2$:

$$\frac{1}{2} \sum_{t=1}^{\infty} \mathbf{E}[a_{t:n_t}^2] \leq \sum_{t=1}^{\infty} \mathbf{E}[d_{t:n_t}] = \sum_{t=1}^{\infty} \sum_{k=t}^{n_t} \mathbf{E}[d_k] \leq h \cdot \sum_{k=1}^{\infty} \mathbf{E}[d_k] = h \cdot D_{\infty} \leq h \cdot \ln w_{\mu}^{-1} < \infty. \quad (3.71)$$

In the second inequality we have used that the number of times $d_k \geq 0$ occurs for some k in the double sum is $\min\{h, k\} \leq h$. The bound implies $a_{t:n_t} \rightarrow 0$ i.m.s., which implies $\xi'_{t:n_t} \xrightarrow{t \rightarrow \infty} \mu'_{t:n_t}$ i.m.s. by dropping the sum in the definition of $a_{t:n_t}$. The bound loosens by a factor of h for h -step prediction as compared to 1-step prediction. The same bound holds for bounded horizon $h_t := n_t - t + 1 \leq h < \infty \forall t$ (increase $\sum_{k=t}^{n_t}$ to $\sum_{k=t}^{t+h-1}$), i.e.

$$\xi(x'_{t:n_t} | x_{<t}) \rightarrow \mu(x'_{t:n_t} | x_{<t}) \text{ i.m.s. for } t \rightarrow \infty \text{ if } h_t := n_t - t + 1 \leq h < \infty \forall t.$$

Delayed sequence prediction. A delayed feedback, where at time t , x is only known up to time $t-h$ for some delay h , is common in many practical problems. This is equivalent to predicting x_{n_t} from $x_{<t}$ with $n_t = t + h - 1$. The probability of x'_{n_t} , given $x_{<t}$, is $\rho(x'_{n_t} | x_{<t}) := \sum_{x'_{t:n_t-1}} \rho'_{t:n_t} = \sum_{x_{t:n_t-1}} \rho(x_{<n_t} x'_{n_t}) / \rho(x_{<t})$. Using $\sum_{x'_{n_t}} |\xi(x'_{n_t} | x_{<t}) - \mu(x'_{n_t} | x_{<t})| \leq \sum_{x'_{t:n_t}} |\xi'_{t:n_t} - \mu'_{t:n_t}| \equiv a_{t:n_t}$ and bound (3.71) we get

$$\sum_{t=1}^{\infty} \mathbf{E} \left(\sum_{x'_{n_t}} \left| \xi(x'_{n_t} | x_{<t}) - \mu(x'_{n_t} | x_{<t}) \right| \right)^2 \leq \sum_{t=1}^{\infty} \mathbf{E}[a_{t:n_t}^2] \leq 2h \cdot \ln w_{\mu}^{-1} < \infty$$

which implies $\xi(x'_{n_t} | x_{<t}) \xrightarrow{t \rightarrow \infty} \mu(x'_{n_t} | x_{<t})$ i.m.s. The loss bounds of Theorem 3.59(i,ii) also generalize to the delayed case. Loss bounds similar to Theorem 3.59(iii) and Theorem 3.48 should also be derivable.

Convergence i.m. for arbitrary horizon. Convergence i.m.s. does generally not hold for unbounded horizon h_t (see Problem 3.15). Remarkably, convergence (i.m.) holds nevertheless: For any limit path $n \geq t \rightarrow \infty$ we have

$$\lim_{t,n \rightarrow \infty} \mathbf{E}[d_{t:n}] = \lim_{t,n \rightarrow \infty} [D_n - D_{t-1}] = \lim_{n \rightarrow \infty} D_n - \lim_{t \rightarrow \infty} D_{t-1} = D_{\infty} - D_{\infty} = 0.$$

So for any $n_t = t + h_t - 1$ we have $\frac{1}{2} \mathbf{E}[a_{t:n_t}^2] \leq \mathbf{E}[d_{t:n_t}] \xrightarrow{t \rightarrow \infty} 0$, which implies

$$\xi(x'_{t:n_t} | x_{<t}) \xrightarrow[t \rightarrow \infty]{i.m.} \mu(x'_{t:n_t} | x_{<t}) \quad \text{and} \quad \xi(x'_{n_t} | x_{<t}) \xrightarrow[t \rightarrow \infty]{i.m.} \mu(x'_{n_t} | x_{<t}) \quad \text{for any } h_t.$$

Convergence i.m. is weaker than convergence w.p.1 and i.m.s. and is potentially slow. We expect cases where convergence is very slow when h_t grows very fast. We do not know whether convergence is reasonably fast for slowly growing horizon, e.g. for $h_t = \log t$ (or $h_t = t$).

3.7.2 Continuous Probability Classes \mathcal{M}

We have considered thus far countable probability classes \mathcal{M} , which makes sense from a computational point of view as emphasized in Subsection 3.2.9. On the other hand in statistical parameter estimation one often has a continuous hypothesis class (e.g. a Bernoulli(θ) process with unknown $\theta \in [0,1]$). Let

$$\mathcal{M} := \{\mu_\theta : \theta \in \Theta \subseteq \mathbb{R}^d\}$$

be a family of probability distributions parameterized by a d -dimensional continuous parameter θ . Let $\mu \equiv \mu_{\theta_0} \in \mathcal{M}$ be the true generating distribution and θ_0 be in the interior of the compact set Θ . We may restrict \mathcal{M} to a countable dense subset, like $\{\mu_\theta\}$ with computable (or rational) θ . If θ_0 is itself a computable real (or rational) then Theorem 3.60 applies. From a practical point of view the assumption of a computable θ_0 is not so serious. It is more from a traditional analysis point of view that one would like quantities and results depending smoothly on θ and not in a weird fashion depending on the computational complexity of θ . For instance, the weight $w(\theta)$ is often a continuous probability density

$$\xi(x_{1:n}) := \int_{\Theta} d\theta w(\theta) \cdot \mu_\theta(x_{1:n}), \quad \int_{\Theta} d\theta w(\theta) = 1, \quad w(\theta) \geq 0. \quad (3.72)$$

The most important property of ξ used in this work was $\xi(x_{1:n}) \geq w_\nu \cdot \nu(x_{1:n})$ which has been obtained from (3.5) by dropping the sum over ν . The analogous construction here is to restrict the integral over Θ to a small vicinity N_δ of θ . For sufficiently smooth μ_θ and $w(\theta)$ we expect $\xi(x_{1:n}) \gtrsim |N_{\delta_n}| \cdot w(\theta) \cdot \mu_\theta(x_{1:n})$, where $|N_{\delta_n}|$ is the volume of N_{δ_n} . This in turn leads to $D_n \lesssim \ln w_\mu^{-1} + \ln |N_{\delta_n}|^{-1}$, where $w_\mu := w(\theta_0)$. N_{δ_n} should be the largest possible region in which $\ln \mu_\theta$ is approximately flat on average. The averaged instantaneous, mean, and total curvature matrices of $\ln \mu$ are

$$\begin{aligned} j_t(x_{<t}) &:= \mathbf{E}_t[\nabla_\theta \ln \mu_\theta(x_t|x_{<t}) \nabla_\theta^T \ln \mu_\theta(x_t|x_{<t})]_{|\theta=\theta_0}, & \bar{J}_n &:= \frac{1}{n} J_n \\ J_n &:= \sum_{t=1}^n \mathbf{E}[j_t(x_{<t})] = \mathbf{E}[\nabla_\theta \ln \mu_\theta(x_{1:n}) \nabla_\theta^T \ln \mu_\theta(x_{1:n})]_{|\theta=\theta_0} \end{aligned} \quad (3.73)$$

They are the Fisher information of μ and may be viewed as measures of the parametric complexity of μ_θ at $\theta = \theta_0$. The last equality can be shown by using the fact that the μ -expected value of $\nabla \ln \mu \cdot \nabla^T \ln \mu$ coincides with $-\nabla \nabla^T \ln \mu$ (since \mathcal{X} is finite) and a similar line of reasoning as in (3.18) for D_n .

Theorem 3.74 (Continuous Entropy Bound) Let μ_θ be twice continuously differentiable at $\theta_0 \in \Theta \subseteq \mathbb{R}^d$ and $w(\theta)$ be continuous and positive at θ_0 . Furthermore we assume that the inverse of the mean Fisher information matrix $(\bar{j}_n)^{-1}$ exists, is bounded for $n \rightarrow \infty$, and is uniformly (in n) continuous at θ_0 . Then the relative Entropy D_n between $\mu \equiv \mu_{\theta_0}$ and ξ (defined in (3.72)) can be bounded by

$$D_n := \mathbf{E} \ln \frac{\mu(x_{1:n})}{\xi(x_{1:n})} \leq \ln w_\mu^{-1} + \frac{d}{2} \ln \frac{n}{2\pi} + \frac{1}{2} \ln \det \bar{j}_n + o(1) =: b_\mu$$

where $w_\mu \equiv w(\theta_0)$ is the weight density (3.72) of μ in ξ and $o(1)$ tends to zero for $n \rightarrow \infty$.

For independent and identically distributed distributions $\mu_\theta(x_{1:n}) = \mu_\theta(x_1) \cdots \mu_\theta(x_n) \forall \theta$ this bound has been proven in [CB90, Theorem 2.3]. In this case $J^{[CB90]}(\theta_0) \equiv \bar{j}_n \equiv j_n$ independent of n . For stationary (k^{th} -order) Markov processes \bar{j}_n is also constant. The proof generalizes to arbitrary μ_θ by replacing $J^{[CB90]}(\theta_0)$ with \bar{j}_n everywhere in their proof. For the proof to go through, the vicinity $N_{\delta_n} := \{\theta : \|\theta - \theta_0\|_{\bar{j}_n} \leq \delta_n\}$ of θ_0 must contract to a point set $\{\theta_0\}$ for $n \rightarrow \infty$ and $\delta_n \rightarrow 0$. \bar{j}_n is always positive semi-definite as can be seen from the definition. The boundedness condition of \bar{j}_n^{-1} implies a strictly positive lower bound independent of n on the Eigenvalues of \bar{j}_n for all sufficiently large n , which ensures $N_{\delta_n} \rightarrow \{\theta_0\}$. The uniform continuity of \bar{j}_n ensures that the remainder $o(1)$ from the Taylor expansion of D_n is independent of n . Note that twice continuous differentiability of D_n at θ_0 [CB90, Condition 2] follows for finite \mathcal{X} from twice continuous differentiability of μ_θ . Under some additional technical conditions one can even prove an equality $D_n = \ln w_\mu^{-1} + \frac{d}{2} \ln \frac{n}{2\pi e} + \frac{1}{2} \ln \det \bar{j}_n + o(1)$ for the i.i.d. case [CB90, (1.4)], which is probably also valid for general μ .

The $\ln w_\mu^{-1}$ part in the bound is the same as for countable \mathcal{M} . The $\frac{d}{2} \ln \frac{n}{2\pi}$ can be understood as follows: Consider $\theta \in [0,1)$ and restrict the continuous \mathcal{M} to θ which are finite binary fractions. Assign a weight $w(\theta) \approx 2^{-l}$ to a θ with binary representation of length l . $D_n \lesssim l \cdot \ln 2$ in this case. But what if θ is not a finite binary fraction? A continuous parameter can typically be estimated with accuracy $O(n^{-1/2})$ after n observations. The data do not allow to distinguish a $\tilde{\theta}$ from the true θ if $|\tilde{\theta} - \theta| < O(n^{-1/2})$. There is such a $\tilde{\theta}$ with binary representation of length $l = \log_2 O(\sqrt{n})$. Hence we expect $D_n \lesssim \frac{1}{2} \ln n + O(1)$ or $\frac{d}{2} \ln n + O(1)$ for a d -dimensional parameter space. In general, the $O(1)$ term depends on the parametric complexity of μ_θ and is explicated by the third $\frac{1}{2} \ln \det \bar{j}_n$ term in Theorem 3.74. See [CB90, p454] for an alternative explanation. Note that a uniform weight $w(\theta) = \frac{1}{|\Theta|}$ does not lead to a uniform bound unlike the discrete case. A uniform bound is obtained for Bernardo's (or in the scalar case Jeffreys') reference prior $w(\theta) \sim \sqrt{\det \bar{j}_\infty(\theta)}$ if j_∞ exists [Ris96].

So Theorems 3.19...3.60 are also applicable to the case of continuously parameterized probability classes. Theorem 3.74 is also valid for a mixture of the discrete and continuous case $\xi = \sum_a \int d\theta w^a(\theta) \mu_\theta^a$ with $\sum_a \int d\theta w^a(\theta) = 1$.

3.7.3 Further Applications

Partial sequence prediction. There are (at least) two ways to treat partial sequence prediction. With this we mean that not every symbol of the sequence need to be predicted, say given sequences of the form $z_1x_1\dots z_nx_n$ we want to predict the x 's only. The first way is to keep the Λ_ρ prediction schemes of the last sections mainly as they are, and use a time dependent loss function, which assigns zero loss $\ell_{zy}^t \equiv 0$ at the z positions. Any dummy prediction y is then consistent with (3.45). The losses for predicting x are generally non-zero. This solution is satisfactory as long as the z 's are drawn from a probability distribution. The second (preferable) way does not rely on a probability distribution over the z . We replace all distributions $\rho(x_{1:n})$ ($\rho = \mu, \nu, \xi$) everywhere by distributions $\rho(x_{1:n}|z_{1:n})$ conditioned on $z_{1:n}$. The $z_{1:n}$ conditions cause nowhere problems as they can essentially be thought of as fixed (or as oracles or spectators). So the bounds in Theorems 3.19...3.74 also hold in this case for all individual z 's.

Independent experiments and classification (CF). A typical experimental situation is a sequence of independent (i.i.d) experiments, predictions and observations. At time t one arranges an experiment z_t (or observes data z_t), then tries to make a prediction, and finally observes the true outcome x_t . Often one has a parameterized class of models (hypothesis space) $\mu_\theta(x_t|z_t)$ and wants to infer the true θ in order to make improved predictions. This is a special case of partial sequence prediction, where the hypothesis space $\mathcal{M} = \{\mu_\theta(x_{1:n}|z_{1:n}) = \mu_\theta(x_1|z_1) \cdot \dots \cdot \mu_\theta(x_n|z_n)\}$ consists of i.i.d. distributions, but note that ξ is not i.i.d. This is the same setting as for on-line learning of classification tasks, where a $z \in \mathcal{Z}$ should be classified as an $x \in \mathcal{X}$ (cf. Problem 3.12).

3.7.4 Prediction with Expert Advice

There are two schools of universal sequence prediction: We considered expected performance bounds for Bayesian prediction based on mixtures of environments, as is common in information theory and statistics [MF98]. The other approach are predictors based on expert advice (PEA) algorithms with worst case loss bounds in the spirit of Littlestone, Warmuth, Vovk and others. The two schools usually do not refer to each other much. We briefly describe PEA and compare both approaches. For a more comprehensive comparison see [MF98]. In the following we focus on topics not covered in [MF98]. PEA was invented in [LW89, LW94] and [Vov92] and further developed in [CB97, HKW98, KW99] and by many others. Many variations known by many names (prediction/learning with expert advice, weighted majority/average, aggregating strategy, boosting, hedge algorithm, ...) have meanwhile been invented. Early works in this direction are [Daw84, Ris89]. See [Vov99] for a review and further references. We describe the setting and basic idea of PEA for binary alphabet. Consider a finite binary sequence $x_1x_2\dots x_n \in \{0,1\}^n$ and a finite set \mathcal{E} of experts $e \in \mathcal{E}$ making predictions x_t^e in the unit interval $[0,1]$ based on past

observations $x_1 x_2 \dots x_{t-1}$. The loss of expert e in step t is defined as $|x_t - x_t^e|$. In the case of binary predictions $x_t^e \in \{0, 1\}$, $|x_t - x_t^e|$ coincides with our error measure (3.34). The PEA algorithm $p_{\beta n}$ combines the predictions of all experts. It forms its own prediction¹² $x_t^p \in [0, 1]$ according to some weighted average of the expert's predictions x_t^e . There are certain update rules for the weights depending on some parameter β . Various bounds for the total loss $L_p(\mathbf{x}) := \sum_{t=1}^n |x_t - x_t^p|$ of PEA in terms of the total loss $L_\varepsilon(\mathbf{x}) := \sum_{t=1}^n |x_t - x_t^\varepsilon|$ of the best expert $\varepsilon \in \mathcal{E}$ have been proven. It is possible to fine tune β and to eliminate the necessity of knowing n in advance. The first bound of this kind has been obtained in [CB97]:

$$L_p(\mathbf{x}) \leq L_\varepsilon(\mathbf{x}) + 2.8 \ln |\mathcal{E}| + 4\sqrt{L_\varepsilon(\mathbf{x}) \ln |\mathcal{E}|}. \quad (3.75)$$

The constants 2.8 and 4 have been improved in [AG00, YYY01]. The last bound in Theorem 3.36 with $S_n \leq D_n \leq \ln |\mathcal{M}|$ for uniform weights and with $E_n^{\Theta_\mu}$ increased to E_n^Θ reads

$$E_n^{\Theta_\xi} \leq E_n^\Theta + 2 \ln |\mathcal{M}| + 2\sqrt{E_n^\Theta \ln |\mathcal{M}|}.$$

It has a quite similar structure as (3.75), although the algorithms, the settings, the proofs, and the interpretation are quite different. Whereas PEA performs well in any environment, but only relative to a given set of experts \mathcal{E} , our Θ_ξ predictor competes with the best possible Θ_μ predictor (and hence with any other Θ predictor), but only in expectation and for a given set of environments \mathcal{M} . PEA depends on the set of expert, Θ_ξ depends on the set of environments \mathcal{M} . The basic $p_{\beta n}$ algorithm has been extended in different directions: incorporation of different initial weights ($|\mathcal{E}| \sim \ln \frac{1}{w_\nu}$) [LW89, Vov92], more general loss functions [HKW98], continuous valued outcomes [HKW98], and multi-dimensional predictions [KW99] (but not yet for the absolute loss). The work of [Yam98] lies somewhat in between PEA and this work; “PEA” techniques are used to prove expected loss bounds (but only for sequences of independent symbols/experiments and limited classes of loss functions). Finally, note that the predictions of PEA are continuous. This is appropriate for weather forecasters which announce the probability of rain, but the *decision* to wear sunglasses or to take an umbrella is binary, and the suffered loss depends on this binary decision, and not on the probability estimate. It is possible to convert the continuous prediction of PEA into a probabilistic binary prediction by predicting 1 with probability $x_t^p \in [0, 1]$. $|x_t - x_t^p|$ is then the probability of making an error. Note that the expectation is taken over the probabilistic prediction, whereas for the deterministic Θ_ξ algorithm the expectation is taken over the environmental distribution μ . The multi-dimensional case [KW99] could then be interpreted as a (probabilistic) prediction of symbols over an alphabet $\mathcal{X} = \{0, 1\}^d$, but error bounds for the absolute loss have yet to be proven. In [FS97] the regret is bounded by $\ln |\mathcal{E}| + \sqrt{2\tilde{L} \ln |\mathcal{E}|}$ for arbitrary unit loss function and alphabet, where \tilde{L} is an upper bound on L_ε ,

¹²The original PEA version [LW89] had discrete prediction $x_t^p \in \{0, 1\}$ with (necessarily) double as many errors as the best expert and is only of historical interest any more.

which has to be known in advance. It would be interesting to generalize PEA and bound (3.75) to arbitrary alphabet and weights and to general loss functions with probabilistic interpretation.

3.7.5 Outlook

In the following we discuss several directions in which the findings of this work may be extended.

Infinite Alphabet. In many cases the basic prediction unit is not a letter, but a number (for inducing number sequences), or a word (for completing sentences), or a real number or vector (for physical measurements). The prediction may either be generalized to a block by block prediction of symbols or, more suitably, the finite alphabet \mathcal{X} could be generalized to countable (numbers, words) or continuous (real or vector) alphabet. The presented Theorems are independent of the size of \mathcal{X} and hence should generalize to countably infinite alphabets by appropriately taking the limit $|\mathcal{X}| \rightarrow \infty$ and to continuous alphabets by a denseness or separability argument. Since the proofs are also independent of the size of \mathcal{X} we may directly replace all finite sums over \mathcal{X} by infinite sums or integrals and carefully check the validity of each operation. We expect all Theorems to remain valid in full generality, except for minor technical existence and convergence constraints.

An infinite prediction space \mathcal{Y} was no problem at all as long as we assumed the existence of $y_t^{\Lambda^p} \in \mathcal{Y}$ (3.45). In case $y_t^{\Lambda^p} \in \mathcal{Y}$ does not exist one may define $y_t^{\Lambda^p} \in \mathcal{Y}$ in a way to achieve a loss at most $\varepsilon_t = o(t^{-1})$ larger than the infimum loss. We expect a small finite correction of the order of $\varepsilon = \sum_{t=1}^{\infty} \varepsilon_t < \infty$ in the loss bounds somehow.

More Active Systems. Prediction means guessing the future, but not influencing it. A small step in the direction to more active systems was to allow the Λ system to act and to receive a loss $\ell_{x_t y_t}$ depending on the action y_t and the outcome x_t . The probability μ is still independent of the action, and the loss function ℓ^t has to be known in advance. This ensures that the greedy strategy (3.45) is optimal. The loss function may be generalized to depend not only on the history $x_{<t}$, but also on the historic actions $y_{<t}$ with μ still independent of the action. It would be interesting to know whether the scheme Λ and/or the loss bounds generalize to this case. The full model of an acting agent influencing the environment has been developed in [Hut01d], but loss bounds have yet to be proven.

Miscellaneous. Another direction is to investigate the learning aspect of universal prediction. Many prediction schemes explicitly learn and exploit a model of the environment. Learning and exploitation are melted together in the framework of universal Bayesian prediction. A separation of these two aspects in the spirit of hypothesis learning with MDL [VL00] could lead to new insights. Also, the separation of noise from useful data, usually an important issue [GTV01], did not play a role here. The attempt at an information theoretic interpretation of Theorem 3.63 may be made more rigorous in this or another way. In the end, this may lead to

a simpler proof of Theorem 3.63 and maybe even for the loss bounds. A unified picture of the loss bounds obtained here and the loss bounds for predictors based on expert advice (PEA) could also be fruitful. Yamanishi [Yam98] used PEA methods to prove expected loss bounds for Bayesian prediction, so maybe the proof technique presented here could be used *vice versa* to prove more general loss bounds for PEA. Maximum-likelihood predictors may also be studied. Since $2^{-K(x)}$ (or some of its variants) is a close approximation of ξ , it is generally believed that predictions based on K are as good as predictions based on ξ . Loss bounds for predictors based on K would prove this conjecture, but we have some heuristic arguments that there are situations where K may fail. Finally, the system should be applied to specific induction problems for specific \mathcal{M} with computable ξ .

3.8 Summary

We compared universal predictions based on Bayes-mixtures ξ to the infeasible informed predictor based on the unknown true generating distribution μ . We have shown that the universal posterior ξ converges to μ and that $\xi/\mu \rightarrow 1$. Our main focus was on a decision-theoretic setting, where each prediction $y_t \in \mathcal{X}$ (or more generally action $y_t \in \mathcal{Y}$) results in a loss $\ell_{x_t y_t}$ if x_t is the true next symbol of the sequence. We have shown that the Λ_ξ predictor suffers only slightly more loss than the Λ_μ predictor. We have shown that the derived error and loss bounds cannot be improved in general, i.e. without making extra assumptions on ℓ , μ , \mathcal{M} , or w_ν . Within a factor of 2 this is also true for any μ independent predictor. We have also shown Pareto-optimality of ξ in the sense that there is no other predictor which performs better or equal in all environments $\nu \in \mathcal{M}$ and strictly better in at least one. Optimal predictors can (in most cases) be based on a mixture distributions ξ . Finally we gave an Occam's razor argument that Solomonoff's prior with weights $w_\nu = 2^{-K(\nu)}$ is optimal, where $K(\nu)$ is the Kolmogorov complexity of ν . Of course, optimality always depends on the setup, the assumptions, and the chosen criteria. For instance, the universal predictor was not always Pareto-optimal, but at least for many popular, and for all decision theoretic performance measures. Bayes predictors are also not necessarily optimal under worst case criteria [CBL01]. We also derived a bound for the relative entropy between ξ and μ in the case of a continuously parameterized family of environments, which allowed us to generalize the loss bounds to continuous \mathcal{M} . Furthermore, we discussed the duality between the Bayes-mixture and expert-mixture approaches and results, classification tasks, games of chances, infinite alphabet, active systems influencing the environment, and others.

3.9 Technical Proofs

3.9.1 How to Deal with $\mu=0$

Some expressions (like conditional or inverse probabilities) are undefined for zero μ . We thought of the following solutions:

Avoid the problem. We may restrict ourselves to $\mu(x_{1:n}) > 0 \forall x_{1:n}$.

- + The treatment in this work is then rigorous and a zero μ can be approximated to an arbitrary precision. From a practical point of view this is a completely satisfactory approach.
- Theoretically unsatisfactory, because deterministic environments (for which $\mu(x_t|x_{<t}) = 0$ for all but one x_t) are of special interest, and not just esoteric limits.

Take the limit. Develop all theorems for $\mu^{(i)} > 0$ and finally perform the limit $\mu^{(i)} \xrightarrow{i \rightarrow \infty} \mu$, where μ might be zero for some strings. For instance $\mu^{(\varepsilon)}(x_{1:n}) := (1-\varepsilon)\mu(x_{1:n}) + \varepsilon/2^n$ and $\varepsilon \rightarrow 0$ will do.

- + Rigorous treatment with the advantage not having to deal with the problem until the end. If all spaces are finite, then interchange of finite sums or maxs with $\lim_{i \rightarrow \infty}$ is safe.
- Problematic for infinite spaces (e.g. alphabet, time, ...), since limits may not be interchangeable.

Face the problem. A way of facing the problem, which is slightly different from Section 3.2.1, is to restrict the set of strings to one with non-zero μ probability. Define the critical set $Z := \bigcup_{x \in \mathcal{X}^*: \mu(x)=0} \Gamma_x$, where $\Gamma_x := \{\omega : \omega_{1:l(x)} = x\}$ is defined as the cylinder set containing all infinite sequences ω starting with x .

- + Since Z is a countable (for countable alphabet) union of cylinder sets $\Gamma_{x_{1:k}}$ of measure zero, Z itself is measurable with μ -measure zero. So all theorems proven with μ probability 1 on $\Gamma_\varepsilon \setminus Z$ still hold on Γ_ε with μ probability 1, since $\mu(Z) = 0$.
- All sums over x have to be restricted appropriately. More seriously, other measures on Γ_ε , especially ξ deteriorate to semimeasures on $\Gamma_\varepsilon \setminus Z$ (see Section 2.4 and Problem 3.1).

Ignore the problem. Address other more severe or interesting problems, first. Why wasting time with bugging with exceptions when everything seems to work well anyway and there are more important problems to solve.

- + Time efficient “physicist” approach (like: always exchange limits and integrals until you get into troubles).
- The approach is risky and a mathematician will turn in his grave.

We often (but not always) faced the problem, but decided to avoid/ignore these subtleties in the presentation (see Problem 3.8).

3.9.2 Entropy Inequalities (3.11)

¹³We show that

$$\frac{1}{2} \sum_{i=1}^N f(y_i - z_i) \leq f\left(\sqrt{\frac{1}{2} \sum_{i=1}^N y_i \ln \frac{y_i}{z_i}}\right) \quad \text{for } y_i \geq 0, \quad z_i \geq 0, \quad \sum_{i=1}^N y_i = 1 = \sum_{i=1}^N z_i. \quad (3.76)$$

for any convex and even ($f(x) = f(-x)$) function with $f(0) \leq 0$. For $f(x) = x^2$ we get inequality Lemma 3.11s, for $f(x) = |x|$ we get inequality (3.20). To prove (3.76) we partition $i \in \{1, \dots, N\} = G^+ \cup G^-$, $G^+ \cap G^- = \{\}$, and define $y^\pm := \sum_{i \in G^\pm} y_i$ and

$z^\pm := \sum_{i \in G^\pm} z_i$. It is well known that the relative Entropy is positive, i.e.

$$\sum_{i \in G^\pm} p_i \ln \frac{p_i}{q_i} \geq 0 \quad \text{for } p_i \geq 0, \quad q_i \geq 0, \quad \sum_{i \in G^\pm} p_i = 1 = \sum_{i \in G^\pm} q_i. \quad (3.77)$$

Note that there are 4 probability distributions (p_i and q_i for $i \in G^+$ and $i \in G^-$). For $i \in G^\pm$, $p_i := y_i / y^\pm$ and $q_i := z_i / z^\pm$ satisfy the conditions on p and q . Inserting this into (3.77) and rearranging the terms we get

$$\sum_{i \in G^\pm} y_i \ln \frac{y_i}{z_i} \geq y^\pm \ln \frac{y^\pm}{z^\pm}.$$

If we sum over \pm and define $y \equiv y^+ = 1 - y^-$ and $z \equiv z^+ = 1 - z^-$ we get

$$\sum_{i=1}^N y_i \ln \frac{y_i}{z_i} \geq \sum_{\pm} y^\pm \ln \frac{y^\pm}{z^\pm} = y \ln \frac{y}{z} + (1-y) \ln \frac{1-y}{1-z} \geq 2(y-z)^2 \quad (3.78)$$

The last inequality is elementary and well known. For the special choice $G^\pm := \{i : y_i \gtrless z_i\}$, we can upper bound $\sum_i f(y_i - z_i)$ as follows

$$\begin{aligned} \sum_{i \in G^\pm} f(y_i - z_i) &\stackrel{(a)}{=} \sum_{i \in G^\pm} f(|y_i - z_i|) \stackrel{(b)}{\leq} f\left(\sum_{i \in G^\pm} |y_i - z_i|\right) \stackrel{(c)}{=} f\left(\left|\sum_{i \in G^\pm} y_i - z_i\right|\right) \stackrel{(d)}{=} \\ &\stackrel{(d)}{=} f(|y^\pm - z^\pm|) \stackrel{(e)}{=} f(|y - z|) \stackrel{(f)}{=} f(\sqrt{(y - z)^2}) \stackrel{(g)}{\leq} f\left(\sqrt{\frac{1}{2} \sum_{i=1}^N y_i \ln \frac{y_i}{z_i}}\right) \end{aligned} \quad (3.79)$$

(a) follows from the symmetry of f . (b) follows from the convexity¹⁴ of f and from $f(0) \leq 0$. (c) is true, since all $y_i - z_i$ are positive/negative for $i \in G^\pm$ due to the special

¹³We will not explicate every subtlety and only sketch the proofs. Subtleties regarding $y, z = 0/1$ have been checked but will be passed over. $0 \ln \frac{0}{z_i} := 0$ even for $z_i = 0$. Positive means ≥ 0 . The probability constraints in (3.76) on y and z apply to all appendices. $z > 0$ if $y > 0$.

¹⁴Inserting $y = 0$ and $x = a + b$ in the convexity definition $\alpha f(x) + (1 - \alpha)f(y) \geq f(\alpha x + (1 - \alpha)y)$ leads to $\alpha f(a + b) + (1 - \alpha)f(0) \geq f(\alpha(a + b))$. Inserting $\alpha = \frac{a}{a+b}$ and $\alpha = \frac{b}{a+b}$ and adding both inequalities gives $f(a + b) + f(0) \geq f(a) + f(b)$ for $a, b \geq 0$. Using $f(0) \leq 0$ we get $f(\sum_i x_i) \geq \sum_i f(x_i)$ for $x_i \geq 0$ by induction.

choice of G^\pm . (d) and (e) follow from the definition of $y^{(\pm)}$ and $z^{(\pm)}$, (f) is obvious. (g) follows from (3.78) and the monotonicity¹⁵ of $\sqrt{\cdot}$ and f for positive arguments. Inequality (3.76) follows by summation of (3.79) over \pm and noting that $f(\sqrt{\cdot})$ is independent of \pm .

This proves Lemma 3.11f. Inserting $f(x) = x^2$ yields Lemma 3.11s, inserting $f(x) = |x|$ yields Lemma 3.11a. Lemma 3.11b follows from

$$\sum_{i=1}^N y_i \left| \ln \frac{y_i}{z_i} \right| - \sum_i y_i \ln \frac{y_i}{z_i} = -2 \sum_{i \in G^-} y_i \ln \frac{y_i}{z_i} \leq 2 \sum_{i \in G^-} z_i - y_i = \sum_{i=1}^N |y_i - z_i|,$$

where we have used $-\ln x \leq \frac{1}{x} - 1$. Lemma 3.11h is proven differently. For arbitrary $y \geq 0$ and $z \geq 0$ we define

$$f(y, z) := y \ln \frac{y}{z} - (\sqrt{y} - \sqrt{z})^2 + z - y = 2yg(\sqrt{z/y}) \quad \text{with} \quad g(t) := -\ln t + t - 1 \geq 0.$$

This shows $f \geq 0$, and hence $\sum_i f(y_i, z_i) \geq 0$, which implies

$$\sum_i y_i \ln \frac{y_i}{z_i} - \sum_i (\sqrt{y_i} - \sqrt{z_i})^2 \geq \sum_i y_i - \sum_i z_i = 1 - 1 = 0.$$

This proves Lemma 3.11h. □

3.9.3 Error Inequality (3.36)

Here we give a direct proof of the second bound in Theorem 3.36. Again, we try to find constants A and B that satisfy the linear inequality

$$E_n^{\Theta_\xi} \leq (A+1)E_n^{\Theta_\mu} + (B+1)S_n. \quad (3.80)$$

If we could show

$$e_t^{\Theta_\xi}(x_{<t}) \leq (A+1)e_t^{\Theta_\mu}(x_{<t}) + (B+1)s_t(x_{<t}) \quad (3.81)$$

for all $t \leq n$ and all $x_{<t}$, (3.80) would follow immediately by summation and the definition of E_n and S_n . With the abbreviations (3.12) and the abbreviations $m = x_t^{\Theta_\mu}$ and $s = x_t^{\Theta_\xi}$ the various error functions can then be expressed by $e_t^{\Theta_\xi} = 1 - y_s$, $e_t^{\Theta_\mu} = 1 - y_m$ and $s_t = \sum_i (y_i - z_i)^2$. Inserting this into (3.81) we get

$$1 - y_s \leq (A+1)(1 - y_m) + (B+1) \sum_{i=1}^N (y_i - z_i)^2. \quad (3.82)$$

¹⁵Inserting $b = y = -x$ and $\alpha = \frac{1}{2}$ into the convexity definition and using the symmetry of f we get $f(b) \geq f(0)$. Inserting this into $f(a+b) + f(0) \geq f(a) + f(b)$ we get $f(a+b) \geq f(a)$ which proves that f is monotonically increasing for positive arguments ($a, b \geq 0$).

By definition of $x_t^{\Theta_\mu}$ and $x_t^{\Theta_\xi}$ we have $y_m \geq y_i$ and $z_s \geq z_i$ for all i . We prove a sequence of inequalities which show that

$$(B+1) \sum_{i=1}^N (y_i - z_i)^2 + (A+1)(1-y_m) - (1-y_s) \geq \dots \quad (3.83)$$

is positive for suitable $A \geq 0$ and $B \geq 0$, which proves (3.82). For $m = s$ (3.83) is obviously positive. So we will assume $m \neq s$ in the following. From the square we keep only contributions from $i = m$ and $i = s$.

$$\dots \geq (B+1)[(y_m - z_m)^2 + (y_s - z_s)^2] + (A+1)(1-y_m) - (1-y_s) \geq \dots$$

By definition of y , z , \mathcal{M} and s we have the constraints $y_m + y_s \leq 1$, $z_m + z_s \leq 1$, $y_m \geq y_s \geq 0$ and $z_s \geq z_m \geq 0$. From the latter two it is easy to see that the square terms (as a function of z_m and z_s) are minimized by $z_m = z_s = \frac{1}{2}(y_m + y_s)$. Furthermore, we define $x := y_m - y_s$ and eliminate y_s .

$$\dots \geq (B+1)\frac{1}{2}x^2 + A(1-y_m) - x \geq \dots \quad (3.84)$$

The constraint on $y_m + y_s \leq 1$ translates into $y_m \leq \frac{x+1}{2}$, hence (3.84) is minimized by $y_m = \frac{x+1}{2}$.

$$\dots \geq \frac{1}{2}[(B+1)x^2 - (A+2)x + A] \geq \dots \quad (3.85)$$

(3.85) is quadratic in x and minimized by $x^* = \frac{A+2}{2(B+1)}$. Inserting x^* gives

$$\dots \geq \frac{4AB - A^2 - 4}{8(B+1)} \geq 0 \quad \text{for} \quad B \geq \frac{1}{4}A + \frac{1}{A}, \quad A > 0, \quad (\Rightarrow B \geq 1). \quad (3.86)$$

Inequality (3.80) therefore holds for any $A > 0$, provided we insert $B = \frac{1}{4}A + \frac{1}{A}$. Thus we might minimize the r.h.s. of (3.80) w.r.t. A leading to the upper bound

$$E_n^{\Theta_\xi} \leq E_n^{\Theta_\mu} + S_n + \sqrt{4E_n^{\Theta_\mu} S_n + S_n^2} \quad \text{for} \quad A^2 = \frac{S_n}{E_n^{\Theta_\mu} + \frac{1}{4}S_n}$$

which completes the proof of Theorem 3.36. \square

3.9.4 Binary Loss Inequality for $z \leq \frac{1}{2}$ (3.57)

With the definition

$$f(y, z) := B' \cdot \left[y \ln \frac{y}{z} + (1-y) \ln \frac{1-y}{1-z} \right] + A' \cdot (1-y) \frac{z}{1-z} - y, \quad z \leq \frac{1}{2} \quad (3.87)$$

we show $f(y, z) \geq 0$ for suitable $A' \equiv A+1$ and $B' \equiv B+1$. We do this by showing that $f \geq 0$ at all extremal values and “at” boundaries. $f \rightarrow +\infty$ for $z \rightarrow 0$, if we choose

$B' > 0$. For the boundary $z = \frac{1}{2}$ we lower bound the relative entropy by the sum over squares (Lemma 3.11s)

$$f(y, \tfrac{1}{2}) \geq 2B'(y - \tfrac{1}{2})^2 + A'(1 - y) - y$$

The r.h.s. is quadratic in y with minimum at $y^* = \frac{A' + 2B' + 1}{4B'}$, which implies

$$f(y, \tfrac{1}{2}) \geq f(y^*, \tfrac{1}{2}) \geq \frac{4AB - A^2 - 4}{8(B + 1)} \geq 0 \quad \text{for } B \geq \tfrac{1}{4}A + \tfrac{1}{A}, \quad A > 0, \quad (\Rightarrow B \geq 1).$$

Furthermore, for $A \geq 4$ and $B \geq 1$ we have $f(y, \frac{1}{2}) \geq 2(1 - y)(3 - 2y) \geq 0$. Hence $f(y, \frac{1}{2}) \geq 0$ for $B \geq \frac{1}{A} + 1$, since for $A \geq 4$ it implies $B \geq 1$ and for $A \leq 4$ it implies $B \geq \frac{1}{4}A + \frac{1}{A}$.

The extremal condition $\partial f / \partial z = 0$ (keeping y fixed) leads to

$$y = y^* := z \cdot \frac{B'(1 - z) + A'}{B'(1 - z) + A'z}.$$

Inserting y^* into the definition of f and, again, replacing the relative entropy by the sum over squares (Lemma 3.11s), we get

$$f(y^*, z) \geq 2B'(y^* - z)^2 + A'(1 - y^*) \frac{z}{1 - z} - y^* = \frac{z(1 - z)}{[B'(1 - z) + A'z]^2} \cdot g(z),$$

$$g(z) := 2B'A'^2 z(1 - z) + [(A' - 1)B'(1 - z) - A'](B' + A' \frac{z}{1 - z}).$$

We have reduced the problem to showing $g \geq 0$. If the bracket [...] is positive, then g is positive. If the bracket is negative, we can decrease g by increasing $\frac{z}{1 - z} \leq 1$ in $(B' + A' \frac{z}{1 - z})$ to 1. The resulting expression is now quadratic in z with minima at the boundary values $z = 0$ and $z = \frac{1}{2}$. It is therefore sufficient to check

$$g(0) \geq (AB - 1)(A + B + 2) \geq 0 \quad \text{and} \quad g(\tfrac{1}{2}) \geq \tfrac{1}{2}(AB - 1)(2A + B + 3) \geq 0$$

which is true for $B \geq \frac{1}{A}$. In summary we have proved (3.87) for $B \geq \frac{1}{A} + 1$ and $A > 0$. \square

3.9.5 Binary Loss Inequality for $z \geq \frac{1}{2}$ (3.58)

With the definition

$$f(y, z) := B' \cdot \left[y \ln \frac{y}{z} + (1 - y) \ln \frac{1 - y}{1 - z} \right] + A' \cdot (1 - y) - y \frac{1 - z}{z}, \quad z \geq \tfrac{1}{2} \quad (3.88)$$

we show $f(y, z) \geq 0$ for suitable $A' \equiv A + 1 > 1$ and $B' \equiv B + 1 > 2$ similarly as in the last paragraph by proving that $f \geq 0$ at all extremal values and “at” boundaries. $f \rightarrow +\infty$ for $z \rightarrow 1$. The boundary $z = \frac{1}{2}$ has already been checked in in the last paragraph. The extremal condition $\partial f / \partial z = 0$ (keeping y fixed) leads to

$$y = y^* := z \cdot \frac{B'z}{(B' + 1)z - 1}.$$

Inserting y^* into the definition of f and replacing the relative entropy by the sum over squares (Lemma 3.11s), we get

$$f(y^*, z) \geq 2B'(y^* - z)^2 + A'(1 - y^*) - y^* \frac{1-z}{z} = \frac{z(1-z)}{[(B'+1)z-1]^2} \cdot g(z),$$

$$g(z) := [(A' - 1)B'z - A' + 2z(1 - z)](B' + 1 - \frac{1}{z}) + 2(1 - z)^2.$$

We have reduced the problem to showing $g \geq 0$. Since $(B' + 1 - \frac{1}{z}) \geq 0$ it is sufficient to show that the bracket is positive. We solve $[\dots] \geq 0$ w.r.t. B and get

$$B \geq \frac{1 - 2z(1 - z)}{z} \cdot \frac{1}{A} + \frac{1 - z}{z}.$$

For $B \geq \frac{1}{A} + 1$ this is satisfied for all $\frac{1}{2} \leq z \leq 1$. In summary we have proved (3.88) for $B \geq \frac{1}{A} + 1$ and $A > 0$. \square

3.9.6 General Loss Inequality (3.53)

We reduce

$$f(\mathbf{y}, \mathbf{z}) := B' \sum_{i=1}^N y_i \ln \frac{y_i}{z_i} + A' \sum_{i=1}^N y_i \ell_{im} - \sum_{i=1}^N y_i \ell_{is} \geq 0 \quad (3.89)$$

$$\text{for } \sum_{i=1}^N z_i d_i \geq 0, \quad d_i := \ell_{im} - \ell_{is} \quad (3.90)$$

to the binary $N=2$ case. We do this by keeping \mathbf{y} fixed and showing that f as a function of \mathbf{z} is positive at all extrema in the interior of the simplex $\Delta := \{\mathbf{z} : \sum_i z_i = 1, z_i \geq 0\}$ of the domain of \mathbf{z} and “at” all boundaries. First, the boundaries $z_i \rightarrow 0$ are safe as $f \rightarrow \infty$ for $B' > 0$. Variation of f w.r.t. to \mathbf{z} leads to a minimum at $\mathbf{z} = \mathbf{y}$. If $\sum_i z_i d_i \geq 0$, we have

$$f(\mathbf{y}, \mathbf{y}) = \sum_i y_i (A' \ell_{im} - \ell_{is}) \geq \sum_i y_i (\ell_{im} - \ell_{is}) = \sum_i z_i d_i \geq 0.$$

In the first inequality we used $A' > 1$. If $\sum_i z_i d_i < 0$, $\mathbf{z} = \mathbf{y}$ is outside the valid domain due to the constraint (3.90) and the valid minima are attained at the boundary $\Delta \cap P$ with $P := \{\mathbf{z} : \sum_i z_i d_i = 0\}$. We implement the constraints with the help of Lagrange multipliers and extremize

$$L(\mathbf{y}, \mathbf{z}) := f(\mathbf{y}, \mathbf{z}) + B' \lambda \sum z_i + B' \mu \sum z_i d_i.$$

$\partial L / \partial z_i = 0$ leads to $y_i = y_i^* := z_i (\lambda + \mu d_i)$. Summing this equation over i we obtain $\lambda = 1$. μ is a function of \mathbf{y} for which a formal expression might be given. If we eliminate y_i in favor of z_i , we get

$$f(\mathbf{y}^*, \mathbf{z}) = \sum_i c_i z_i, \quad c_i := (1 + \mu d_i)(B' \ln(1 + \mu d_i) + A' \ell_{im} - \ell_{is}).$$

In principle μ is a function of \mathbf{y} but we can treat μ directly as an independent variable, since \mathbf{y} has been eliminated.

The next step is to determine the extrema of the function $f = \sum c_i z_i$ for $\mathbf{z} \in \Delta \cap P$. For clearness we state the line of reasoning for $N=3$. In this case Δ is a triangle. As f is linear in \mathbf{z} it assumes its extrema at the vertices of the triangle, where all $z_i = 0$ except one. But we have to take into account a further constraint $\mathbf{z} \in P$. The plane P intersects triangle Δ in a finite line (for $\Delta \cap P = \{\}$ the only boundaries are $z_i \rightarrow 0$ which have already been treated). Again, as f is linear, it assumes its extrema at the ends of the line, i.e. at edges of the triangle Δ on which all but two z_i are zero. With a similar line of arguments for $N > 3$ we conclude that a necessary condition for a minimum of f at the boundary is that at most two z_i are non-zero. But this implies that all but two y_i are zero. If we had eliminated \mathbf{z} in favor of \mathbf{y} , we could not have made the analogous conclusion because $y_i = 0$ does not necessarily imply $z_i = 0$. We have effectively reduced the problem of showing $f(\mathbf{y}^*, \mathbf{z}) \geq 0$ to the case $N=2$. We can go back one step further and prove (3.89) for $N=2$, which implies $f(\mathbf{y}^*, \mathbf{z}) \geq 0$ for $N=2$. A proof of (3.89) for $N=2$ implies, by the arguments given above, that it holds for all N . This is what we set out to show here. \square

The $N=2$ case has been proven in Section 3.4 and the two previous paragraphs.

3.10 History & References

There are good introductions and surveys of Solomonoff sequence prediction [LV92a, LV97], inductive inference in general [AS83, Sol97, MF98], competitive online statistics [Vov99], and reasoning under uncertainty [Grü98], the latter also containing a more serious discussion of the case $\mu \notin \mathcal{M}$ in a related context. The convergence $\xi \rightarrow \mu$ of Theorem 3.19(iii) has first been proven in [BD62] in a more general framework with the help of martingales, but the martingale proof does not provide the speed of convergence. Solomonoff's contribution was to focus on the set of all (semi)computable distributions [Sol64] and to prove Theorem 3.19(i) for binary alphabet, which shows that convergence (iii) of ξ to μ is rapid [Sol78]. The generalization of (i) to arbitrary finite alphabet has probably first been shown by the author in [Hut01a], but may have occurred earlier somewhere in the statistics literature. Convergence (v) of the ratio ξ/μ to 1 w. μ .p.1 has first been shown by Gács with the help of martingales [LV97], again not allowing to estimate the speed of convergence. The elementary proof of $\xi/\mu \rightarrow 1$ i.m.s., i.e. of (iv) and (v) is new, showing that convergence is rapid. (vi) directly follows from Lemma 3.11(a). Lemma 3.11(a) is due Pinsker [Pin64] and Csiszàr [Csi67], and can be found in [CT91, Lem.12.6.1]. A different proof of Lemma 3.11(h) can be found in [BM98, p178]. Lemma 3.11(b) is also known [Bar00]. All other results are new except when mentioned otherwise.

3.11 Problems

3.1 (Semimeasures) [C30u/C40o] All results in this chapter have been obtained for probability measures μ , and ξ and w_ν , i.e. $\sum_{x_{1:t}} \xi(x_{1:t}) = \sum_{x_{1:t}} \mu(x_{1:t}) = \sum_\nu w_\nu = 1$.

On the other hand, the primary class \mathcal{M} of interest in this work is the class of all enumerable semimeasures and $\sum_\nu w_\nu \leq 1$, (see Section 2.4). In general, each of the following 4 items could be semi ($<$) or not ($=$): $(\xi, \mu, \mathcal{M}, w_\nu)$, where \mathcal{M} is semi if some elements are semi.

Which of the 2^4 combinations makes sense? (Hint: 6 of the 16). Show that the entropy inequalities (3.11) holds for $(<, =, <, <)$, but not for $(<, <, <, <)$. Nevertheless, show that $\xi \rightarrow \mu$ (Th.3.19 iii) for $(<, <, <, <)$ with maximal μ semi-probability, i.e. fails with μ semi-probability 0. Generalize all other theorems in this chapter as far as possible to the semi case.

3.2 (Dominance of Speed Prior) [C40oi] Which (semi)measures are multiplicatively dominated by the Speed prior defined in [Sch02c]. Show that the Speed prior dominates every computable deterministic environment. This implies that also every probabilistic environment, assigning positive probability only to those strings each computable within time t with the same t for all, is dominated by the Speed prior. Is there an easily characterizable class (or at least a less trivial subclass than the one above) of dominated probabilistic environments?

3.3 (Comparing two mixtures) [C05u/C40o] Consider two mixtures ξ and ξ' over \mathcal{M} . Show that $\sum_{t=1}^{\infty} \mathbf{E}[\sum_{x_t} (\xi(x_t|x_{<t}) - \xi'(x_t|x_{<t}))^2] \leq \ln w_\mu^{-1} + \ln w'_\mu^{-1} < \infty$, i.e. $\xi(x_t|x_{<t}) \rightarrow \xi'(x_t|x_{<t})$ for $t \rightarrow \infty$ with μ probability 1 for all $\mu \in \mathcal{M}$. (compare with Theorem 3.19(i)). Furthermore, show that $L_n^{\Lambda_\xi} - L_n^{\Lambda_{\xi'}} \leq O(\sqrt{L_n^{\Lambda_\mu}})$ (compare with Theorem 3.48). Is the stronger result $L_\infty^{\Lambda_\xi} - L_\infty^{\Lambda_{\xi'}} < \infty$ also true?

3.4 (Convergence and loss bounds with high probability) [C30oi] Show that $\mathbf{P}[\sum_{t=1}^n (\mu(x_t|x_{<t}) - \xi(x_t|x_{<t}))^2 \geq \frac{1}{\varepsilon} \ln w_\mu^{-1}] \leq \varepsilon$ and $\mathbf{P}[\sum_{t=1}^n (l_t^{\Lambda_\xi} - l_t^{\Lambda_\mu})^2 \geq \frac{2}{\varepsilon} \ln w_\mu^{-1}] \leq \varepsilon$, where \mathbf{P} denotes μ -probability. Use Theorem 3.19(i) and Theorem 3.59(i) and Markov's inequality. Is it possible to prove similar high probability bounds for the ratio $l_t^{\Lambda_\xi}/l_t^{\Lambda_\mu}$, possibly exploiting Theorem 3.48 and Corollary 3.49(iii). High probability bounds on $l_t^{\Lambda_\xi}(x_{<t})$ still involve an expectation over x_t (see definition of $l_t^{\Lambda_\xi}$). Is it possible to prove high probability bounds on the difference or ratio of $\ell_{x_t y_t}^{\Lambda_\xi}$ and $\ell_{x_t y_t}^{\Lambda_\mu}$, which do not involve any expectations?

3.5 (Pareto-optimality) [C30u] Show that ξ is not Pareto-optimal w.r.t. the α -norm $\mathcal{F}(\nu, \xi) = \|\nu - \xi\|_\alpha = \sqrt[\alpha]{\sum_{x_t} |\nu(x_t|x_{<t}) - \xi(x_t|x_{<t})|^\alpha}$ if $\alpha \neq 2$. Further, Pareto-optimality of ξ w.r.t. \mathcal{F}_1 and \mathcal{F}_2 is neither a necessary nor a sufficient condition for ξ being Pareto-optimal w.r.t. their sum $\mathcal{F}_1 + \mathcal{F}_2$. Finally, if ξ is Pareto-optimality w.r.t. \mathcal{F} , then ξ is also Pareto-optimality w.r.t. any monotone increasing function of \mathcal{F} (e.g. \mathcal{F}^α , $\alpha \geq 0$).

Hint: Intuition on this problem can be gained by considering probability vectors $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \Delta \subset \mathbb{R}^3$, where Δ is the 2d probability triangle, and $\mathbf{z} = w\mathbf{x} + (1-w)\mathbf{y}$ is a mixture of \mathbf{x} and \mathbf{y} . Consider the sets $M_{\mathbf{x}} := \{\mathbf{r} : \mathcal{F}(\mathbf{x}, \mathbf{r}) \leq \mathcal{F}(\mathbf{x}, \mathbf{z})\}$ and analogously $M_{\mathbf{y}}$. $M_{\mathbf{x}} \cap M_{\mathbf{y}}$ is not empty; it contains \mathbf{z} . If $M_{\mathbf{x}} \cap M_{\mathbf{y}}$ has an interior, then \mathbf{z} is not

Pareto-optimal. Visualize the 1d boundaries of the 2d areas M_x and M_y qualitatively for the various performance measures \mathcal{F} . Now consider mixtures of three vectors in \mathbb{R}^4 . This should give you enough intuition to prove Pareto-optimality and to construct counter-examples.

3.6 (Lower error bound) [C30u] It is possible to derive good (but not tight) lower error bounds with a fixed (n independent) set \mathcal{M} and weights, as opposed to the n dependent set \mathcal{M} chosen in the proof of Theorem 3.64. For instance, choose $\mathcal{M} = \{\mu_1, \mu_2\}$ with $\mu_{1/2}(1|x_{<t}) = \frac{1}{2} \pm \varepsilon_t$ with $\varepsilon_t = \min\{\frac{1}{2}, \sqrt{\ln w_{\mu_1}^{-1}} / \sqrt{t \ln t}\}$. For $w_{\mu_1} \ll w_{\mu_2}$ and $n \rightarrow \infty$ one can show that $E_n^{\Theta_\varepsilon} - E_n^{\Theta_\mu} \sim \frac{1}{\ln n} \sqrt{E_n^{\Theta_\mu} D_n}$. Is it possible to derive a tight(er) lower bound for different, but n independent \mathcal{M} ?

3.7 ((Non)uniqueness of universal weights) [C25ui] Section 3.6.4 showed that the universal weights $w_\nu = 2^{-K(\nu)}$ are optimal in a sense precisely stated in Theorem 3.70. Show that this choice for w_ν is not unique (even not within a constant factor). For instance, for $v_\nu = O(1)$ for $\nu = \xi_w$ and v_ν arbitrary (e.g. 0) for all other ν , the obvious dominance $\xi_\nu \geq v_\nu \nu$ can be improved to $\xi_\nu \geq w_\nu \nu$. Indeed, formally every choice of weights $v_\nu > 0 \forall \nu$ leads within a multiplicative constant to the same universal distribution, but this constant is not necessarily of “acceptable” size. Suitably define “acceptable size” by considering the implications for the loss bounds. Construct (counter)examples and necessary/sufficient conditions for weights to be acceptable.

3.8 (Deal with zero μ) [C35uo] Verify that all results in this chapter remain valid even if μ is allowed to take value zero for some arguments. Take the limit or face the problem as discussed in Section 3.9.1. Note, that when facing the problem, ξ deteriorates to a semimeasure (see Section 2.4 and Problem 3.1).

3.9 (Relations between random convergence criteria) [C30sm] Prove the relations in Lemma 3.9 between the various convergence criteria given in Definition 3.8. Show that no other implications hold by constructing example random sequences. More precisely, implications are strict with reverse being wrong and disconnected criteria (in the transitive hull) are incomparable. Show that convergence i.m.s implies that z_t deviates from z_* by more than ε only finitely many times and give a bound on the number.

3.10 (Individual $\xi \rightarrow \mu$ convergence) [C45om] In Problem 2.3 the open question whether $\xi_U(x_t|x_{<t})$ converges to $\mu(x_t|x_{<t})$ (in ratio or difference sense) individually for all Martin-Löf random sequences has been posed (short $\xi_U \xrightarrow{\text{M.L.}} \mu$). Theorem 3.22 shows that $\xi_U \xrightarrow{\text{M.L.}} \mu$ cannot be decided from ξ_U being a mixture distribution (3.5) or from the dominance property (3.6) alone. $\xi \notin \mathcal{M}$ for the classes used in Theorem 3.22. Construct $\xi_{\mathcal{M}} \in \mathcal{M}$ and prove a Theorem analogous to Theorem 3.22 for these \mathcal{M} 's (Start with two-element classes \mathcal{M} , then enlarge \mathcal{M} as far as possible). Hence, $\xi_U \in \mathcal{M}_U$ is also not sufficient to resolve $\xi_U \xrightarrow{\text{M.L.}} \mu$. Convert Theorem 3.19(ii/iv)

to potential μ .M.L. randomness tests, but show that they are not effective. Try also to generalize Vovk's result [Vov87] to non-recursive distributions. Where is the problem? With these insights try again to solve Problem 2.3.

3.11 (Speed of $\xi \rightarrow \mu$ convergence) [C35o] Theorem 3.19(i) shows that $\sum_{t=1}^{\infty} s_t < \infty$. If s_t were monotone decreasing ($s_t \geq s_{t+1}$) this would imply that s_t tends to zero faster than $1/t$, i.e. $s_t = o(1/t)$. Show (or refute) that this monotonicity is generally wrong. Prove or disprove the following conjecture: For $\xi = \xi_U$ and for every computable function $f > 0$ with $\lim_{t \rightarrow \infty} f_t = 0$ (however slow it converges to zero) there is a μ such that $\limsup_{t \rightarrow \infty} s_t / f_t > 0$. Validity of this result would imply that s_t can tend to zero arbitrarily slow, although it converges fast to zero in an average sense, e.g. $s_t > \varepsilon$ at most $\frac{1}{\varepsilon} \ln w_{\mu}^{-1}$ times.

3.12 (Learnability of the Universal Turing Machine) [C05u] Consider the problem of learning a function $f: \mathcal{Z} \rightarrow \mathcal{X}$. A sequence of sample pairs $(z_1, x_1), (z_2, x_2), \dots, (z_{n-1}, x_{n-1})$ with $x_i = f(z_i)$ is given. The task is to predict $x_n = f(z_n)$ from z_n . This setup is a special case of the one described in Section 3.7.3. Show that if f is a recursive function, then the Θ_{ξ} predictor makes at most a finite number of prediction errors, more precisely, at most $2 \ln 2 \cdot K(f) + O(1)$ errors. Consider now the universal (partial) function $f(z) := U(z)$, where U is some universal Turing machine. Show that U is learnable in the sense that after finite time n , Θ_{ξ} correctly predicts $x_i = U(z_i)$ for all $i > n$ as long as all z_i are in the domain of f . What happens if for some z_i , U does not halt?

3.13 (Posterization) [C30ui] Show that many properties of Kolmogorov complexity, Solomonoff's prior, and (policies based on) Bayes-mixtures remain valid after "posterization". With posterization we mean replacing $L_n, w_{\nu}, K(\nu), \nu(x_{1:n})$, etc. by the posteriors $L_{kn} := \sum_{t=k}^n \mathbf{E}[\ell_{x_t y_t} | x_{<k}], w_{\nu}(x_{<k}), K(\nu | x_{<k}), \nu(x_{k:n} | x_{<k})$, etc. Show that, strangely enough, for $\mathcal{M} = \mathcal{M}_U$ and $w_{\nu} = 2^{-K(\nu)}$ it is not true that $w_{\nu}(x_{<k}) \stackrel{\times}{\approx} 2^{-K(\nu | x_{<k})}$ (not even $\log_2 w_{\nu}(x_{<k}) = -K(\nu | x_{<k}) + O(\log)$ holds). The important \geq direction fails. Is the other direction \leq true? So, bounds of, e.g. $L_{kn}^{\Lambda_{\xi}}$ in terms of $\ln w_{\mu}(x_{<k})$ cannot be converted to bounds in terms of $K(\mu | x_{<k})$ unlike the $k=1$ case. But if we go one step back we see that a bound on $\ln[\mu(x_{k:n} | x_{<k}) / \xi(x_{k:n} | x_{<k})]$ is sufficient to bound $L_{kn}^{\Lambda_{\xi}}$. Use Problem 2.6(iii) to bound this expression by $\ln 2 \cdot \tilde{K}(\mu | x_{<k}) + O(?)$. The more information the history $x_{<k}$ contains about the environment μ , the smaller the bounds get.

3.14 (Probabilistic error bounds) [C35s/C35o] Instead of making (deterministic) predictions which minimize the ρ -expected loss or error, we may define probabilistic ρ -predictors which predict x_t from $x_{<t}$ with probability $\rho(x_t | x_{<t})$, and compare the performance of the ξ -predictor with the μ - or other ρ -predictors. For simplicity, consider binary sequences (drawn from μ) and the error-loss. Let $e_t^{\rho}(x_{<t}) := \mathbf{E}_t[1 - \rho(x_t | x_{<t})]$ be the error probability in the t^{th} prediction and $E_n^{\rho} := \sum_{t=1}^n \mathbf{E}[e_t^{\rho}(x_{<t})]$ be the μ -expected total number of errors in the first n predictions.

Prove the following error relations between universal ($\rho=\xi$), informed ($\rho=\mu$), and general (ρ) predictors:

$$\begin{aligned}
i) \quad & |E_n^\xi - E_n^\mu| \leq \frac{1}{2}A_n \leq D_n + \sqrt{2E_n^\mu D_n} \\
ii) \quad & E_n^\xi \geq \frac{1}{2}[S_n + E_n^\mu] \\
iii) \quad & E_n^\xi \geq E_n^\mu + D_n - \sqrt{2E_n^\mu D_n} \geq D_n \quad \text{for } E_n^\mu \geq 2D_n \\
iv) \quad & E_n^\mu \leq 2E_n^\rho, \quad e_n^\mu \leq 2e_n^\rho \quad \text{for any } \rho \\
v) \quad & E_n^\xi \leq 2E_n^\rho + D_n + \sqrt{4E_n^\rho D_n} \quad \text{for any } \rho,
\end{aligned}$$

where $A_n, S_n, D_n \leq \ln w_\mu^{-1}$ are defined in (3.13)–(3.16), and w_μ is the weight (3.5) of μ in ξ ((i)–(v) has been proven in [Hut01c]). This shows that the ξ -predictor performs well as compared to the μ -predictor ($E_n^\xi - E_n^\mu = O(\sqrt{E_n^\mu})$), but does not exclude the possibility that it makes twice as many errors as other (better) predictors (there is a factor two in $E_n^\xi/E_n^\rho \leq 2 + O((E_n^\rho)^{-1/2})$). Give an example for $E_n^\xi/E_n^{\Theta_\mu} \geq 2$ showing that the factor 2 in (iv) and (v) cannot be improved in general. Generalize (i)–(v) to nonbinary alphabet and arbitrary loss functions.

3.15 (Posterior convergence for unbounded horizon) [C15ui/C30o] Show that for unbounded horizon $h_t \rightarrow \infty$, there exist \mathcal{M} , μ , and w_ν such that $\sum_{t=1}^\infty \mathbf{E}[d_{t:t+h_t-1}] = \infty$ (whereas $\sum_{t=1}^\infty \mathbf{E}[d_t] < \infty$). Show that condition $\sup_t h_t = \infty$ is not sufficient. Show that convergence can be (very) slow when h_t grows (very) fast? For this question to answer one has to define ‘(very) slows/fast’, e.g. as logarithmical/exponentially increasing, or slower/faster than any unbounded computable function. Is convergence reasonably fast for slowly growing horizon, e.g. for $h_t = \log t$? What about $h_t = t$?



Alan Turing
(1912–1954)

Napoleon: “How is it that, although you say so much about the Universe, you say nothing about its Creator?”

Laplace: “No, Sire, I had no need of that hypothesis.”

Lagrange: “Ah, but it is such a good hypothesis: it explains so many things!”

Laplace: “Indeed, Sire, Monsieur Lagrange has, with his usual sagacity, put his finger on the precise difficulty with the hypothesis: it explains everything, but predicts nothing.”

(Conversation between Laplace and Lagrange mediated by Napoleon)

Chapter 4

Agents in Known Probabilistic Environments

4.1	The $AI\mu$ Model in Functional Form	402
4.1.1	The Cybernetic Agent Model	402
4.1.2	Strings	403
4.1.3	AI model for Known Deterministic Environment	405
4.1.4	AI Model for Known Prior Probability	406
4.2	The $AI\mu$ Model in Recursive and Iterative Form	408
4.2.1	Probability Distributions	408
4.2.2	Explicit Form of the $AI\mu$ Model	409
4.2.3	Equivalence of Functional and Explicit AI model	410
4.3	Special Aspects of the $AI\mu$ Model	412
4.3.1	Factorizable Environments	412
4.3.2	Constants and Limits	414
4.3.3	Sequential Decision Theory	415
4.4	Problems	416

The general framework for AI might be viewed as the design and study of intelligent agents [RN95]. An agent is a cybernetic system with some internal state, which acts with output y_k on some environment in cycle k , perceives some input x_k from the environment and updates its internal state. Then the next cycle follows. We split the input x_k into a regular part x'_k and a reward r_k , often called reinforce-

ment feedback. From time to time the environment provides non-zero reward to the agent. The task of the agent is to maximize its utility, defined as the sum of future rewards. A probabilistic environment can be described by the conditional probability μ for the inputs $x_1 \dots x_n$ to the agent under the condition that the agent outputs $y_1 \dots y_n$. Most, if not all environments are of this type. We give formal expressions for the outputs of the agent, which maximize the total μ -expected reward sum, called value. This model is called the $\text{AI}\mu$ model. As every AI problem can be brought into this form, the problem of maximizing utility is hence being formally solved, if μ is known. Furthermore, we study some special aspects of the $\text{AI}\mu$ model. We introduce factorizable probability distributions describing environments with independent episodes. They occur in several problem classes studied in Chapter 6 and are a special case of more general separable probability distributions defined in Section 5.3. We also clarify the connection to the Bellman equations of sequential decision theory and discuss similarities and differences. We discuss minor parameters of our model, including (the size of) the input and output spaces \mathcal{X} and \mathcal{Y} and the lifetime of the agent, and their universal choice, which we have in mind. There is nothing remarkable in this chapter, it is the essence of sequential decision theory [NM44, Bel57, BT96, SB98], presented in a more general form. Notation and formulas needed in later sections are simply developed. There are two major remaining problems. The problem of the unknown true probability distribution μ , which is solved in Chapter 5, and computational aspects, which are addressed in Chapter 7.

4.1 The $\text{AI}\mu$ Model in Functional Form

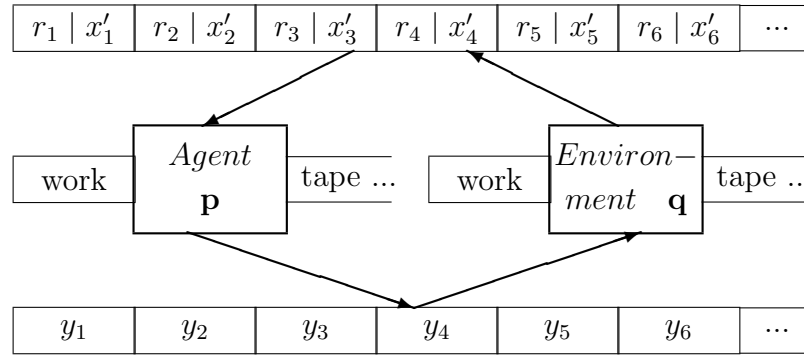
4.1.1 The Cybernetic Agent Model

A good way to start thinking about intelligent systems is to consider more generally cybernetic systems, in AI usually called agents. This avoids having to struggle with the meaning of intelligence from the very beginning. A cybernetic system is a control circuit with input x and output y and an internal state. From an external input and the internal state the agent calculates deterministically or stochastically an output. This output (action) modifies the environment and leads to a new input (perception). This continues ad infinitum or for a finite number of cycles.

Definition 4.1 (The Agent Model) An agent is a system which interacts with an environment in cycles $k = 1, 2, 3, \dots$. In cycle k the action (output) $y_k \in \mathcal{Y}$ of the agent is determined by a policy p which depends on the I/O-history $y_1 x_1 \dots y_{k-1} x_{k-1}$. The environment reacts to this action and leads to a new perception (input) $x_k \in \mathcal{X}$ determined by a deterministic function q or probability distribution μ which depends on the history $y_1 x_1 \dots y_{k-1} x_{k-1} y_k$. Then the next cycle $k+1$ starts.

There is significant overlap between control theory studied by engineers and agent theory studied in AI, but both fields differ in notation and emphasis. Table 4.2 compares notation and emphasis in both fields. Only the interchange of input \leftrightarrow output can cause confusion. With few exceptions we use the notion common in AI.

As explained in the last section, we need some reward assignment to the cybernetic system. The input x is divided into two parts, the standard input x' and some reward input r . If input and output are represented by strings, a deterministic cybernetic system can be modeled by a Turing machine p . p is called the policy of the agent, which determines the (re)action to a perception. If the environment is also computable it might be modeled by a Turing machine q as well. The interaction of the agent with the environment can be illustrated as follows:



p as well as q have unidirectional input and output tapes and bidirectional work tapes. What entangles the agent with the environment, is the fact that the upper tape serves as input tape for p , as well as output tape for q , and that the lower tape serves as output tape for p as well as input tape for q . Further, the reading head must always be left of the writing head, i.e. the symbols must first be written, before they are read. p and q have their own mutually inaccessible work tapes containing their own ‘secrets’. The heads move in the following way. In the k^{th} cycle p writes y_k , q reads y_k , q writes $x_k \equiv r_k x'_k$, p reads $x_k \equiv r_k x'_k$, followed by the $(k+1)^{th}$ cycle and so on. The whole process starts with the first cycle, all heads on tape start and work tapes being empty. We want to call Turing machines behaving in this way, *chronological Turing machines*. Before continuing, some notations on strings and probability distributions are appropriate.

4.1.2 Strings

We will denote strings over the alphabet \mathcal{X} by $s = x_1 x_2 \dots x_n$, with $x_k \in \mathcal{X}$, where \mathcal{X} is alternatively interpreted as a non-empty subset of \mathbb{N} or itself as a prefix free set of binary strings. $l(s) = l(x_1) + \dots + l(x_n)$ is the length of s . Analogous definitions hold for $y_k \in \mathcal{Y}$. We call x_k the k^{th} input word and y_k the k^{th} output word (rather than letter). The string $s = y_1 x_1 \dots y_n x_n$ represents the input/output in chronological order. Due to the prefix property of the x_k and y_k , s can be uniquely separated into its

Table 4.2 (Notation and emphasis in AI versus control theory) The upper part ($\hat{=}$) of the table compares notation used in AI or reinforcement learning to notation used in control theory. The lower part (\leftrightarrow) compares the objectives of both fields.

artificial intelligence		control theory
agent	$\hat{=}$	controller
environment	$\hat{=}$	system
policy	$\hat{=}$	control=policy
transition matrix	$\hat{=}$	transition matrix?
observation=input=perception	$\hat{=}$	output
action=output	$\hat{=}$	input
(instantaneous) reward	$\hat{=}$	immediate or one-period cost
cumulative reward=value	$\hat{=}$	expected (total) cost(-to-go)
model learning	$\hat{=}$	system identification
exploitation	$\hat{=}$?	(optimal?) stochastic control?
reactive agent	$\hat{=}$	closed loop control
prewired agent?	$\hat{=}$	open loop control
Markov decision process	$\hat{=}$	controlled Markov chain
belief state	$\hat{=}$	information state
Bellman equationBellman equations	$\hat{=}$	Bellman equation (Ricatti eq.?)
reinforcement learning	$\hat{=}$	sequential decision theory
reinforcement learning	$\hat{=}$	adaptive control
?	$\hat{=}$	consistent control
?	$\hat{=}$	self-tuning control
?	$\hat{=}$	self-optimizing control
exploration \leftrightarrow exploitation problem	$\hat{=}$	estimation \leftrightarrow control problem
qualitative solution	\leftrightarrow	high precision
complex environment	\leftrightarrow	simple machine
temporal difference learning	\leftrightarrow	value/policy iteration
temporal difference learning?	\leftrightarrow	dynamic programming

words. The words appearing in strings are always in chronological order. We further introduce the following abbreviations: ϵ is the empty string, $x_{n:m} := x_n x_{n+1} \dots x_{m-1} x_m$ for $n \leq m$ and ϵ for $n > m$. $x_{<n} := x_1 \dots x_{n-1}$. Analogously for y . Further, $y x_n := y_n x_n$, $y x_{n:m} := y_n x_n \dots y_m x_m$, and so on.

4.1.3 AI model for Known Deterministic Environment

Let us define for the chronological Turing machine p a partial function also named $p : \mathcal{X}^* \rightarrow \mathcal{Y}^*$ with $y_{1:k} = p(x_{<k})$ where $y_{1:k}$ is the output of Turing machine p on input $x_{<k}$ in cycle k , i.e. where p has read up to x_{k-1} but no further¹. In an analogous way, we define $q : \mathcal{Y}^* \rightarrow \mathcal{X}^*$ with $x_{1:k} = q(y_{1:k})$. Conversely, for every partial recursive chronological function we can define a corresponding chronological Turing machine. Each (agent, environment) pair (p, q) produces a unique I/O sequence $\omega^{pq} := y_1^{pq} x_1^{pq} y_2^{pq} x_2^{pq} \dots$. When we look at the definitions of p and q we see a nice symmetry between the cybernetic system and the environment. Until now, not much intelligence is in our agent. Now the credit assignment comes into the game and removes the symmetry somewhat. We split the input $x_k \in \mathcal{X} := \mathcal{R} \times \mathcal{X}'$ into a regular part $x'_k \in \mathcal{X}'$ and a reward $r_k \in \mathcal{R} \subset \mathbb{R}$. We define $x_k \equiv r_k x'_k$ and $r_k \equiv r(x_k)$. The goal of the agent should be to maximize received rewards. This is called reinforcement learning. The reason for the asymmetry is, that eventually we (humans) will be the environment with which the agent will communicate and *we* want to dictate what is good and what is wrong, not the other way round. This one way learning, the agent learns from the environment, and not conversely, neither prevents the agent from becoming more intelligent than the environment, nor does it prevent the environment learning from the agent because the environment can itself interpret the outputs y_k as a regular and a reward part. The environment is just not forced to learn, whereas the agent is. In cases where we restrict the reward to two values $r \in \mathcal{R} = \mathcal{B} := \{0, 1\}$, $r = 1$ is interpreted as a positive feedback, called *good* or *correct* and $r = 0$ a negative feedback, called *bad* or *error*. Further, let us restrict for a while the lifetime (number of cycles) m of the agent to a large, but finite value. Let $V_{km}^{pq} := \sum_{i=k}^m r(x_i)$ be the future total reward (called future utility), the agent p receives from the environment q in the cycles k to m . It is now natural to call the agent p^* , which maximizes V_{1m} (called total utility), the *best* one:²

$$p^* := \arg \max_p V_{1m}^{pq} \quad \Rightarrow \quad V_{km}^{p^*q} \geq V_{km}^{pq} \quad \forall p : y_{<k}^{pq} = y_{<k}^{p^*q}. \quad (4.3)$$

For $k = 1$ the condition on p is nil. For $k > 1$ it states that p shall be consistent with p^* in the sense that they have the same history. If \mathcal{X} , \mathcal{Y} and m are finite, the number of different behaviors of the agent, i.e. the search space is finite. Therefore,

¹Note that a possible additional dependence of p on $y_{<k}$ as mentioned in Definition 4.1 can be eliminated by recursive substitution; see below. Similarly for q .

² $\arg \max_p V(p)$ is the p which maximizes $V(\cdot)$. If there is more than one maximum we might choose the lexicographically smallest one for definiteness.

because we have assumed that q is known, p^* can effectively be determined (by pre-analyzing all behaviors). The main reason for restricting to finite m was not to ensure computability of p^* but that the limit $m \rightarrow \infty$ might not exist. The ease with which we defined and computed the optimal policy p^* is not remarkable. Just the (unrealistic) assumption of a completely known deterministic environment q has trivialized everything.

4.1.4 AI Model for Known Prior Probability

Let us now weaken our assumptions by replacing the environment q with a probability distribution $\mu(q)$ over chronological functions. μ might be interpreted in two ways. Either the environment itself behaves stochastically defined by μ or the true environment is deterministic, but we only have subjective (probabilistic) information, of which environment being the true environment. Combinations of both cases are also possible. We assume here that μ is known and describes the true stochastic behavior of the environment. The case of unknown μ with the agent having some beliefs about the environment lies at the heart of the AI ξ model described in Chapter 5.

The *best* agent is now the one which maximizes the *expected* utility (called value function) $V_\mu^p \equiv V_{1m}^{p\mu} := \sum_q \mu(q) V_{1m}^{pq}$. This defines the AI μ model.

Definition 4.4 (The AI μ model) The AI μ model is the agent with policy p^μ which maximizes the μ -expected total reward $r_1 + \dots + r_m$, i.e. $p^* \equiv p^\mu := \arg\max_p V_\mu^p$. $V_\mu^* := V_\mu^{p^\mu}$.

We need the concept of a *value function* in a slightly more general form.

Definition 4.5 (The μ /true/generating value function) The agent's perception x consists of a regular input $x' \in \mathcal{X}'$ and a reward $r \in \mathcal{R} \subset \mathbb{R}$. In cycle k the *value* $V_{km}^{p\mu}(y_{<k})$ is defined as the μ -expectation of the future reward sum $r_k + \dots + r_m$ with actions generated by policy p , and fixed history $y_{<k}$. We say that $V_{km}^{p\mu}(y_{<k})$ is the (future) *value* of policy p in environment μ given history $y_{<k}$, or shorter, the μ or true or generating value of p given $y_{<k}$. $V_\mu^p := V_{1m}^{p\mu}$ is the (total) value of p .

We now give a more formal definition for $V_{km}^{p\mu}$. Let us assume we are in cycle k with history $\dot{y}_1 \dots \dot{y}_{k-1}$ and ask for the *best* output y_k . Further, let $\dot{Q}_k := \{q : q(\dot{y}_{<k}) = \dot{x}_{<k}\}$ be the set of all environments producing the above history. We say that $q \in \dot{Q}_k$ is *consistent* with history $\dot{y}_{<k}$. The expected reward for the next $m - k + 1$ cycles (given the above history) is called the value of policy p and is given by a conditional probability:

$$V_{km}^{p\mu}(\dot{y}_{<k}) := \frac{\sum_{q \in \dot{Q}_k} \mu(q) V_{km}^{pq}}{\sum_{q \in \dot{Q}_k} \mu(q)}. \quad (4.6)$$

Policy p and environment μ do not determine history $\dot{y}_{<k}$ unlike the deterministic case because the history is no longer deterministically determined by p and q , but depends on p and μ *and* on the outcome of a stochastic process. Every new cycle adds new information (\dot{x}_i) to the agent. This is indicated by the dots over the symbols. In cycle k we have to maximize the expected future rewards, taking into account the information in the history $\dot{y}_{<k}$. This information is not already present in p and q/μ at the agent's start unlike in the deterministic case.

Furthermore, we want to generalize the finite lifetime m to a dynamical (computable) farsightedness $h_k \equiv m_k - k + 1 \geq 1$, called horizon. For $m_k = m$ we have our original finite lifetime, for $h_k = h$ the agent maximizes in every cycle the next h expected rewards. A discussion of the choices for m_k is delayed to Section 5.7.

The next h_k rewards are maximized by

$$p_k^* := \arg \max_{p \in \dot{P}_k} V_{km_k}^{p\mu}(\dot{y}_{<k}),$$

where $\dot{P}_k := \{p : \exists y_k : p(\dot{x}_{<k}) = \dot{y}_{<k} y_k\}$ is the set of policies *consistent* with the current history. p_k^* depends on k and is used only in step k to determine \dot{y}_k by $p_k^*(\dot{x}_{<k} | \dot{y}_{<k}) = \dot{y}_{<k} \dot{y}_k$. After writing \dot{y}_k the environment replies with \dot{x}_k with (conditional) probability $\mu(\dot{Q}_{k+1})/\mu(\dot{Q}_k)$. This probabilistic outcome provides new information to the agent. The cycle $k+1$ starts with determining \dot{y}_{k+1} from p_{k+1}^* (which differs from p_k^* as \dot{x}_k is now fixed) and so on. Note that p_k^* implicitly depends also on $\dot{y}_{<k}$ because \dot{P}_k and \dot{Q}_k do so. But recursively inserting p_{k-1}^* and so on, we can define

$$p^*(\dot{x}_{<k}) := p_k^*(\dot{x}_{<k} | p_{k-1}^*(\dot{x}_{<k-1} | \dots p_1^*)). \quad (4.7)$$

It is a chronological function and computable if \mathcal{X} , \mathcal{Y} and m_k are finite and μ is computable. For constant m one can show that the policy (4.7) coincides with the AI μ model (Definition 4.4). This also proves

$$V_{km}^{*p\mu}(\dot{y}_{<k}) \geq V_{km}^{p\mu}(\dot{y}_{<k}) \quad \forall p \in \dot{P}_k, \quad \text{i.e. for all } p \text{ consistent with } \dot{y}_{<k} \quad (4.8)$$

similarly to (4.3) (see Problem 4.1). For $k=1$ this is obvious. We also call (4.7) AI μ model. For deterministic³ μ this model reduces to the deterministic case discussed in the last subsection.

It is important to maximize the sum of future rewards and not, for instance, to be greedy and only maximize the next reward, as is done e.g. in sequence prediction. For example, let the environment be a sequence of chess games and each cycle corresponds to one move. Only at the end of each game a positive reward $r=1$ is given to the agent if it won the game (and made no illegal move). For the agent, maximizing all future rewards means trying to win as many games in as short as possible time (and avoiding illegal moves). The same performance is reached, if we choose h_k much larger than the typical game lengths. Maximization of only the next

³We call a probability distribution deterministic if it assumes values 0 and 1 only.

reward would be a very bad chess playing agent. Even if we would make our reward r finer, e.g. by evaluating the number of chessmen, the agent would play very bad chess for $h_k=1$, indeed.

The $\text{AI}\mu$ model still depends on μ and m_k . m_k is addressed in Section 5.7. To get our final universal AI model the idea is to replace μ by the universal probability ξ , defined later. This is motivated by the fact that ξ converges to μ in a certain sense for any μ . With ξ instead of μ our model no longer depends on any parameters, so it is truly universal. It remains to show that it behaves intelligently. But let us continue step by step. In the following we develop an alternative but equivalent formulation of the $\text{AI}\mu$ model. Whereas the functional form presented above is more suitable for theoretical considerations, especially for the development of a time bounded version in Section 7.2, the iterative/recursive formulation of the next subsections will be more appropriate for the explicit calculations in most of the other sections.

4.2 The $\text{AI}\mu$ Model in Recursive and Iterative Form

4.2.1 Probability Distributions

We use Greek letters for probability distributions and underline their arguments to indicate that they are probability arguments. Let $\rho_n(\underline{x}_1 \dots \underline{x}_n)$ be the probability that an (infinite) string starts with $x_1 \dots x_n$. We drop the index on ρ if it is clear from its arguments:

$$\sum_{x_n \in \mathcal{X}} \rho(\underline{x}_{1:n}) \equiv \sum_{x_n} \rho_n(\underline{x}_{1:n}) = \rho_{n-1}(\underline{x}_{<n}) \equiv \rho(\underline{x}_{<n}), \quad \rho(\epsilon) \equiv \rho_0(\epsilon) = 1. \quad (4.9)$$

We also need conditional probabilities. We prefer a notation which preserves the chronological order of the words, in contrast to the standard notation $\rho(\cdot|\cdot)$ which flips it. We extend the definition of ρ to the conditional case with the following convention for its arguments: An underlined argument \underline{x}_k is a probability variable and other non-underlined arguments x_k represent conditions. With this convention, the conditional probability has the form $\rho(x_{<n} \underline{x}_n) = \rho(\underline{x}_{1:n}) / \rho(\underline{x}_{<n})$. The equation states that the probability that a string $x_1 \dots x_{n-1}$ is followed by x_n is equal to the probability of $x_1 \dots x_n^*$ divided by the probability of $x_1 \dots x_{n-1}^*$. We use x^* as an abbreviation for ‘strings starting with x ’.

The introduced notation is also suitable for defining the conditional probability $\rho(y_1 \underline{x}_1 \dots y_n \underline{x}_n)$ that the environment reacts with $x_1 \dots x_n$ under the condition that the output of the agent is $y_1 \dots y_n$. The environment is chronological, i.e. input x_i depends on $y_{<i}$ only. In the probabilistic case this means that $\rho(y_{<k} \underline{x}_k) := \sum_{x_k} \rho(y_{1:k})$ is independent of y_k , hence a tailing y_k in the arguments of ρ can be dropped. Probability distributions with this property will be called *chronological*. The y are always conditions, i.e. never underlined, whereas additional conditioning for the x

can be obtained with the chain rule

$$\rho(\underline{y}_{<n}\underline{y}_n) = \rho(\underline{y}_{1:n})/\rho(\underline{y}_{<n}) \quad \text{and} \quad (4.10)$$

$$\rho(\underline{y}_{1:n}) = \rho(\underline{y}_1) \cdot \rho(\underline{y}_1\underline{y}_2) \cdot \dots \cdot \rho(\underline{y}_{<n}\underline{y}_n). \quad (4.11)$$

The second equation is the first equation applied n times.

4.2.2 Explicit Form of the AI μ Model

Let us define the AI μ model p^* in a different way. In the next subsection we will show that the p^* model defined here is identical to the functional definition of p^* given in the last section.

Let $\mu(\underline{y}_{<k}\underline{y}_k)$ be the true probability of input x_k in cycle k , given the history $\underline{y}_{<k}y_k$. $\mu(\underline{y}_{1:k})$ is the true chronological prior probability that the environment reacts with $x_{1:k}$ if provided with actions $y_{1:k}$ from the agent. We assume the cybernetic model depicted on page 403 to be valid. Next we define the value $V_{k+1,m}^{*\mu}(\underline{y}_{1:k})$ to be the μ -expected reward sum $r_{k+1} + \dots + r_m$ in cycles $k+1$ to m with outputs y_i generated by agent p^* , which maximizes the expected reward sum, and responses x_i from the environment, drawn according to μ . Adding $r(x_k) \equiv r_k$ we get the reward including cycle k . The probability of x_k , given $\underline{y}_{<k}y_k$, is given by the conditional probability $\mu(\underline{y}_{<k}\underline{y}_k)$. So the expected reward sum in cycles k to m given $\underline{y}_{<k}y_k$ is

$$V_{km}^{*\mu}(\underline{y}_{<k}y_k) := \sum_{x_k} [r(x_k) + V_{k+1,m}^{*\mu}(\underline{y}_{1:k})] \cdot \mu(\underline{y}_{<k}\underline{y}_k) \quad (4.12)$$

Now we ask about how p^* chooses y_k . It should choose y_k as to maximize the future rewards. So the expected reward in cycles k to m given $\underline{y}_{<k}$ and y_k chosen by p^* is $V_{km}^{*\mu}(\underline{y}_{<k}) := \max_{y_k} V_{km}^{*\mu}(\underline{y}_{<k}y_k)$ (see Figure 4.13). Together with the induction start

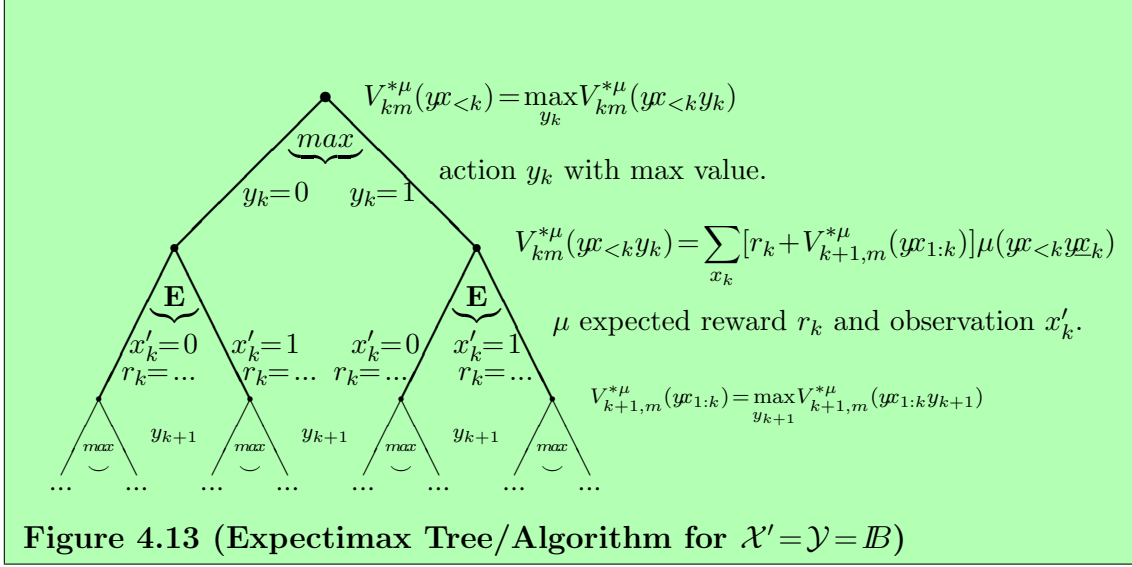
$$V_{m+1,m}^{*\mu}(\underline{y}_{1:m}) := 0 \quad (4.14)$$

$V_{km}^{*\mu}$ is completely defined. We might summarize one cycle into the formula

$$V_{km}^{*\mu}(\underline{y}_{<k}) = \max_{y_k} \sum_{x_k} [r(x_k) + V_{k+1,m}^{*\mu}(\underline{y}_{1:k})] \cdot \mu(\underline{y}_{<k}\underline{y}_k) \quad (4.15)$$

We introduce a dynamical (computable) farsightedness $h_k \equiv m_k - k + 1 \geq 1$, called horizon. For $m_k = m$, where m is the lifetime of the agent, we achieve optimal behavior, for limited farsightedness $h_k = h$ ($m = m_k = h + k - 1$), the agent maximizes in every cycle the next h expected rewards. A discussion of the choices for m_k is delayed to Section 5.7. If m_k is our horizon function of p^* and $\dot{\underline{y}}_{<k}$ is the actual history in cycle k , the output \dot{y}_k of the agent is explicitly given by

$$\dot{y}_k = \arg \max_{y_k} V_{km_k}^{*\mu}(\dot{\underline{y}}_{<k}y_k) \quad (4.16)$$

Figure 4.13 (Expectimax Tree/Algorithm for $\mathcal{X}' = \mathcal{Y} = \mathcal{B}$)

which in turn defines the policy p^* . Then the environment responds \dot{x}_k with probability $\mu(\dot{y}_{<k} \dot{y}_k)$. Then cycle $k+1$ starts. We might unfold the recursion (4.15) further and give \dot{y}_k non-recursively as

$$\dot{y}_k \equiv \dot{y}_k^\mu := \arg \max_{y_k} \sum_{x_k} \max_{y_{k+1}} \sum_{x_{k+1}} \dots \max_{y_{m_k}} \sum_{x_{m_k}} (r(x_k) + \dots + r(x_{m_k})) \cdot \mu(\dot{y}_{<k} \underline{y}_{k:m_k}) \quad (4.17)$$

This has a direct interpretation: the probability of inputs $x_{k:m_k}$ in cycle k when the agent outputs $y_{k:m_k}$ with actual history $\dot{y}_{<k}$ is $\mu(\dot{y}_{<k} \underline{y}_{k:m_k})$. The future reward in this case is $r(x_k) + \dots + r(x_{m_k})$. The best expected reward is obtained by averaging over the x_i (\sum_{x_i}) and maximizing over the y_i . This has to be done in chronological order to correctly incorporate the dependency of x_i and y_i on the history. This is essentially the expectimax algorithm/tree [Mic66, RN95]. The $\text{AI}\mu$ model is *optimal* in the sense that no other policy leads to higher expected reward. The value for a general policy p can be written in the form

$$V_{km}^{p\mu}(y_{<k}) := \sum_{x_{1:m}} (r_k + \dots + r_m) \mu(y_{<k} \underline{y}_{k:m})|_{y_{1:m}=p(x_{<m})} \quad (4.18)$$

4.2.3 Equivalence of Functional and Explicit AI model

As is clear from their interpretations, the iterative environmental probability μ relates to the functional form in the following way:

$$\mu(\underline{y}_{1:k}) = \sum_{q: q(y_{1:k})=x_{1:k}} \mu(q) \quad (4.19)$$

Theorem 4.20 (Equivalence of functional and explicit AI model) The actions of the functional AI model (4.7) coincide with the actions of the explicit (iterative/recursive) AI model (4.15-4.17) with environments identified by (4.19).

Proof. We will prove the equivalence of (4.7) and (4.16) only for $k=2$ and $m_2=3$. The proof of the general case is completely analogously except that the notation becomes quite messy.

Let us first evaluate (4.6) for fixed $\dot{y}_1\dot{x}_1$ and some $p \in \dot{P}_2$, i.e. $p(\dot{x}_1) = \dot{y}_1 y_2$ for some y_2 . If the next input to the agent is x_2 , p will respond with $p(\dot{x}_1 x_2) = \dot{y}_1 y_2 y_3$ for some y_3 depending on x_2 . We write $y_3(x_2)$ in the following⁴. The numerator of (4.6) simplifies to

$$\begin{aligned} \sum_{q \in \dot{Q}_2} \mu(q) V_{23}^{pq} &= \sum_{q: q(\dot{y}_1) = \dot{x}_1} \mu(q) V_{23}^{pq} = \sum_{x_2 x_3} (r(x_2) + r(x_3)) \sum_{q: q(\dot{y}_1 y_2 y_3(x_2)) = \dot{x}_1 x_2 x_3} \mu(q) = \\ &= \sum_{x_2 x_3} (r(x_2) + r(x_3)) \cdot \mu(\dot{y}_1 \dot{x}_1 y_2 \underline{x}_2 y_3(x_2) \underline{x}_3) \end{aligned}$$

In the first equality we inserted the definition of \dot{Q}_2 . In the second equality we split the sum over q by first summing over q with fixed $x_2 x_3$. This allows us to pull $V_{23} = r(x_2) + r(x_3)$ out of the inner sum. Then we sum over $x_2 x_3$. Further, we have inserted p , i.e. replaced p by y_2 and $y_3(\cdot)$. In the last equality we used (4.19). The denominator reduces to

$$\sum_{q \in \dot{Q}_2} \mu(q) = \sum_{q: q(\dot{y}_1) = \dot{x}_1} \mu(q) = \mu(\dot{y}_1 \dot{x}_1).$$

For the quotient we get

$$V_{23}^{p\mu}(\dot{y}_1 \dot{x}_1) = \sum_{x_2 x_3} (r(x_2) + r(x_3)) \cdot \mu(\dot{y}_1 \dot{x}_1 y_2 \underline{x}_2 y_3(x_2) \underline{x}_3).$$

We have seen that the relevant behavior of $p \in \dot{P}_2$ in cycle 2 and 3 is completely determined by y_2 and the function $y_3(\cdot)$.

$$\begin{aligned} \max_{p \in \dot{P}_2} V_{23}^{p\mu}(\dot{y}_1 \dot{x}_1) &= \max_{y_2} \max_{y_3(\cdot)} \sum_{x_2 x_3} (r(x_2) + r(x_3)) \cdot \mu(\dot{y}_1 \dot{x}_1 y_2 \underline{x}_2 y_3(x_2) \underline{x}_3) = \\ &= \max_{y_2} \sum_{x_2} \max_{y_3} \sum_{x_3} (r(x_2) + r(x_3)) \cdot \mu(\dot{y}_1 \dot{x}_1 y_2 \underline{x}_2 y_3 \underline{x}_3) \end{aligned}$$

In the last equality we have used the fact that the functional minimization over $y_3(\cdot)$ reduces to a simple minimization over the word y_3 when interchanging with the sum over its arguments ($\max_{y_3(\cdot)} \sum_{x_2} \equiv \sum_{x_2} \max_{y_3}$). In the functional case \dot{y}_2 is therefore determined by

$$\dot{y}_2 = \arg \max_{y_2} \sum_{x_2} \max_{y_3} \sum_{x_3} (r(x_2) + r(x_3)) \cdot \mu(\dot{y}_1 \dot{x}_1 y_2 \underline{x}_2 y_3 \underline{x}_3)$$

This is identical to the iterative definition (4.17) with $k=2$ and $m_2=3$. \square

⁴Dependency on dotted words like \dot{x}_1 is not shown as the dotted words are fixed.

4.3 Special Aspects of the AI_μ Model

4.3.1 Factorizable Environments

Up to now we have made no restrictions on the form of the prior probability μ apart from being a chronological probability distribution. On the other hand, we will see that, in order to prove rigorous reward bounds, the prior probability must satisfy some separability condition to be defined later. Here we introduce a very strong form of separability, when μ factorizes into products. We start with a factorization into two factors. Let us assume that μ is of the form

$$\mu(\underline{y}_{1:n}) = \mu_1(\underline{y}_{<l}) \cdot \mu_2(\underline{y}_{l:n}) \quad (4.21)$$

for some fixed l and sufficiently large $n \geq m_k$. For this μ the output \dot{y}_k in cycle k of the AI_μ agent (4.17) for $k \geq l$ depends on $\dot{y}_{l:k-1}$ and μ_2 only and is independent of $\dot{y}_{<l}$ and μ_1 . This is easily seen when inserting

$$\mu(\dot{y}_{<k} \underline{y}_{k:m_k}) = \underbrace{\mu_1(\dot{y}_{<l})}_{\equiv 1} \cdot \mu_2(\dot{y}_{l:k-1} \underline{y}_{k:m_k}) \quad (4.22)$$

into (4.17). For $k < l$ the output \dot{y}_k depends on $\dot{y}_{<k}$ (this is trivial) and μ_1 only (trivial if $m_k < l$) and is independent of μ_2 . The non-trivial case, where the horizon $m_k \geq l$ reaches into the region μ_2 , can be proven as follows (we abbreviate $m := m_k$ in the following). Inserting (4.21) into the definition of $V_{lm}^{*\mu}(\underline{y}_{<l})$, the factor μ_1 is 1 as in (4.22). We abbreviate $V_{lm}^{*\mu} := V_{lm}^{*\mu}(\underline{y}_{<l})$ as it is independent of its arguments. One can decompose

$$V_{km}^{*\mu}(\underline{y}_{<k}) = V_{k,l-1}^{*\mu}(\underline{y}_{<k}) + V_{lm}^{*\mu} \quad (4.23)$$

For $k=l$ this is true because the first term on the r.h.s. is zero. For $k < l$ we prove the decomposition by induction from $k+1$ to k .

$$\begin{aligned} V_{km}^{*\mu}(\underline{y}_{<k}) &= \max_{y_k} \sum_{x_k} [r(x_k) + V_{k+1,l-1}^{*\mu}(\underline{y}_{1:k}) + V_{lm}^{*\mu}] \cdot \mu_1(\underline{y}_{<k} \underline{y}_k) = \\ &= \max_{y_k} \left[\sum_{x_k} (r(x_k) + V_{k+1,l-1}^{*\mu}(\underline{y}_{<k})) \cdot \mu_1(\underline{y}_{<k} \underline{y}_k) + V_{lm}^{*\mu} \right] = \\ &= V_{k,l-1}^{*\mu}(\underline{y}_{<k}) + V_{lm}^{*\mu} \end{aligned}$$

Inserting (4.23), valid for $k+1$ by induction hypothesis, into (4.15) gives the first equality. In the second equality we have performed the x_k sum for the $V_{lm}^{*\mu} \cdot \mu_1$ term which is now independent of y_k . It can therefore be pulled out of \max_{y_k} . In the last equality we used again the definition (4.15). This completes the induction step and proves (4.23) for $k < l$. \dot{y}_k can now be represented as

$$\dot{y}_k = \arg \max_{y_k} V_{km}^{*\mu}(\dot{y}_{<k} y_k) = \arg \max_{y_k} V_{k,l-1}^{*\mu}(\dot{y}_{<k} y_k) \quad (4.24)$$

where (4.16) and (4.23) and the fact that an additive constant $V_{lm}^{*\mu}$ does not change $\operatorname{argmax}_{y_k}$ has been used. $V_{k,l-1}^{*\mu}(\dot{y}_{<k}y_k)$ and hence \dot{y}_k is independent of μ_2 for $k < l$. Note, that \dot{y}_k is also independent of the choice of m , as long as $m \geq l$.

In the general case of an (infinite) sequence of consecutive episodes one can show an analogous result:

Theorem 4.25 (Factorizable environments μ) Assume that the cycles are grouped into independent episodes $r=1,2,3,\dots$, where each episode r consists of the cycles $k=n_r+1,\dots,n_{r+1}$ for some $0=n_0 < n_1 < \dots < n_s=n$:

$$\mu(\underline{y}_{1:n}) = \prod_{r=0}^{s-1} \mu_r(\underline{y}_{n_r+1:n_{r+1}}) \quad (4.26)$$

(In the simplest case, when all episodes have the same length l then $n_r=r \cdot l$) \dot{y}_k depends on μ_r and x and y of episode r only, with r such that $n_r < k \leq n_{r+1}$.

$$\dot{y}_k = \operatorname{argmax}_{y_k} \sum_{x_k} \dots \max_{y_t} \sum_{x_t} (r(x_k) + \dots + r(x_t)) \cdot \mu_r(\dot{y}_{n_r+1:k-1} \underline{y}_{k:t}) \quad (4.27)$$

with $t := \min\{m_k, n_{r+1}\}$. The different episodes are completely independent in the sense that the inputs x_k of different episodes are statistically independent and depend only on y_k of the same episode. The outputs y_k depend on the x and y of the corresponding episode r only, and are independent of the actual I/O of the other episodes.

If all episodes have a length of at most l , i.e. $n_{r+1} - n_r \leq l$ and if we choose the horizon h_k to be at least l , then $m_k \geq k + l - 1 \geq n_r + l \geq n_{r+1}$ and hence $t = n_{r+1}$ independent of m_k . This means that for factorizable μ there is no problem in taking the limit $m_k \rightarrow \infty$. Maybe this limit can also be performed in the more general case of a sufficiently separable μ . The (problem of the) choice of m_k will be discussed in more detail later.

Although factorizable μ are too restrictive to cover all AI problems, it often occurs in practice in the form of repeated problem solving, and hence, is worth being studied. For example, if the agent has to play games like chess repeatedly, or has to minimize different functions, the different games/functions might be completely independent, i.e. the environmental probability factorizes, where each factor corresponds to a game/function minimization. For details, see the appropriate sections on strategic games and function minimization.

Further, for factorizable μ it is probably easier to derive suitable reward bounds for the universal $AI\xi$ model defined in the next section, than for the separable cases which will be introduced later. This could be a first step toward a definition and proof for the general case of separable problems. One goal of this paragraph was to show, that the notion of a factorizable μ could be the first step toward a definition and analysis of the general case of separable μ .

4.3.2 Constants and Limits

We have in mind a universal agent with complex interactions that is as least as intelligent and complex as a human being. One might think of an agent whose input y_k comes from a digital video camera, the output x_k is some image to a monitor⁵, only for the rewards we might restrict to the most primitive binary one, i.e. $r_k \in \mathbb{B}$. So we think of the following constant sizes:

$$\begin{array}{ccccccc} 1 & \ll & \langle l(y_k x_k) \rangle & \ll & k & \leq & m & \ll & |\mathcal{Y} \times \mathcal{X}| \\ 1 & \ll & 2^{16} & \ll & 2^{24} & \leq & 2^{32} & \ll & 2^{65536} \end{array}$$

The first two limits say that the actual number k of inputs/outputs should be reasonably large, compared to the typical size $\langle l \rangle$ of the input/output words, which itself should be rather sizeable. The last limit expresses the fact that the total lifetime m (number of I/O cycles) of the agent is far too small to allow every possible input to occur, or to try every possible output, or to make use of identically repeated inputs or outputs. We do not expect any useful outputs for $k \lesssim \langle l \rangle$. More interesting than the lengths of the inputs is the complexity $K(x_1 \dots x_k)$ of all inputs until now, to be defined later. The environment is usually not “perfect”. The agent could either interact with a non-perfect human or tackle a non-deterministic world (due to quantum mechanics or chaos)⁶. In either case, the sequence contains some noise, leading to $K \propto \langle l \rangle \cdot k$. The complexity of the probability distribution of the input sequence is something different. We assume that this noisy world operates according to some simple computable rules. $K(\mu_k) \ll \langle l \rangle \cdot k$, i.e. the rules of the world can be highly compressed. We may allow environments in which new aspects appear for $k \rightarrow \infty$ causing a non-bounded $K(\mu_k)$.

In the following we never use these limits, except when explicitly stated. In some simpler models and examples the size of the constants will even violate these limits (e.g. $l(x_k) = l(y_k) = 1$), but it is the limits above that the reader should bear in mind. We are only interested in theorems which do not degenerate under the above limits. In order to avoid cumbersome convergence and existence considerations we make the following assumptions throughout this work.

Assumption 4.28 (Finiteness) We assume that

- the input/perception space \mathcal{X} is finite,
- the output/action space \mathcal{Y} is finite,
- the rewards are non-negative and bounded, i.e. $r_k \in \mathcal{R} \subseteq [0, r_{max}]$,
- the horizon m is finite.

Finite \mathcal{X} and bounded \mathcal{R} (each separately) ensure existence of μ -expectations, but are sometimes needed together, finite \mathcal{Y} ensures that $\arg\max_{y_k \in \mathcal{Y}} [\dots]$ exists,

⁵Humans can only simulate a screen as output device by drawing pictures.

⁶Whether there exist truly stochastic processes at all is a difficult question. At least the quantum indeterminacy comes very close to it.

i.e. that maxima are attained, finite m avoids various technical and philosophical problems (Section 5.7), positive rewards are needed for the time bounded AI ξ^{tl} model (Section 7.2). Many theorems can be generalized by relaxing some or all of the above finiteness assumptions.

4.3.3 Sequential Decision Theory

In the following we clarify the connection of (4.15) and (4.16) to the Bellman equations [Bel57, BT96] of sequential decision theory and discuss similarities and differences. We consider an MDP, where with probability M_{ij}^a , the agent under consideration should reach (environmental) state $j \in \mathcal{S}$ when taking action $a \in \mathcal{A}$ in (the current) state $i \in \mathcal{S}$. If the agent receives reward $R(i)$, the optimal policy p^* , maximizing expected utility (defined as sum of future rewards), and the utility $U(i)$ of policy p^* are

$$p^*(i) = \arg \max_a \sum_j M_{ij}^a U(j) \quad , \quad U(i) = R(i) + \max_a \sum_j M_{ij}^a U(j) \quad (4.29)$$

See [RN95, BT96] for details and further references. Let us identify

$$\begin{aligned} \mathcal{S} &= (\mathcal{Y} \times \mathcal{X})^*, \quad \mathcal{A} = \mathcal{Y}, \quad a = y_k, \quad M_{ij}^a = \mu(\mathcal{Y}_{<k} \mathcal{Y}_k), \quad p^*(i) = y_k, \\ i &= \mathcal{Y}_{<k}, \quad R(i) = r(x_{k-1}), \quad U(i) = V_{k-1,m}^*(\mathcal{Y}_{<k}) = r(x_{k-1}) + V_{km}^*(\mathcal{Y}_{<k}), \\ j &= \mathcal{Y}_{1:k}, \quad R(j) = r(x_k), \quad U(j) = V_{km}^*(\mathcal{Y}_{1:k}) = r(x_k) + V_{k+1,m}^*(\mathcal{Y}_{1:k}), \end{aligned}$$

where we further set $M_{ij}^a = 0$ if i is not a starting substring of j or if $a \neq y_k$. This ensures the sum over j in (4.29) to reduce to a sum over x_k . If we set $m_k = m$ and insert (4.12) into (4.16), it is easy to see that (4.29) coincides with (4.15) and (4.16).

Note that despite of this formal equivalence, we were forced to use the complete history $\mathcal{Y}_{<k}$ as environmental state i . The AI μ model neither assumes stationarity, nor Markov property, nor complete accessibility of the environment, as any assumption would restrict the applicability of AI μ . The consequence is that every state occurs at most once in the lifetime of the agent. Every moment in the universe is unique! Even if the state space could be identified with the input space \mathcal{X} , inputs would usually not occur twice by the assumption $k \ll |\mathcal{X}|$, made in the last subsection. Further, there is no (obvious) universal similarity relation on $(\mathcal{X} \times \mathcal{Y})^*$ allowing an effective reduction of the size of the state space. Although many algorithms (like value and policy iteration) have problems in solving (4.29) for huge or infinite state spaces in practice, there is no principle problem in determining p^* and U , as long as μ is known and computable and $|\mathcal{X}|$, $|\mathcal{Y}|$ and m are finite.

Things drastically change if μ is unknown. Reinforcement learning algorithms [KLM96] are commonly used in this case to learn the unknown μ . They succeed if the state space is either small or has effectively been made small by so called generalization techniques. In any case, the solutions are either ad hoc, or work in restricted domains only, or have serious problems with state space exploration versus

exploitation, are prone to diverge, or have non-optimal learning rate. There is no universal and optimal solution to this problem so far. In the next section we present a new model and argue that it formally solves all these problems in an optimal way. It will not concern with learning of μ directly. All we do is to replace the true prior probability μ by a universal probability ξ , which is shown to converge to μ in a sense.

4.4 Problems

4.1 (Value Dominance ρ) [C20s] In (4.5) $V_{km}^{p\mu}(y_{<k})$ has been defined as the μ -expected future reward sum $r_k + \dots + r_m$ with actions generated by policy p , and fixed history $y_{<k}$. The optimal policy p^μ was defined as the one with maximal total μ -value, i.e. $p^\mu := \operatorname{argmax}_p V_{1m}^{p\mu}(\epsilon)$. The corresponding value function is $V_{km}^{*\mu}(y_{<k}) := V_{km}^{p^\mu}(y_{<k})$. Obviously $V_{1m}^{*\mu} \geq V_{1m}^{p\mu}(\epsilon) \forall p$. Show that this implies $V_{km}^{*\mu}(y_{<k}) \geq V_{km}^{p\mu}(y_{<k}) \forall p$ consistent with $y_{<k}$. The derivation of a result of this form goes in hand with the derivation of Bellman's equations [Ber95a].

4.2 (Probabilistic policies) [C15usi] In this chapter we only gave formal definitions of the value function for deterministic policies, but allowed probabilistic environments. Generalize the definition of the value function (4.18) to probabilistic policies π in the following way:

$$\begin{aligned} V_\mu^\pi &= \sum_{y_{1:m}} (r_1 + \dots + r_m) \mu(x_m | y_{<m} y_m) \pi(y_m | y_{<m}) \dots \mu(x_1 | y_1) \pi(y_1) = \\ &= \sum_{y_{1:m}} (r_k + \dots + r_m) \mu(x_{1:m} | y_{1:m}) \pi(y_{1:m} | x_{<m}) = \sum_{y_{1:m}} (r_k + \dots + r_m) \nu(y_{1:m} x_{1:m}) \end{aligned}$$

and similarly for $V_{km}^{\pi\mu}(y_{<k})$. We used here the conventional notation $\mu(\cdot|\cdot)$ for conditional probabilities to emphasis the following oddity. The last equality seems to say something like $p(x|y)p(y|x) = p(x\&y)$, which would contradict Bayes law $p(x|y)p(y) = p(x\&y)$. Show that everything goes right here. What important property (defined in the main text) is used to arrive at the above expressions? Show that $\max_\pi V_\mu^\pi = \max_p V_\mu^p$ and that among optimal policies there is always a deterministic one. Since we are mainly interested in optimal policies, the restriction to deterministic policies is not serious. Generalize the theorems in this and later chapters involving V_μ^p to V_μ^π . A serious treatment of probabilistic policies can be found in game theory, where optimal policies can also be truly probabilistic [OR94].



Isaac Asimov
(1920-1992)

The Three Laws of Robotics:

1. *A robot may not injure a human being, or, through inaction, allow a human being to come to harm.*
2. *A robot must obey the orders given it by human beings except where such orders would conflict with the First Law.*
3. *A robot must protect its own existence as long as such protection does not conflict with the First or Second Law.*

(Asimov, 1940)

Chapter 5

The Universal Algorithmic Agent AIXI

5.1	The Universal AIXI Model	502
5.1.1	Definition of the AIXI Model	503
5.1.2	Universality of ξ^{AI}	504
5.1.3	Convergence of ξ^{AI} to μ^{AI}	505
5.1.4	Intelligence Order Relation	506
5.2	On the Optimality of AIXI	507
5.3	Value Bounds and Separability Concepts	509
5.3.1	Introduction	509
5.3.2	(Pseudo) Passive μ and the HeavenHell Example	509
5.3.3	The OnlyOne Example	510
5.3.4	Asymptotic Learnability	511
5.3.5	Uniform μ	512
5.3.6	Other Concepts	512
5.3.7	Summary	512
5.4	Value Related Optimality Results	513
5.4.1	The $AI\rho$ Models: Preliminaries	513
5.4.2	Pareto Optimality of $AI\xi$	514
5.4.3	Self-optimizing Policy p^ξ w.r.t. Average Value	515

5.5	Discounted Future Value Function	518
5.6	Markov Decision Processes (MDP)	524
5.7	The Choice of the Horizon	528
5.8	Outlook	531
5.9	Conclusions	532
5.10	Functions \leadsto Chronological Semimeasures	532
5.11	Proof of the Entropy Inequality	534
5.12	History & References	535
5.13	Problems	536

Active systems, like game playing (SG) and optimization (FM), cannot be reduced to induction systems. The *main idea of this work* is to generalize universal induction to the general agent model described in Chapter 4. For this, we generalize ξ to include conditions and replace μ by ξ in the rational agent model, resulting in the $\text{AI}\xi$ model. In this way the problem that the true prior probability μ is usually unknown is solved. Convergence of $\xi \rightarrow \mu$ will be shown indicating that the $\text{AI}\xi$ model could behave optimally in any computable but unknown environment with reinforcement feedback.

The main focus of Chapter 5 is to investigate what we can expect from a universally optimal agent and to clarify the meanings of *universal*, *optimal*, etc. Similarly to the induction case it is convenient to consider a general mixture distribution ξ of a weighted sum of distributions $\nu \in \mathcal{M}$, where \mathcal{M} is any class of distributions including the true environment μ . We show that the Bayes-optimal policy p^ξ based on the mixture ξ is self-optimizing in the sense that the average value converges asymptotically for all $\mu \in \mathcal{M}$ to the optimal value achieved by the (infeasible) Bayes-optimal policy p^μ which knows μ in advance. We show that the necessary condition that \mathcal{M} admits self-optimizing policies at all, is also sufficient. No other structural assumptions are made on \mathcal{M} .

We show that Bandits, i.i.d. processes, classification tasks, certain classes of POMDPs, (k^{th} order) ergodic MDPs, factorizable environments, repeated games, and prediction problems admit self-optimizing policies. Unfortunately, the class \mathcal{M}_U of all enumerable semi-measures does *not* admit self-optimizing policies. This forces us to lower our expectation about universally optimal agents and to introduce other (weaker) performance measures. Finally, we show that p^ξ is Pareto-optimal in the sense that there is no other policy yielding higher or equal value in *all* environments $\nu \in \mathcal{M}$ and a strictly higher value in at least one. Pareto-optimality holds for any choice of \mathcal{M} (including \mathcal{M}_U).

5.1 The Universal AIXI Model

5.1.1 Definition of the AIXI Model

We have developed enough formalism to suggest our universal AIXI model. All we have to do is to suitably generalize the universal semimeasure $\xi = \xi_U$ from the last subsection and replace the true but unknown prior probability μ^{AI} in the $AI\mu$ model by this generalized $\xi^{AI} = \xi_U^{AI}$. In what sense this AIXI model is universal will be discussed later.

In the functional formulation we define the universal probability ξ^{AI} of an environment q just as $2^{-l(q)}$

$$\xi(q) := 2^{-l(q)}$$

The definition could not be easier!^{1,2} Collecting the formulas of Section 4.1 and replacing $\mu(q)$ by $\xi(q)$ we get the definition of the AIXI agent in functional form. Given the history $\dot{y}_{<k}$ the functional AIXI agent outputs

$$\dot{y}_k := \arg \max_{y_k} \max_{p: p(\dot{x}_{<k}) = \dot{y}_{<k} y_k} \sum_{q: q(\dot{y}_{<k}) = \dot{x}_{<k}} 2^{-l(q)} \cdot V_{km_k}^{pq} \quad (5.1)$$

in cycle k , where $V_{km_k}^{pq}$ is the total reward of cycles k to m_k when agent p interacts with environment q . We have dropped the denominator $\sum_q \mu(q)$ from (4.6) as it is independent of the $p \in \dot{P}_k$ and a constant multiplicative factor does not change $\arg \max_{y_k}$.

For the iterative formulation the universal probability ξ can be obtained by inserting the functional $\xi(q)$ into (4.19)

$$\xi(\underline{y}_{1:k}) = \sum_{q: q(y_{1:k}) = x_{1:k}} 2^{-l(q)} \quad (5.2)$$

Replacing μ by ξ in (4.17) the iterative AIXI agent outputs

$$\dot{y}_k \equiv \dot{y}_k^\xi := \arg \max_{y_k} \sum_{x_k} \max_{y_{k+1}} \sum_{x_{k+1}} \dots \max_{y_{m_k}} \sum_{x_{m_k}} (r(x_k) + \dots + r(x_{m_k})) \cdot \xi(\dot{y}_{<k} \underline{y}_{k:m_k}) \quad (5.3)$$

in cycle k given the history $\dot{y}_{<k}$.

One subtlety has been passed over. Like in the SP case, ξ is not a probability distribution but satisfies only the weaker inequalities

$$\sum_{x_n} \xi(\underline{y}_{1:n}) \leq \xi(\underline{y}_{<n}), \quad \xi(\epsilon) \leq 1 \quad (5.4)$$

Note, that the sum on the l.h.s. is *not* independent of y_n unlike for chronological probability distributions. Nevertheless, it is bounded by something (the r.h.s) which

¹It is not necessary to use $2^{-K(q)}$ or something similar as some readers may expect at this point. The reason is that for every program q there exists a functionally equivalent program q' with $K(q) = l(q')$.

²Here and later we identify objects with their coding relative to some fixed Turing machine U . For example, if q is a function $K(q) := K(\lceil q \rceil)$ with $\lceil q \rceil$ being a binary coding of q such that $U(\lceil q \rceil, y) = q(y)$. On the other hand, if q already is a binary string we define $q(y) := U(q, y)$.

is independent of y_n . The reason is that the sum in (5.2) runs over (partial recursive) chronological functions only and the functions q which satisfy $q(y_{1:n}) = x_{<n}$ are a subset of the functions satisfying $q(y_{<n}) = x_{<n}$. We will in general call functions satisfying (5.4) *chronological semimeasures*. The important point is that the conditional probabilities (4.10) are ≤ 1 like for true probability distributions.

The equivalence of the functional and iterative AI model proven in Section 4.2 is true for every chronological semimeasure ρ , especially for ξ , hence we can talk about *the* AIXI model in this respect. It (slightly) depends on the choice of the universal Turing machine. $l(\lceil q \rceil)$ is defined only up to an additive constant. The AIXI model also depends on the choice of $\mathcal{X} = \mathcal{R} \times \mathcal{X}'$ and \mathcal{Y} , but we do not expect any bias when the spaces are chosen sufficiently simple, e.g. all strings of length 2^{16} . Choosing \mathcal{N} as the word space would be ideal, but whether the maxima (suprema) exist in this case, has to be shown beforehand. The only non-trivial dependence is on the horizon function m_k which will be discussed later. So apart from m_k and unimportant details the AIXI agent is uniquely defined by (5.1) or (5.3). It doesn't depend on any assumption about the environment apart from being generated by some computable (but unknown!) probability distribution.

5.1.2 Universality of ξ^{AI}

In which sense the AIXI model is optimal will be clarified later. In this and the next subsection we show that ξ^{AI} defined in (5.2) is universal and converges to μ^{AI} analogous to the SP case (2.25) and (2.23). The proofs are generalizations from the SP case. The y are pure spectators and cause no difficulties in the generalization. The replacement of the binary alphabet \mathcal{B} used in SP by the (possibly infinite) alphabet \mathcal{X} is possible, but needs to be done with care. In (2.25) $U(p) = x*$ produces strings starting with x , whereas in (5.2) we can demand q to output exactly n words $x_{1:n}$ as q knows n from the number of input words $y_1 \dots y_n$. For proofs of (2.25) and (2.23) see [Sol78] and [LV92a].

There is an alternative definition of ξ which coincides with (5.2) within a multiplicative constant of $O(1)$,

$$\xi(y_{1:n}) \stackrel{\times}{=} \sum_{\rho} 2^{-K(\rho)} \rho(y_{1:n}) \quad (5.5)$$

where the sum runs over all enumerable chronological semimeasures. The $2^{-K(\rho)}$ weighted sum over probabilistic environments ρ , coincides with the sum over $2^{-l(q)}$ weighted deterministic environments q , as will be proven below. In Appendix 5.10 we show that an enumeration of all enumerable functions can be converted into an enumeration of enumerable chronological semimeasures ρ . $K(\rho)$ is co-enumerable, therefore ξ defined in (5.5) is itself enumerable. The representation (5.2) is also enumerable. As $\sum_{\rho} 2^{-K(\rho)} \leq 1$ and the ρ 's satisfy (5.4), ξ is a chronological semimeasure as well. If we pick one ρ in (5.5) we get the universality property “for free”

$$\xi(y_{1:n}) \stackrel{\times}{\geq} 2^{-K(\rho)} \rho(y_{1:n}) \quad (5.6)$$

ξ is a universal element in the sense of (5.6) in the set of all enumerable chronological semimeasures.

To prove universality of ξ in the form (5.2) we have to show that for every enumerable chronological semimeasure ρ there exists a Turing machine T with

$$\rho(\underline{y}_{1:n}) = \sum_{q: T(q, \underline{y}_{1:n}) = x_{1:n}} 2^{-l(q)} \quad \text{and} \quad l(T) \stackrel{\pm}{=} K(\rho). \quad (5.7)$$

A proof of (5.7) will be given elsewhere (see Problem 5.3). Given T the universality of ξ follows from

$$\xi(\underline{y}_{1:n}) = \sum_{q: U(q, \underline{y}_{1:n}) = x_{1:n}} 2^{-l(q)} \geq \sum_{q': U(Tq', \underline{y}_{1:n}) = x_{1:n}} 2^{-l(Tq')} = 2^{-l(T)} \sum_{q': T(q', \underline{y}_{1:n}) = x_{1:n}} 2^{-l(q')} \stackrel{\times}{=} 2^{-K(\rho)} \rho(\underline{y}_{1:n})$$

The first equality and (5.2) are identical by definition. In the inequality we have restricted the sum over all q to q of the form $q = Tq'$. The third relation is true as running U on Tz is a simulation of T on z . The last equality follows from (5.7). All enumerable, universal, chronological semimeasures coincide up to a multiplicative constant, as they mutually dominate each other. Hence, definitions (5.2) and (5.5) are, indeed, equivalent.

5.1.3 Convergence of ξ^{AI} to μ^{AI}

In Section 3.9.2 we have proven the following entropy inequality

$$\sum_{i=1}^N (y_i - z_i)^2 \leq \sum_{i=1}^N y_i \ln \frac{y_i}{z_i} \quad \text{with} \quad \sum_{i=1}^N y_i = 1, \quad \sum_{i=1}^N z_i \leq 1 \quad (5.8)$$

In Section 5.11 we give a different proof since it contains some ideas which could be interesting on their own sake.³ If we identify $N = |\mathcal{X}|$, $i = x_k$, $y_i = \mu(\underline{y}_{<k} \underline{y}_k)$ and $z_i = \xi(\underline{y}_{<k} \underline{y}_k)$, multiply both sides with $\mu(\underline{y}_{<k})$, take the sum over $x_{<k}$ and k and use the chain rule $\mu(\underline{y}_{<k}) \cdot \mu(\underline{y}_{<k} \underline{y}_k) = \mu(\underline{y}_{1:k})$ we get

$$\sum_{k=1}^n \sum_{x_{1:k}} \mu(\underline{y}_{<k}) \left(\mu(\underline{y}_{<k} \underline{y}_k) - \xi(\underline{y}_{<k} \underline{y}_k) \right)^2 \leq \sum_{k=1}^n \sum_{x_{1:k}} \mu(\underline{y}_{1:k}) \ln \frac{\mu(\underline{y}_{<k} \underline{y}_k)}{\xi(\underline{y}_{<k} \underline{y}_k)} = \dots \quad (5.9)$$

In the r.h.s. we can replace $\sum_{x_{1:k}} \mu(\underline{y}_{1:k})$ by $\sum_{x_{1:n}} \mu(\underline{y}_{1:n})$ as the argument of the logarithm is independent of $x_{k+1:n}$. The k sum can now be brought into the logarithm and converts to a product. Using the chain rule (4.10) for μ and ξ we get

$$\dots = \sum_{x_{1:n}} \mu(\underline{y}_{1:n}) \ln \prod_{k=1}^n \frac{\mu(\underline{y}_{<k} \underline{y}_k)}{\xi(\underline{y}_{<k} \underline{y}_k)} = \sum_{x_{1:n}} \mu(\underline{y}_{1:n}) \ln \frac{\mu(\underline{y}_{1:n})}{\xi(\underline{y}_{1:n})} \stackrel{+}{\leq} \ln 2 \cdot K(\mu) \quad (5.10)$$

³Actually a proof similar to 5.11 of an inequality similar to (5.8) was found first (in 1999 already) before the one of Section 3.9.2 and, in a sense, initiated the whole thesis.

where we have used the universality property (5.6) of ξ in the last step. The line of reasoning is identically to (3.18); the y are, again, pure spectators. This will change when we analyze loss/reward bounds analogous to Theorem 3.48.

(5.9,5.10) shows that the μ -expected squared difference of μ and ξ is finite for computable μ . This, in turn, shows that $\xi(yx_{<k}\underline{y}'_k)$ converges to $\mu(yx_{<k}\underline{y}'_k)$ for $k \rightarrow \infty$ with μ probability 1. If we take a finite product of ξ 's and use the chain rule, we see that also $\xi(yx_{<k}\underline{y}'_{k:k+r})$ converges to $\mu(yx_{<k}\underline{y}'_{k:k+r})$. More generally, by supplementing the results on multi-step predictions in Section 3.7.1 with action-conditions y we have⁴

$$\xi(yx_{<k}\underline{y}'_{k:m_k}) \xrightarrow{k \rightarrow \infty} \mu(yx_{<k}\underline{y}'_{k:m_k}) \begin{cases} \text{i.m.s.} & \text{if } h_k \equiv m_k - k + 1 \leq h_{\max} < \infty, \\ \text{i.m.} & \text{for general } h_k \equiv m_k - k + 1. \end{cases} \quad (5.11)$$

This gives hope that the outputs \dot{y}_k of the AIXI model (5.3) could converge to the outputs \dot{y}_k from the $\text{AI}\mu$ model (4.17).

We want to call an AI model *universal*, if it is μ independent (unbiased, model-free) and is able to solve any solvable problem and learn any learnable task. Further, we call a universal model *universally optimal*, if there is no program, which can solve or learn significantly faster (in terms of interaction cycles). As the AIXI model is parameter-free, ξ converges to μ (5.11), the $\text{AI}\mu$ model is itself optimal, and we expect no other model to converge faster to $\text{AI}\mu$ by analogy to SP (Theorem 3.48):

Claim 5.12 (We expect AIXI to be universally optimal)

This is our main claim. In a sense, the intention of the remaining sections is to define this statement more rigorously and to give further support.

5.1.4 Intelligence Order Relation

We define the ξ -expected reward in cycles k to m of a policy p similar to (4.6) and (5.1). We extend the definition to programs $p \notin \dot{P}_k$ which are not consistent with the current history.

$$V_{km}^{p\xi}(\dot{y}_{<k}) := \frac{1}{\mathcal{N}} \sum_{q:q(\dot{y}_{<k})=\dot{x}_{<k}} 2^{-l(q)} \cdot V_{km}^{\tilde{p}q} \quad (5.13)$$

The normalization \mathcal{N} is again only necessary for interpreting V_{km} as the expected reward but otherwise unneeded. For consistent policies $p \in \dot{P}_k$ we define $\tilde{p} := p$. For $p \notin \dot{P}_k$, \tilde{p} is a modification of p in such a way that its outputs are consistent with the current history $\dot{y}_{<k}$, hence $\tilde{p} \in \dot{P}_k$, but unaltered for the current and future cycles $\geq k$. Using this definition of V_{km} we could take the maximum over all policies p in (5.1), rather than only the consistent ones.

⁴Here, and everywhere else, with $\xi_k \rightarrow \mu_k$ we mean $\xi_k - \mu_k \rightarrow 0$, and not that μ_k (and ξ_k) itself converge to a limiting value.

Definition 5.14 (Intelligence order relation) We call a policy p *more or equally intelligent* than p' and write

$$p \succeq p' \quad :\Leftrightarrow \quad \forall k \forall \dot{\mathbf{x}}_{<k} : V_{km_k}^{p\xi}(\dot{\mathbf{x}}_{<k}) \geq V_{km_k}^{p'\xi}(\dot{\mathbf{x}}_{<k}).$$

i.e. if p yields in any circumstance higher ξ -expected reward than p' .

As the algorithm p^ξ behind the AIXI agent maximizes $V_{km_k}^{p\xi}$ we have $p^\xi \succeq p$ for all p . The AIXI model is hence the most intelligent agent w.r.t. \succeq . \succeq is a universal order relation in the sense that it is free of any parameters (except m_k) or specific assumptions about the environment. A proof, that \succeq is a reliable intelligence order (what we believe to be true), would prove that AIXI is universally optimal. We could further ask: how useful is \succeq for ordering policies of practical interest with intermediate intelligence, or how can \succeq help to guide toward constructing more intelligent systems with reasonable computation time. An effective intelligence order relation \succeq^c will be defined in Section 7.2, which is more useful from a practical point of view.

5.2 On the Optimality of AIXI

In this section we outline ways towards an optimality proof of AIXI, which will be followed thereafter, but not always to the end. Sources of inspiration are the SP loss bounds proven in Chapter 3 and optimality criteria from the adaptive control literature (mainly) for linear systems [KV86]. The value bounds for AIXI are expected to be, in a sense, weaker than the SP loss bounds because the problem class covered by AIXI is much larger than the class of induction problems. Convergence of ξ to μ has already been proven, but is not sufficient to establish convergence of the behavior of the AIXI model to the behavior of the $\text{AI}\mu$ model. We will focus on three approaches towards a general optimality proof:

1) What is meant by universal optimality. The first step is to investigate what we can expect from AIXI, i.e. what is meant by *universal optimality*. A “learner” (like AIXI) may converge to the optimal informed decision maker (like $\text{AI}\mu$) in several senses. Possibly relevant concepts from statistics are consistency, self-tuningness, self-optimizingness, efficiency, unbiasedness, asymptotic or finite convergence [KV86], Pareto-optimality, and some more defined in Section 5.3. Some concepts are stronger than necessary, others are weaker than desirable but suitable to start with. Self-optimizingness is defined as the asymptotic convergence of the average true value $\frac{1}{m} V_{1m}^{p^\xi \mu}$ of AIXI to the optimal value $\frac{1}{m} V_{1m}^{*\mu}$. Apart from convergence-speed, self-optimizingness of AIXI would most closely correspond to the loss bounds proven for SP. We investigate which properties are desirable and under which circumstances the AIXI model satisfies these properties. We will show that no universal model, including AIXI, can in general be self-optimizing. On the other

hand we show that AIXI is Pareto-optimal in the sense that there is no other policy which performs better or equal in all environments, and strictly better in at least one.

2) Limited environmental classes. The problem of defining and proving general value bounds becomes more feasible by considering, in a first step, restricted concept classes. We analyze AIXI for known classes (Markovian, factorizable, predictive, game, optimization, and supervised learning environments) in Chapter 6 and for the new classes (forgetful, relevant, asymptotically learnable, farsighted, uniform and (pseudo-)passive) defined in Section 5.3.

3) Generalization of AIXI to general Bayes-mixtures. Another approach is to generalize AIXI to $\text{AI}\zeta$, where $\zeta() = \sum_{\nu \in \mathcal{M}} w_{\nu} \nu()$ is a general Bayes-mixture of distributions ν in some class \mathcal{M} . If \mathcal{M} is the multi-set of enumerable semi-measures enumerated by a Turing-machine, then $\text{AI}\zeta$ coincides with AIXI. If \mathcal{M} is the (multi)set of passive effective environments, then AIXI reduces to the Λ_{ξ} predictor which has been shown to perform well. We show that these loss/value bounds generalize to wider classes, at least asymptotically. Promising classes, are again, the ones described in Section 5.3. Especially for ergodic MDPs we show $\text{AI}\zeta$ is self-optimizing. Obviously, the least we must demand from \mathcal{M} to have a chance of finding a self-optimizing policy is that there exists some self-optimizing policy at all. This necessary condition will also be sufficient. More generally, the key is not to prove absolute results for specific problem classes, but to prove relative results of the form “if there exists a policy with certain desirable properties, then $\text{AI}\zeta$ also possesses these desirable properties”. If there are tasks which cannot be solved by any policy, $\text{AI}\zeta$ cannot be blamed for failing. Note that in this second approach we have for each environmental class a corresponding model $\text{AI}\zeta$, whereas in approach (2) the same AIXI model is analyzed for all environmental classes.

4) Optimality by construction. A possible 3rd approach towards an optimality “proof” is to regard AIXI as *optimal by construction*. This perspective is common in various (simpler) settings. For instance, in Bandit problems, where pulling arm i leads to reward 1 (0) with unknown probability p_i ($1-p_i$), the traditional Bayesian solution to the uncertainty about p_i is to assume a uniform (or Dirichlet) prior over p_i and to maximize the (subjectively) expected reward sum over multiple trials. The exact solution (in terms of Gittins indices) is widely regarded as “optimal”, although justified alternative approaches exist. Similarly, but simpler, assuming a uniform subjective prior over the Bernoulli parameter $p_{(i)} \in [0,1]$ one arrives at the (reasonable, but more controversial) Laplace rule for predicting i.i.d. sequences. AIXI is similar in the sense that the unknown $\mu \in \mathcal{M}$ is the analogue of the unknown $p \in [0,1]$ and the prior beliefs $w_{\nu} = 2^{-K(\nu)}$ justified by Occam’s razor are the analogue of a uniform distribution over $[0,1]$. In the same sense as Gittins’ solution to the Bandit problem and Laplace rule for Bernoulli sequences, AIXI may be too, regarded as optimal by construction. Theorems relating AIXI to $\text{AI}\mu$ would not be regarded

as optimality proofs of AIXI but just as how much harder it becomes to operate when μ is unknown, i.e. the achievements of approaches 1–3 are simply reinterpreted.

5.3 Value Bounds and Separability Concepts

5.3.1 Introduction

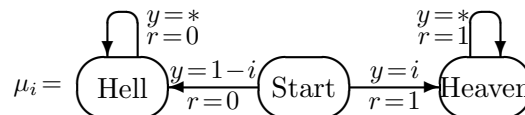
The values V_{km} associated with the AI systems correspond roughly to the negative error measure $-E_n^\rho$ of the SP systems. In SP, we were interested in small bounds for the error excess $E_n^{\Theta_\xi} - E_n^\rho$. Unfortunately, simple value bounds for AIXI in terms of V_{km} analogous to the loss bound in Theorem 3.48 do not hold. We even have difficulties in specifying what we can expect to hold for AIXI or any AI system which claims to be universally optimal. Consequently, we cannot have a proof if we don't know what to prove. In SP, the only important property of μ for proving error bounds was its complexity $K(\mu)$. We will see that in the AI case, there are no useful bounds in terms of $K(\mu)$ only. We either have to study restricted problem classes or consider bounds depending on other properties of μ , rather than on its complexity only. In the following, we will exhibit the difficulties by two examples and introduce concepts which may be useful for proving value bounds. Despite the difficulties in even claiming useful value bounds, we nevertheless, firmly believe that the order relation (5.14) correctly formalizes the intuitive meaning of intelligence and, hence, that the AIXI agent is universally optimal.

5.3.2 (Pseudo) Passive μ and the HeavenHell Example

In the following, we choose $m_k = m$. We want to compare the true, i.e. μ -expected value V_{1m}^μ of a μ independent universal policy p^{best} with any other policy p . Naively, we might expect the existence of a policy p^{best} which maximizes V_{1m}^μ , apart from additive corrections of lower order for $m \rightarrow \infty$

$$V_{1m}^{p^{best}\mu} \geq V_{1m}^{p\mu} - o(\dots) \quad \forall \mu, p \quad (5.15)$$

Such policies are sometimes called self-optimizing [KV86]. Note, that $V_{1m}^{*\mu} \geq V_{1m}^{p\mu} \forall p$, but p^μ is not a candidate for (a universal) p^{best} as it depends on μ . On the other hand, the policy p^ξ of the AIXI agent maximizes V_{1m}^ξ by definition ($p^\xi \succeq p$). As V_{1m}^ξ is thought to be a guess of V_{1m}^μ , we might expect $p^{best} = p^\xi$ to approximately maximize V_{1m}^μ , i.e. (5.15) to hold. Let us consider the problem class (set of environments) $\{\mu_0, \mu_1\}$ with $\mathcal{Y} = \mathcal{R} = \{0, 1\}$ and $r_k = \delta_{iy_1}$ in environment μ_i . The first action y_1 decides whether you go to heaven with all future rewards r_k being 1 (good) or to hell with all future rewards being 0 (bad). Note that μ_i are (deterministic, non-ergodic) MDPs:



It is clear, that if μ_i , i.e. i is known, the optimal policy p^{μ_i} is to output $y_1 = i$ in the first cycle with $V_{1m}^{p^{\mu_i}\mu} = m$. On the other hand, any unbiased policy p^{best} independent of the actual μ either outputs $y_1 = 1$ or $y_1 = 0$. Independent of the actual choice y_1 , there is always an environment ($\mu = \mu_{1-y_1}$) for which this choice is catastrophic ($V_{1m}^{p^{best}\mu} = 0$). No single agent can perform well in both environments μ_0 and μ_1 . The r.h.s. of (5.15) equals $m - o(m)$ for $p = p^\mu$. For all p^{best} there is a μ for which the l.h.s. is zero (The situation remains the same if we add a third action allowing to stay in the start state. See also Problem 5.2 for a similar two-state environment). We have shown that no p^{best} can satisfy (5.15) for all μ and p , so we cannot expect p^ξ to do so. Nevertheless, there are problem classes for which (5.15) holds, for instance SP and CF and ergodic MDPs. For SP, (5.15) is just a reformulation of Theorem 3.48 with an appropriate choice for p^{best} (which differs from p^ξ , see next section). We expect (5.15) to hold for all inductive problems in which the environment is not influenced⁵ by the output of the agent. We want to call these μ , *passive* or *inductive* environments. Further, we want to call \mathcal{M} and $\mu \in \mathcal{M}$ satisfying (5.15) with $p^{best} = p^\xi$ *pseudo passive*. So we expect inductive μ to be pseudo passive.

5.3.3 The OnlyOne Example

Let us give a further example to demonstrate the difficulties in establishing value bounds. Let $\mathcal{R} = \{0, 1\}$ and $|\mathcal{Y}|$ be large. We consider all (deterministic) environments in which a single complex output y^* is correct ($r = 1$) and all others are wrong ($r = 0$). The problem class \mathcal{M} is defined by

$$\mathcal{M} := \{\mu : \mu(y_{<k} y_k \underline{1}) = \delta_{y_k y^*} \forall k, y^* \in \mathcal{Y}, K(y^*) = \lfloor \log_2 |\mathcal{Y}| \rfloor\}$$

There are $N \stackrel{\times}{=} |\mathcal{Y}|$ such y^* . The only way a μ -independent policy p can find the correct y^* , is by trying one y after the other in a certain order. In the first $N - 1$ cycles at most, $N - 1$ different y are tested. As there are N different possible y^* , there is always a $\mu \in \mathcal{M}$ for which p gives erroneous outputs in the first $N - 1$ cycles. The number of errors are $E_\infty^p \geq N - 1 \stackrel{\times}{=} |\mathcal{Y}| \stackrel{\times}{=} 2^{K(y^*)} \stackrel{\times}{=} 2^{K(\mu)}$ for this μ . As this is true for any p , it is also true for the AIXI model, hence $E_k^{p^\xi} \leq 2^{K(\mu)}$ is the best possible error bound we can expect, which depends on $K(\mu)$ only. Actually, we will derive such a bound in Section 6.2 for SP. Unfortunately, as we are mainly interested in the cycle region $k \ll |\mathcal{Y}| \stackrel{\times}{=} 2^{K(\mu)}$ (see Section 4.3.2) this bound is trivial. There are no interesting bounds for deterministic μ depending on $K(\mu)$ only, unlike the SP case. Bounds must either depend on additional properties of μ or we have to consider specialized bounds for restricted problem classes. The case of probabilistic μ is similar. Whereas for SP there are useful bounds in terms of $E_k^{\Theta_\mu}$ and $K(\mu)$, there are no such bounds for AIXI. Again, this is not a drawback of AIXI since for

⁵Of course, the reward feedback r_k depends on the agent's output. What we have in mind is, like in sequence prediction, that the true sequence is not influenced by the agent.

no unbiased AI system the errors/rewards could be bound in terms of $K(\mu)$ and the errors/rewards of $\text{AI}\mu$ only.

There is a way to make use of gross (e.g. $2^{K(\mu)}$) bounds. Assume that after a reasonable number of cycles k , the information $\dot{x}_{<k}$ perceived by the AIXI agent contains a lot of information about the true environment μ . The information in $\dot{x}_{<k}$ might be coded in any form. Let us assume that the complexity $K(\mu|\dot{x}_{<k})$ of μ under the condition that $\dot{x}_{<k}$ is known, is of order 1. Consider a theorem, bounding the sum of rewards or of other quantities over cycles $1\dots\infty$ in terms of $f(K(\mu))$ for a function f with $f(O(1))=O(1)$, like $f(n)=2^n$. Then, there will be a bound for cycles $k\dots\infty$ in terms of $\approx f(K(\mu|\dot{x}_{<k}))=O(1)$. Hence, a bound like $2^{K(\mu)}$ can be replaced by small bound $\approx 2^{K(\mu|\dot{x}_{<k})}=O(1)$ after k cycles. All one has to show/ensure/assume is that enough information about μ is presented (in any form) in the first k cycles. In this way, even a gross bound could become useful. In Section 6.5 we use a similar argument to prove that AIXI is able to learn supervised (cf. Problems 3.13 and 6.3).

5.3.4 Asymptotic Learnability

In the following, we weaken (5.15) in the hope of getting a bound applicable to wider problem classes than the passive one. Consider the I/O sequence $\dot{y}_1\dot{x}_1\dots\dot{y}_n\dot{x}_n$ caused by AIXI. On history $\dot{y}\dot{x}_{<k}$, AIXI will output $\dot{y}_k \equiv \dot{y}_k^\xi$ in cycle k . Let us compare this to \dot{y}_k^μ what $\text{AI}\mu$ would output, still on the same history $\dot{y}\dot{x}_{<k}$ produced by AIXI. As $\text{AI}\mu$ maximizes the μ -expected value, AIXI causes lower (or at best equal) $V_{km_k}^\mu$, if \dot{y}_k^ξ differs from \dot{y}_k^μ . Let $D_{n\mu\xi} := \langle \sum_{k=1}^n 1 - \delta_{\dot{y}_k^\mu, \dot{y}_k^\xi} \rangle_\mu$ be the μ -expected number of suboptimal choices of AIXI, i.e. outputs different from $\text{AI}\mu$ in the first n cycles. One might weigh the deviating cases by their severity. Especially when the μ -expected rewards $V_{km_k}^{p\mu}$ for \dot{y}_k^ξ and \dot{y}_k^μ are equal or close to each other, this should be taken into account in a definition of $D_{n\mu\xi}$, e.g. by a weight factor $[V_{km_k}^{*\mu}(\dot{y}\dot{x}_{<k}) - V_{km_k}^{p\xi\mu}(\dot{y}\dot{x}_{<k})]$. These details do not matter in the following qualitative discussion. The important difference to (5.15) is that here we stick on the history produced by AIXI and count a wrong decision as, at most, one error. The wrong decision in the HeavenHell example in the first cycle no longer counts as losing m rewards, but counts as one wrong decision. In a sense, this is fairer. One shouldn't blame somebody too much who makes a single wrong decision for which he just has too little information available, in order to make a correct decision. The AIXI model would deserve to be called asymptotically optimal, if the probability of making a wrong decision tends to zero, i.e. if

$$D_{n\mu\xi}/n \rightarrow 0 \quad \text{for } n \rightarrow \infty, \quad \text{i.e. } D_{n\mu\xi} = o(n). \quad (5.16)$$

We say that μ can be *asymptotically learned* (by AIXI) if (5.16) is satisfied. We claim that AIXI (for $m_k \rightarrow \infty$) can asymptotically learn every problem μ of relevance, i.e. AIXI is asymptotically optimal. We included the qualifier *of relevance*, as we are not sure whether there could be strange μ spoiling (5.16) but we expect those μ to be irrelevant from the perspective of AI. In the field of Learning, there are many

asymptotic learnability theorems, often not too difficult to prove. So a proof of (5.16) might also be feasible. Unfortunately, asymptotic learnability theorems are often too weak to be useful from a practical point of view. Nevertheless, they point in the right direction.

5.3.5 Uniform μ

From the convergence (5.11) of $\xi \rightarrow \mu$ we might expect $V_{km_k}^{\xi} \rightarrow V_{km_k}^{\mu}$ and hence, \dot{y}_k^{ξ} defined in (5.3) to converge to \dot{y}_k^{μ} defined in (4.17) with μ probability 1 for $k \rightarrow \infty$. The first problem is, that if the V_{km_k} for the different choices of y_k are nearly equal, then even if $V_{km_k}^{\xi} \approx V_{km_k}^{\mu}$, $\dot{y}_k^{\xi} \neq \dot{y}_k^{\mu}$ is possible due to the non-continuity of $\arg\max_{y_k}$. This can be cured by a weighted $D_{n\mu\xi}$ as described above. More serious is the second problem we explain for $h_k = 1$ and $\mathcal{X} = \mathcal{R} = \{0, 1\}$. For $\dot{y}_k^{\xi} \equiv \arg\max_{y_k} \xi(\dot{y}_{< k} y_k \underline{1})$ to converge to $\dot{y}_k^{\mu} \equiv \arg\max_{y_k} \mu(\dot{y}_{< k} y_k \underline{1})$, it is not sufficient to know that $\xi(\dot{y}_{< k} \dot{y}_k) \rightarrow \mu(\dot{y}_{< k} \dot{y}_k)$ as proven in (5.11). We need convergence not only for the true output \dot{y}_k and reward \dot{r}_k , but also for alternate outputs y_k and reward 1. \dot{y}_k^{ξ} converges to \dot{y}_k^{μ} if ξ converges uniformly to μ , i.e. if in addition to (5.11)

$$|\mu(y_{< k} y'_k \underline{x}'_k) - \xi(y_{< k} y'_k \underline{x}'_k)| < c \cdot |\mu(y_{< k} y_k \underline{x}_k) - \xi(y_{< k} y_k \underline{x}_k)| \quad \forall y'_k, x'_k \quad (5.17)$$

holds for some constant c (at least in a μ -expected sense). We call μ satisfying (5.17) *uniform*. For uniform μ one can show (5.16) with appropriately weighted $D_{n\mu\xi}$ and bounded horizon $h_k \leq h_{max}$. Unfortunately there are relevant μ which are not uniform.

5.3.6 Other Concepts

In the following, we briefly mention some further concepts. A *Markovian* μ is defined as depending only on the last cycle, i.e. $\mu(y_{< k} y_k \underline{x}_k) = \mu_k(x_{k-1} y_k)$. We say μ is *generalized* (l^{th} order) *Markovian*, if $\mu(y_{< k} y_k \underline{x}_k) = \mu_k(x_{k-l} y_{k-l+1:k-1} y_k)$ for fixed l . This property has some similarities to *factorizable* μ defined in (4.26). If further $\mu_k \equiv \mu_1 \forall k$, μ is called *stationary*. *Ergodic* Markov decision processes are defined in Section 5.6. Further, we call μ (ξ) *forgetful* if $\mu(y_{< k} y_k \underline{x}_k)$ ($\xi(y_{< k} y_k \underline{x}_k)$) become(s) independent of $y_{< l}$ for fixed l and $k \rightarrow \infty$ with μ probability 1 (cf. Problems 2.5 and 5.13). Further, we say μ is *farsighted*, if $\lim_{m_k \rightarrow \infty} \dot{y}_k^{(m_k)}$ exists. More details will be given in Section 5.7, where we also give an example of a farsighted μ for which nevertheless the limit $m_k \rightarrow \infty$ makes no sense.

5.3.7 Summary

We have introduced several concepts, which might be useful for proving value bounds, including forgetful, relevant, asymptotically learnable, farsighted, uniform, (generalized) Markovian, factorizable and (pseudo) passive μ . We have sorted them here, approximately in the order of decreasing generality. We want to call them

separability concepts. The more general (like relevant, asymptotically learnable and farsighted) μ will be called weakly separable, the more restrictive (like (pseudo) passive and factorizable) μ will be called strongly separable, but we will use these qualifiers in a more qualitative, rather than rigid sense. Other (non-separability) concepts are deterministic μ and, of course, the class of all chronological μ .

5.4 Value Related Optimality Results

5.4.1 The AI_ρ Models: Preliminaries

In Chapter 4 we gave verbal definitions of the AI_μ model (Definition 4.4) and of value functions (Definition 4.5) and derived different mathematical expressions for them. The AIXI model involved similar expressions with μ replaced by $\xi = \xi_U$. The following definitions summarize all formulas for general environment ρ , and general environmental classes \mathcal{M} and general weights w_ν in iterative form.

Definition 5.18 (ρ -Value function) We define the *value* of policy p in environment ρ given history $y_{<k}$, or shorter, the ρ -value of p given $y_{<k}$, as

$$V_{km}^{pp}(y_{<k}) := \sum_{x_{k:m}} (r_k + \dots + r_m) \rho(y_{<k} \underline{y}_{k:m})|_{y_{1:m}=p(x_{<m})}$$

Definition 5.19 (Functional AI_ρ model) The AI_ρ model is defined as the policy p^ρ which maximizes the (total) value $V_\rho^p := V_{1m}^{pp}(\epsilon)$:

$$p^\rho := \arg \max_p V_\rho^p, \quad V_{km}^{*\rho}(y_{<k}) := V_{km}^{p^\rho \rho}(y_{<k}).$$

Theorem 5.20 (Iterative AI_ρ model) The ρ -optimal policy p^ρ and its value $V_{km}^{*\rho}(y_{<k})$ can be explicitly written as

$$\begin{aligned} y_k &= \arg \max_{y_k} \sum_{x_k} \max_{y_{k+1}} \sum_{x_{k+1}} \dots \max_{y_m} \sum_{x_m} (r_k + \dots + r_m) \cdot \rho(y_{<k} \underline{y}_{k:m}), \\ V_{km}^{*\rho}(y_{<k}) &= \max_{y_k} \sum_{x_k} \max_{y_{k+1}} \sum_{x_{k+1}} \dots \max_{y_m} \sum_{x_m} (r_k + \dots + r_m) \cdot \rho(y_{<k} \underline{y}_{k:m}). \end{aligned}$$

Furthermore $V_{km}^{*\rho}(y_{<k}) \geq V_{km}^{pp}(y_{<k}) \forall p$.

The proof is the same as in Chapter 4 with μ replaced by ρ . The following property of V_ρ is crucial.

Theorem 5.21 (Linearity and convexity of V_ρ in ρ) V_ρ^p is a linear function in ρ and V_ρ^* is a convex function in ρ in the sense that

$$V_\xi^p = \sum_{\nu \in \mathcal{M}} w_\nu V_\nu^p \quad \text{and} \quad V_\xi^* \leq \sum_{\nu \in \mathcal{M}} w_\nu V_\nu^* \quad \text{where} \quad \xi(\underline{y}_{1:m}) = \sum_{\nu \in \mathcal{M}} w_\nu \nu(\underline{y}_{1:m})$$

Proof. Linearity is obvious from the Definition 5.18 of V_ρ^p . Convexity follows from $V_\xi^* \equiv V_\xi^{p^\xi} = \sum_{\nu} w_\nu V_\nu^{p^\xi} \leq \sum_{\nu} w_\nu V_\nu^*$, where the first equality is just the definition (5.19), the second equality uses linearity of $V_\rho^{p^\xi}$ just proven, and the last inequality follows from the dominance (5.20) and positivity of the weights w_ν . \square

One loose interpretation of the convexity is that a mixture can never increase performance.

5.4.2 Pareto Optimality of AI ξ

This subsection shows Pareto-optimality of AI ξ analogous to SP. The total μ -expected reward $V_\mu^{p^\xi}$ of policy p^ξ of the AI ξ model is of central interest in judging the performance of AI ξ . We know that there *are* policies (e.g. p^μ of AI μ) with higher μ -value ($V_\mu^* \geq V_\mu^{p^\xi}$). In general, every policy based on an estimate ρ of μ which is closer to μ than ξ is, outperforms p^ξ in environment μ , simply because it is more tailored toward μ . On the other hand, such a system probably performs worse than p^ξ in other environments. Since we do not know μ in advance we may ask whether there exists a policy p with better or equal performance than p^ξ in *all* environments $\nu \in \mathcal{M}$ and a strictly better performance for one $\nu \in \mathcal{M}$. This would clearly render p^ξ suboptimal. We show that there is no such p .

Definition 5.22 (Pareto Optimality) A policy \tilde{p} is called Pareto-optimal if there is no other policy p with $V_\nu^p \geq V_\nu^{\tilde{p}}$ for all $\nu \in \mathcal{M}$ and strict inequality for at least one ν .

Theorem 5.23 (Pareto Optimality) AI ξ alias p^ξ is Pareto-optimal.

Proof. We want to arrive at a contradiction by assuming that p^ξ is not Pareto-optimal, i.e. by assuming the existence of a policy p with $V_\nu^p \geq V_\nu^{p^\xi}$ for all $\nu \in \mathcal{M}$ and strict inequality for at least one ν :

$$V_\xi^p = \sum_{\nu} w_\nu V_\nu^p > \sum_{\nu} w_\nu V_\nu^{p^\xi} = V_\xi^{p^\xi} \equiv V_\xi^* \geq V_\xi^p$$

The two equalities follow from linearity of V_ρ (5.21). The strict inequality follows from the assumption and $w_\nu > 0$. The last inequality follows from the fact that

p^ξ maximizes by definition the universal value (5.20). The contradiction $V_\xi^p > V_\xi^{p^\xi}$ proves Pareto-optimality of $\text{AI}\xi$. \square

Pareto-optimality should be regarded as a necessary condition for an agent aiming to be optimal. From a practical point of view a significant increase of V for many environments ν may be desirable even if this causes a small decrease of V for a few other ν . The impossibility of such a “balanced” improvement is a more demanding condition on p^ξ than pure Pareto-optimality. The next theorem shows that p^ξ is also balanced-Pareto-optimal in the following sense:

Theorem 5.24 (Balanced Pareto Optimality)

$$\Delta_\nu := V_\nu^{p^\xi} - V_\nu^{\tilde{p}}, \quad \Delta := \sum_{\nu \in \mathcal{M}} w_\nu \Delta_\nu \quad \Rightarrow \quad \Delta \geq 0.$$

This implies the following: Assume \tilde{p} has lower value than p^ξ on environments \mathcal{L} by a total weighted amount of $\Delta_{\mathcal{L}} := \sum_{\lambda \in \mathcal{L}} w_\lambda \Delta_\lambda$. Then \tilde{p} can have higher value on $\eta \in \mathcal{H} := \mathcal{M} \setminus \mathcal{L}$, but the improvement is bounded by $\Delta_{\mathcal{H}} := |\sum_{\eta \in \mathcal{H}} w_\eta \Delta_\eta| \leq \Delta_{\mathcal{L}}$. Especially $|\Delta_\eta| \leq w_\eta^{-1} \max_{\lambda \in \mathcal{L}} \Delta_\lambda$.

This means that a weighted value increase $\Delta_{\mathcal{H}}$ by using \tilde{p} instead of p^ξ is compensated by an at least as large weighted decrease $\Delta_{\mathcal{L}}$ on other environments. If the decrease is small, the increase can also only be small. In the special case of only a single environment with decreased value Δ_λ , the increase is bound by $\Delta_\eta \leq \frac{w_\lambda}{w_\eta} |\Delta_\lambda|$, i.e. a decrease by an amount Δ_λ can only cause an increase by at most the same amount times a factor $\frac{w_\lambda}{w_\eta}$. For the choice of the weights $w_\nu \sim 2^{-K(\nu)}$, a decrease can only cause a smaller increase in simpler environments, but a scaled increase in more complex environments. Finally note that pure Pareto-optimality (5.23) follows from balanced Pareto-optimality in the special case of no decrease $\Delta_{\mathcal{L}} \equiv 0$.

Proof. $\Delta \geq 0$ follows from $\Delta = \sum_\nu w_\nu [V_\nu^{p^\xi} - V_\nu^{\tilde{p}}] = V_\xi^{p^\xi} - V_\xi^{\tilde{p}} \geq 0$, where we have used linearity of V_ρ (5.21) and dominance $V_\xi^{p^\xi} \geq V_\xi^p$ (5.20). The remainder of Theorem 5.24 is obvious from $0 \leq \Delta = \Delta_{\mathcal{L}} - \Delta_{\mathcal{H}}$ and by bounding the weighted average Δ_η by its maximum. \square

5.4.3 Self-optimizing Policy p^ξ w.r.t. Average Value

We have argued in Section 5.3 Eq. (5.15) that there is no (universal) policy \tilde{p} (independent of the actual environment μ) for which⁶

$$\frac{1}{m} V_{1m}^{\tilde{p}\mu} \rightarrow \frac{1}{m} V_{1m}^{*\mu} \quad \text{for } m \rightarrow \infty \quad (5.25)$$

⁶Here and elsewhere we interpret $a_m \rightarrow b_m$ as an abbreviation for $a_m - b_m \rightarrow 0$. $\lim_{m \rightarrow \infty} b_m$ may not exist.

and hence also p^ξ cannot converge to the optimal policy p^μ in this sense. On the other hand, we know from Section 2.4 that convergence of this type (and even stronger) holds for SP. Section 5.3 suggested to investigate restricted environmental class. In the following we consider the generalized AI ξ model (with $\xi() = \sum_{\nu \in \mathcal{M}} w_\nu \nu()$) and restricted classes \mathcal{M} . The least we must demand from \mathcal{M} to have a chance that

$$\frac{1}{m} V_{1m}^{p^\xi \mu} \rightarrow \frac{1}{m} V_{1m}^{*\mu} \quad \text{for } m \rightarrow \infty \quad (5.26)$$

is that there exists some policy \tilde{p} at all with this property (5.25). Luckily, this necessary condition will also be sufficient. This is another (asymptotic) optimality property of (generalized) AI ξ . If universal convergence in the sense of (5.25) is possible at all in a class of environments \mathcal{M} , then AI ξ converges in the sense of (5.26). We will call policies \tilde{p} with a property like (5.26) *self-optimizing* [KV86].

The following two Lemmas pave the way for proving the convergence Theorem.

Lemma 5.27 (Value difference relation)

$$0 \leq V_\nu^* - V_\nu^{\tilde{p}} =: \Delta_\nu \quad \Rightarrow \quad 0 \leq V_\nu^* - V_\nu^{p^\xi} \leq \frac{1}{w_\nu} \Delta \quad \text{with} \quad \Delta := \sum_{\nu \in \mathcal{M}} w_\nu \Delta_\nu$$

Proof. The following sequence of inequalities proves the lemma:

$$0 \leq w_\nu [V_\nu^* - V_\nu^{p^\xi}] \leq \sum_\nu w_\nu [V_\nu^* - V_\nu^{p^\xi}] \leq \sum_\nu w_\nu [V_\nu^* - V_\nu^{\tilde{p}}] = \sum_\nu w_\nu \Delta_\nu \equiv \Delta$$

In the first and second inequality we used $w_\nu \geq 0$ and $V_\nu^* - V_\nu^{p^\xi} \geq 0$. The last inequality follows from $\sum_\nu w_\nu V_\nu^{p^\xi} = V_\xi^{p^\xi} \equiv V_\xi^* \geq V_\xi^{\tilde{p}} = \sum_\nu w_\nu V_\nu^{\tilde{p}}$. \square

We also need some results for averages of functions $\delta_\nu(m)$ converging to zero.

Lemma 5.28 (Convergence of averages) For $\delta(m) := \sum_{\nu \in \mathcal{M}} w_\nu \delta_\nu(m)$ the following holds (we only need $\sum_\nu w_\nu \leq 1$):

- i) $\delta_\nu(m) \leq f(m) \quad \forall \nu \quad \text{implies} \quad \delta(m) \leq f(m).$
- ii) $\delta_\nu(m) \xrightarrow{m \rightarrow \infty} 0 \quad \forall \nu \quad \text{implies} \quad \delta(m) \xrightarrow{m \rightarrow \infty} 0 \quad \text{if } 0 \leq \delta_\nu(m) \leq c.$
- iii) $\delta(m) \leq \max_\nu \delta_\nu(m).$
- iv) $\delta_\nu(m) = O(f(m)) \quad \forall \nu \quad \text{implies} \quad \delta(m) = O(f(m)) \quad \text{if } \mathcal{M} \text{ is finite.}$

Proof. (i) immediately follows from $\delta(m) = \sum_\nu w_\nu \delta_\nu(m) \leq \sum_\nu w_\nu f(m) \leq f(m).$

(ii) We choose some order on \mathcal{M} and some $\nu_0 \in \mathcal{M}$ large enough such that $\sum_{\nu \geq \nu_0} w_\nu \leq \frac{\varepsilon}{c}$. Using $\delta_\nu(m) \leq c$ this implies

$$\sum_{\nu \geq \nu_0} w_\nu \delta_\nu(m) \leq \sum_{\nu \geq \nu_0} w_\nu c \leq \varepsilon.$$

Furthermore, the assumption $\delta_\nu(m) \rightarrow 0$ means that there is an $m_{\nu\varepsilon}$ depending on ν and ε such that $\delta_\nu(m) \leq \varepsilon$ for all $m \geq m_{\nu\varepsilon}$. This implies

$$\sum_{\nu \leq \nu_0} w_\nu \delta_\nu(m) \leq \sum_{\nu \leq \nu_0} w_\nu \varepsilon \leq \varepsilon \quad \text{for all } m \geq \max_{\nu \leq \nu_0} \{m_{\nu\varepsilon}\} =: m_\varepsilon.$$

$m_\varepsilon < \infty$, since the maximum is over a finite set. Together we have

$$\delta(m) \equiv \sum_{\nu \in \mathcal{M}} w_\nu \delta_\nu(m) \leq 2\varepsilon \quad \text{for } m \geq m_\varepsilon \quad \Rightarrow \quad \delta(m) \rightarrow 0 \quad \text{for } m \rightarrow \infty$$

since ε was arbitrary and $\delta(m) \geq 0$.

$$(iii) \quad \delta(m) \equiv \sum_{\nu} w_\nu \delta_\nu(m) \leq \sum_{\nu} w_\nu \max_{\nu} \delta_\nu(m) \leq \max_{\nu} \delta_\nu(m).$$

(iv) From $\delta_\nu(m) \leq c_\nu f(m)$ it follows that

$$\delta(m) \leq \sum_{\nu} w_\nu c_\nu f(m) \leq (\max_{\nu} c_\nu) f(m) (\sum_{\nu} w_\nu) \leq c_{\max} f(m)$$

with $c_{\max} := \max_{\nu \in \mathcal{M}} c_\nu$ being finite, since \mathcal{M} is finite by assumption. \square

The boundedness assumption of (ii) without the finiteness assumption (iv) is not sufficient to prove (iv). For instance, for $\mathcal{M} \cong \mathbb{N}$, $\delta_\nu(m) := e^{-m/\nu} \leq 1$ decays exponentially in $m \geq 1$ for every ν , but for $w_\nu = \frac{1}{\nu(\nu+1)}$, $\delta(m) \geq \frac{1}{2m}$ decays only harmonically (Problem 5.6).

Theorem 5.29 (Self-optimizing policy p^ξ w.r.t. average value) If there exists a sequence of policies \tilde{p}_m , $m=1,2,3,\dots$ with value within $\Delta(m)$ to optimum for all environments $\nu \in \mathcal{M}$, then, save for a constant factor this also holds for the sequence of universal policies p_m^ξ , i.e.

$$i) \quad \text{If } \exists \tilde{p}_m \forall \nu : V_{1m}^{*\nu} - V_{1m}^{\tilde{p}_m \nu} \leq \Delta(m) \quad \Rightarrow \quad V_{1m}^{*\mu} - V_{1m}^{p_m^\xi \mu} \leq \frac{1}{w_\mu} \Delta(m)$$

If there exists a sequence of self-optimizing policies \tilde{p}_m in the sense that their expected average reward $\frac{1}{m} V_{1m}^{\tilde{p}_m \nu}$ converges to the optimal average $\frac{1}{m} V_{1m}^{*\nu}$ for all environments $\nu \in \mathcal{M}$, then this also holds for the sequence of universal policies p_m^ξ , i.e.

$$ii) \quad \text{If } \exists \tilde{p}_m \forall \nu : \frac{1}{m} V_{1m}^{\tilde{p}_m \nu} \xrightarrow{m \rightarrow \infty} \frac{1}{m} V_{1m}^{*\nu} \quad \Rightarrow \quad \frac{1}{m} V_{1m}^{p_m^\xi \mu} \xrightarrow{m \rightarrow \infty} \frac{1}{m} V_{1m}^{*\mu}.$$

The beauty of this theorem is that if universal convergence in the sense of (5.25) is possible at all in a class of environments \mathcal{M} , then $\text{AI}\xi$ converges (in the sense of (5.26)). The necessary condition of convergence is also sufficient. The unattractive point is that this is not an asymptotic convergence statement for $V_{km}^{p^\xi \mu}$ of a single policy p^ξ for $k \rightarrow \infty$ for some fixed m , and in fact no such theorem could be true, since always $k \leq m$. The theorem merely says that under the stated conditions the

average value of $\text{AI}\xi(m)$ can be arbitrarily close to optimum for sufficiently large (pre-chosen) horizon m . This weakness will be resolved in the next subsection.

Proof. (i) $\Delta_\nu(m) \leq f(m)$ implies $\Delta(m) \leq f(m)$ by Lemma 5.28(i). Inserting this in Lemma 5.27 we get 5.29(i) (recovering the m dependence and finally renaming $f \rightsquigarrow \Delta$).

(ii) We define $\delta_\nu(m) := \frac{1}{m}\Delta_\nu(m) = \frac{1}{m}[V_\nu^* - V_\nu^{\tilde{p}}]$. Since we generally assumed bounded rewards $0 \leq r \leq r_{\max}$ (4.28) we have

$$V_\nu^* \leq mr_{\max} \quad \text{and} \quad V_\nu^{\tilde{p}} \geq 0 \quad \Rightarrow \quad \Delta_\nu \leq mr_{\max} \quad \Rightarrow \quad 0 \leq \delta_\nu(m) \leq c := r_{\max}.$$

The premise in 5.29(ii) is that $\delta_\nu(m) = \frac{1}{m}[V_{1m}^{*\nu} - V_{1m}^{\tilde{p}\nu}] \rightarrow 0$ which implies

$$0 \leq \frac{1}{m}[V_{1m}^{*\nu} - V_{1m}^{p^\xi\nu}] \leq \frac{1}{w_\nu} \frac{\Delta(m)}{m} = \frac{1}{w_\nu} \delta(m) \rightarrow 0.$$

The inequalities follow from Lemma 5.27 and convergence to zero from Lemma 5.28(ii). This proves 5.29(ii). \square

Lemma 5.28(i), and hence Theorem 5.29(i) may be generalized to infinite \mathcal{M} by demanding $\delta_\nu(m) \leq f(m)$. But usually the convergence rate of policies depends on ν , at least in form of a factor or additive constant, so this generalization is probably vacuous. In Section 5.6 we show that a converging \tilde{p} exists for ergodic MDPs, and hence p^ξ converges in this environmental class too (in the sense of Theorem 5.29).

5.5 Discounted Future Value Function

We now shift our focus from the total value V_{1m} , $m \rightarrow \infty$ to the future value (value-to-go) $V_{k?}$, $k \rightarrow \infty$. The reasons are at least twofold.

i) We want to compare the future value of the optimal informed policy p^μ to the universal learner p^ξ . We regard the first k cycles as a grace period in which p^ξ learns and after which it performs well. The HeavenHell example of Section 5.3 shows that one cannot avoid that a learner gets trapped. We do not want to exclude trapping environments in our analysis from the very beginning, since there could be interesting structure and behavior in the traps itself (Hell or Heaven may reward intelligent behavior). One possibility is to compare future values $V_{k?}^{p^\xi\mu}$ with $V_{k?}^{*\mu}$ on the same (fictitious) history $y_{<k}$. This addresses questions like: If p^ξ gets trapped in a (structured) Hell, does it perform as well as p^μ when put in Hell?

ii) We want to get rid of the horizon parameter m . In the last subsection we have shown a convergence theorem for $m \rightarrow \infty$, but a specific policy p^ξ is defined for all times relative to a fixed horizon m . Current time k is moving, but m is fixed⁷. Actually, to use $k \rightarrow \infty$ arguments we *have* to get rid of m , since $k \leq m$. This is the reason for the question mark in $V_{k?}$ above.

⁷The dynamic horizon m_k introduced early was convenient to discuss qualitative properties, but does not lead to a consistent model.

We eliminate the horizon by discounting the rewards $r_k \rightsquigarrow \gamma_k r_k$ with $\sum_{i=1}^{\infty} \gamma_i < \infty$ and letting $m \rightarrow \infty$. The analogue of m is now an effective horizon h_k^{eff} which may be defined by $\sum_{i=k}^{k+h_k^{eff}} \gamma_i \sim \sum_{i=k+h_k^{eff}}^{\infty} \gamma_i$ (see Section 5.7 for a detailed discussion of the horizon problem). Furthermore, we renormalize $V_{k\infty}$ by $\sum_{i=k}^{\infty} \gamma_i$ and denote it by $V_{k\gamma}$. It can be interpreted as a future expected weighted-average reward. Furthermore we extend the definition to probabilistic policies π (see Problem 4.2).

Definition 5.30 (Discounted AI ρ model and value) We define the γ discounted weighted-average future *value* of (probabilistic) policy π in environment ρ given history $y_{<k}$, or shorter, the ρ -value of π given $y_{<k}$, as

$$V_{k\gamma}^{\pi\rho}(y_{<k}) := \frac{1}{\Gamma_k} \lim_{m \rightarrow \infty} \sum_{y_{k:m}} (\gamma_k r_k + \dots + \gamma_m r_m) \rho(y_{<k} y_{k:m}) \pi(y_{<k} y_{k:m})$$

with $\Gamma_k := \sum_{i=k}^{\infty} \gamma_i$. The discounted AI ρ model is defined as the policy p^ρ which maximizes the future value $V_{k\gamma}^{\pi\rho}$:

$$p^\rho := \arg \max_{\pi} V_{k\gamma}^{\pi\rho}, \quad V_{k\gamma}^{*\rho} := V_{k\gamma}^{p^\rho\rho} = \max_{\pi} V_{k\gamma}^{\pi\rho} \geq V_{k\gamma}^{\pi\rho} \forall \pi.$$

Remarks.

- $\pi(y_{<k} y_{k:m})$ is actually independent of x_m .
- Normalization of $V_{k\gamma}$ by Γ_k does not affect the policy p^ρ .
- The definition of p^ρ is independent of k (in the sense of Problem 5.7).
- Without normalization by Γ_k the future values would converge to zero $k \rightarrow \infty$ in every environment for every policy.
- For an MDP environment, a stationary policy, and geometric discounting, the future value is independent of k and reduces to the well-known MDP value function.
- There is always a deterministic optimizing policy p^ρ (which we use).
- For a deterministic policy there is exactly one $y_{k:m}$ for each $x_{k:m}$ with $\pi \neq 0$. The sum over $y_{k:m}$ drops in this case.
- An iterative representation as in Theorem 5.20 is possible.
- Setting $\gamma_k = 1$ for $k \leq m$ and $\gamma_k = 0$ for $k > m$ gives back the undiscounted AI ρ model (5.19) with $V_{1\gamma}^{p\rho} = \frac{1}{m} V_{1m}^{p\rho}$.
- $V_{k\gamma}$ and w_k^ν (see below) depend on the realized history $y_{<k}$.

Similarly to the previous subsections one can prove the following properties:

Theorem 5.31 (Linearity and convexity of V_ρ in ρ) $V_{k\gamma}^{\pi\rho}$ is a linear function in ρ and $V_{k\gamma}^{*\rho}$ is a convex function in ρ in the sense that

$$V_{k\gamma}^{\pi\xi} = \sum_{\nu \in \mathcal{M}} w_k^\nu V_{k\gamma}^{\pi\nu} \quad \text{and} \quad V_{k\gamma}^{*\xi} \leq \sum_{\nu \in \mathcal{M}} w_k^\nu V_{k\gamma}^{*\nu}$$

$$\text{where } \xi(y_{<k} \underline{y}_{k:m}) = \sum_{\nu \in \mathcal{M}} w_k^\nu \nu(y_{<k} \underline{y}_{k:m}) \quad \text{with} \quad w_k^\nu := w_\nu \frac{\nu(\underline{y}_{<k})}{\xi(\underline{y}_{<k})}$$

The conditional representation of ξ can easily be proven by dividing the definition of $\xi(\underline{y}_{1:m})$ (5.21) by $\xi(\underline{y}_{<k})$ and by using the chain rule. The posterior weight w_k^ν may be interpreted as the posterior belief in ν and is related to learning aspects of policy p^ξ .

Theorem 5.32 (Pareto Optimality) For every k and history $y_{<k}$ the following holds: p^ξ is Pareto-optimal in the sense that there is no other policy π with $V_{k\gamma}^{\pi\nu} \geq V_{k\gamma}^{p^\xi\nu}$ for all $\nu \in \mathcal{M}$ and strict inequality for at least one ν .

Lemma 5.33 (Value difference relation)

$$0 \leq V_{k\gamma}^{*\nu} - V_{k\gamma}^{\tilde{\pi}_k\nu} =: \Delta_k^\nu \Rightarrow 0 \leq V_{k\gamma}^{*\nu} - V_{k\gamma}^{p^\xi\nu} \leq \frac{1}{w_k^\nu} \Delta_k^\nu$$

$$\text{with } \Delta_k := \sum_{\nu \in \mathcal{M}} w_k^\nu \Delta_k^\nu, \quad \text{where all quantities depend on history } y_{<k}.$$

The proof of Theorem 5.32 and Lemma 5.33 follows the same steps as for Theorem 5.23 and Lemma 5.27 with appropriate replacements. The proof of the analogue of the convergence Theorem 5.29 involves one additional step.

Theorem 5.34 (Self-optimizing policy p^ξ w.r.t. discounted value) For any \mathcal{M} , if there exists a sequence of self-optimizing policies $\tilde{\pi}_k$ $k = 1, 2, 3, \dots$ in the sense that their expected weighted-average reward $V_{k\gamma}^{\tilde{\pi}_k\nu}$ converges for $k \rightarrow \infty$ with μ probability one to the optimal value $V_{k\gamma}^{*\nu}$ for all environments $\nu \in \mathcal{M}$, then this also holds for the universal policy p^ξ in the μ -environment, i.e.

$$\text{If } \exists \tilde{\pi}_k \forall \nu : V_{k\gamma}^{\tilde{\pi}_k\nu} \xrightarrow{k \rightarrow \infty} V_{k\gamma}^{*\nu} \text{ w.}\nu\text{.p.1} \Rightarrow V_{k\gamma}^{p^\xi\mu} \xrightarrow{k \rightarrow \infty} V_{k\gamma}^{*\mu} \text{ w.}\mu\text{.p.1.}$$

The probability qualifier refers to the historic perceptions $x_{<k}$. The historic actions $y_{<k}$ are arbitrary.

The conclusion is valid for action histories $y_{<k}$ if the condition is satisfied for this action history. Since we need the conclusion for the p^ξ -action history, which is hard to characterize, we usually need to prove the condition for *all* action histories. Theorem

5.34 is a powerful result: An inconsistent sequence of probabilistic policies $\tilde{\pi}_k$ suffices to prove the existence of a consistent deterministic policy p^ξ . A result similar to 5.29(i) also holds for the discounted case, roughly saying that $V^{\tilde{\pi}} - V^* = O(\Delta(k))$ implies $V^{p^\xi} - V^* = \frac{1}{\varepsilon} O(\Delta(k))$ with μ probability $1 - \varepsilon$ for finite \mathcal{M} .

Proof. We define $\delta_\nu(k) := \Delta_k^\nu = V_{k\gamma}^{*\nu} - V_{k\gamma}^{\tilde{\pi}\nu}$. Since we generally assumed bounded rewards $0 \leq r \leq r_{max}$ (4.28) and $V_{k\gamma}^{*\nu}$ is a weighted average of rewards we have

$$V_{k\gamma}^{*\nu} \leq r_{max} \quad \text{and} \quad V_{k\gamma}^{\tilde{\pi}\nu} \geq 0 \quad \Rightarrow \quad 0 \leq \delta_\nu(k) = \Delta_k^\nu \leq c := r_{max}.$$

The premise in (5.34) is that $\delta_\nu(k) = V_{k\gamma}^{*\nu} - V_{k\gamma}^{\tilde{\pi}\nu} \rightarrow 0$ for $k \rightarrow \infty$ which implies

$$0 \leq V_{k\gamma}^{*\nu} - V_{k\gamma}^{p^\xi\nu} \leq \frac{1}{w_k^\mu} \Delta_k = \frac{1}{w_k^\mu} \delta(k)$$

The inequalities follow from Lemma 5.33. $\delta(k)$ converges to zero (w. μ .p.1) by Lemma 5.28(ii). What is new and what remains to be shown is that w_k^μ is bounded from below. We show that $z_{k-1} := \frac{w_\mu}{w_k^\mu} = \frac{\xi(\underline{y}_{<k})}{\mu(\underline{y}_{<k})} \geq 0$ converges to a finite value, which completes the proof. Let \mathbf{E} denote the μ expectation. Then

$$\mathbf{E}[z_k | x_{<k}] = \sum'_{x_k} \mu(\underline{y}_{<k} \underline{y}_k) \frac{\xi(\underline{y}_{1:k})}{\mu(\underline{y}_{1:k})} = \frac{\sum'_{x_k} \xi(\underline{y}_{<k} \underline{y}_k) \xi(\underline{y}_{<k})}{\mu(\underline{y}_{<k})} \leq \frac{\xi(\underline{y}_{<k})}{\mu(\underline{y}_{<k})} = z_{k-1}$$

\sum'_{x_k} runs over all x_k with $\mu(\underline{y}_{1:k}) \neq 0$. The first equality holds w. μ .p.1. In the second equality we have used the chain rule twice. $\mathbf{E}[z_k | x_{<k}] \leq z_{k-1}$ shows that $-z_k$ is a semi-martingale. Since $-z_k$ is non-positive, [Doo53, Th.4.1s(i), p324] implies that $-z_k$ converges for $k \rightarrow \infty$ to a finite value w. μ .p.1. (If μ and ξ are lower semi-computable, then boundedness of z_{k-1} follows without the use of Martingales from $z_{k-1} = \frac{\xi(\underline{y}_{<k})}{\mu(\underline{y}_{<k})} \leq \frac{\xi_U(\underline{y}_{<k})}{\mu(\underline{y}_{<k})} \leq c < \infty$, where the first inequality follows from the universality of ξ_U and the second inequality holds for all μ .M.L.-random sequences.) \square

We want to give an intuitive reason for the necessity of the probability qualifier. Assume that the true environment is μ , but choose a history $x_{<k}$ sampled from $\nu \neq \mu$. This implies that ξ converges to ν for $k \rightarrow \infty$. This means that at a fixed but large time k , ξ is very close to ν . It is very hard (takes large h_k^{eff}) to get rid of this wrong bias and to become close to μ later. In the limit this is impossible at all.

The following continuity properties for the (discounted) values hold:

Theorem 5.35 (Continuity of discounted value) The values $V_{k\gamma}^{\pi\mu}$ and $V_{k\gamma}^{*\mu}$ are continuous in μ , and $V_{k\gamma}^{p^{\hat{\mu}}\mu}$ is continuous in $\hat{\mu}$ at $\hat{\mu} = \mu$ w.r.t. a (conditional) maximums norm in the following sense: If $|\mu(\underline{y}_{<k} \underline{y}_k) - \hat{\mu}(\underline{y}_{<k} \underline{y}_k)| \leq \varepsilon \quad \forall \underline{y}_{1:k}$ $\forall k \geq k_0$, then (i) $|V_{k\gamma}^{\pi\mu} - V_{k\gamma}^{\pi\hat{\mu}}| \leq \delta(\varepsilon)$, (ii) $|V_{k\gamma}^{*\mu} - V_{k\gamma}^{*\hat{\mu}}| \leq \delta(\varepsilon)$, and (iii) $|V_{k\gamma}^{*\mu} - V_{k\gamma}^{p^{\hat{\mu}}\mu}| \leq 2\delta(\varepsilon)$ for all $k \geq k_0$ and $\underline{y}_{<k}$, where $\delta(\varepsilon) = r_{max} \cdot \min_{n \geq k} \{|\mathcal{X}|(n-k)\varepsilon + \frac{\Gamma_n}{\Gamma_k}\} \xrightarrow{\varepsilon \rightarrow 0} 0$.

Care has to be taken in the interpretation and use of this theorem: It cannot be used to conclude $V_{k\gamma}^{p^\xi\mu} \rightarrow V_{k\gamma}^{*\mu}$, since $\xi \rightarrow \mu$ does not hold for all $y_{1:\infty}$, but only for μ -random ones (more precisely w. μ .p.1). The condition in Theorem 5.35 cannot be weakened, since p^ξ is not self-optimizing if \mathcal{M} does not admit self-optimizing policies. Furthermore continuity is not uniform in k which prevents using this theorem in the proof of Theorem 5.38. Finally, note that $V_{k\gamma}^{p^\mu\mu}$ can be discontinuous in $\hat{\mu}$ at $\hat{\mu} \neq \mu$. On the positive side, continuity holds for any μ and γ , no structural assumptions have to be made. By setting $\gamma_k = 1$ for $k \leq m$ and $\gamma_k = 0$ for $k > m$ we also get continuity of $V_{km}^{p^\mu\mu}$, $V_{km}^{*\mu}$, and $V_{km}^{p^\mu\mu}$ with $\delta(\varepsilon) \leq r_{max}|\mathcal{X}|\varepsilon(m-k+1)^2$ (set $n = m+1$). For geometric discount $\gamma_k = \gamma^k$ the theorem holds with $\delta(\varepsilon) = \frac{r_{max}|\mathcal{X}|\varepsilon}{1-\gamma}$ (which follows from the second last bound on Δ_k in the proof below for $n = \infty$).

Proof. (i) $V_{k\gamma}$ can be represented recursively like in the undiscounted case as

$$\Gamma_k V_{k\gamma}^{\pi\rho}(y_{<k}) = \sum_{y_k} \pi(y_{<k} y_k) \rho(y_{<k} y_k) [\gamma_k r_k + \Gamma_{k+1} V_{k+1,\gamma}^{\pi\rho}(y_{1:k})]$$

which can easily be verified by induction. The absolute difference of two values can be written as

$$\begin{aligned} \Delta_k &:= \Gamma_k |V_{k\gamma}^{\pi\mu} - V_{k\gamma}^{\pi\hat{\mu}}| = \left| \sum_{yx} \pi \cdot \mu \cdot (\gamma r + \Gamma V) - \sum_{yx} \pi \cdot \hat{\mu} \cdot (\gamma r + \Gamma \hat{V}) \right| \\ &\leq \sum_y \pi \left| \sum_x (\mu - \hat{\mu}) \gamma r + \Gamma \sum_x (\mu V - \hat{\mu} \hat{V}) \right| \\ &= \sum_y \pi \left| \sum_x (\mu - \hat{\mu}) \gamma r + \frac{1}{2} \Gamma \sum_x (\mu - \hat{\mu})(V + \hat{V}) + \frac{1}{2} \Gamma \sum_x (\mu + \hat{\mu})(V - \hat{V}) \right| \leq \dots \end{aligned}$$

where we have suppressed all indices and arguments of all variables and functions. We upper bound the last expression by pulling in the absolute bars. Using (in this order) $0 \leq r \leq r_{max}$, $\sum_x |\mu - \hat{\mu}| \leq \varepsilon |\mathcal{X}|$ (by assumption), $0 \leq V + \hat{V} \leq 2r_{max}$, $|V - \hat{V}| \leq \max_{yx} |V - \hat{V}|$, $\sum_x (\mu + \hat{\mu}) = 2$, $\sum_y \pi = 1$ we get

$$\begin{aligned} \dots &\leq \varepsilon |\mathcal{X}| \gamma r_{max} + \Gamma \varepsilon |\mathcal{X}| r_{max} + \Gamma \max_{yx} |V - \hat{V}| \\ &= \varepsilon |\mathcal{X}| \Gamma_k r_{max} + \max_{y_{k:n-1}} \Delta_{k+1} \leq \dots \leq \\ &\leq \varepsilon |\mathcal{X}| r_{max} \sum_{i=k}^{n-1} \Gamma_i + \max_{y_{k:n-1}} \Delta_n \leq \varepsilon |\mathcal{X}| r_{max} (n-k) \Gamma_k + \Gamma_n r_{max}. \end{aligned}$$

In the second line we used $\gamma_k + \Gamma_{k+1} = \Gamma_k$ and the definition of Δ_k . In the third line we recursively inserted the bound for Δ_i , $i = k+1, \dots, n-1$ we have just derived. In the final expression we used $\Gamma_i \leq \Gamma_k$ for $i \geq k$ and $|V - \hat{V}| \leq r_{max}$. This bound on Δ_k is valid for all $n \geq k$, so we may take the minimum over $n \geq k$. This leads to $\Delta_k \leq \Gamma_k \delta(\varepsilon)$ where $\delta(\varepsilon)$ has been defined in Theorem 5.35. This proves (i).

(ii) For any two real-valued functions f, \hat{f} over some domain \mathcal{D} with $|f(x) - \hat{f}(x)| \leq \delta \forall x \in \mathcal{D}$ we also have $|f_{max} - \hat{f}_{max}| \leq \delta$, where $f_{max} := \max_{x \in \mathcal{D}} f(x)$ and

$\hat{f}_{max} := \max_{x \in \mathcal{D}} \hat{f}(x)$, since $f(x) \leq \hat{f}(x) + \delta \leq f_{max} + \delta \forall x$, hence $f_{max} \leq \hat{f}_{max} + \delta$, and similarly $\hat{f}_{max} \leq f_{max} + \delta$. For $x = \pi$, $\mathcal{D} = \{\pi\}$, $f(x) = V_{k\gamma}^{\pi\mu}$, $\hat{f}(x) = V_{k\gamma}^{\pi\hat{\mu}}$ we get (ii) from (i).

(iii) With abbreviation $V_{\rho}^p := V_{k\gamma}^{p\rho}$ and noting that $V_{\hat{\mu}}^* \equiv V_{\hat{\mu}}^{p\hat{\mu}}$ we get $|V_{\mu}^* - V_{\mu}^{p\hat{\mu}}| = |V_{\mu}^* - V_{\hat{\mu}}^* + V_{\hat{\mu}}^{p\hat{\mu}} - V_{\mu}^{p\hat{\mu}}| \leq |V_{\mu}^* - V_{\hat{\mu}}^*| + |V_{\hat{\mu}}^{p\hat{\mu}} - V_{\mu}^{p\hat{\mu}}| \leq 2\delta(\varepsilon)$ by (ii) and (i), which proves (iii).

To prove $\delta(\varepsilon) \rightarrow 0$ we replace \min_n by $n \sim \varepsilon^{-1/2}$ and get $0 \leq \delta(\varepsilon) \leq r_{max} \cdot (|\mathcal{X}|(n - k)\varepsilon + \frac{\Gamma_n}{\Gamma_k}) \xrightarrow{\varepsilon \rightarrow 0} 0$, since $n \rightarrow \infty$, $\Gamma_n \xrightarrow{n \rightarrow \infty} 0$, and $n\varepsilon \rightarrow 0$. \square

The next theorem shows that, for a given policy p and history generated by p and μ , i.e. on-policy, the future universal value $V_{k..}^{p\xi}$ converges to the true value $V_{k..}^{p\mu}$.

Theorem 5.36 (Convergence of universal to true Value) If the history $y_{<k}$ is generated by policy p (and environment μ), and $V_{k..}^{p..} = V_{k..}^{p..}(y_{<k})$, then the universal undiscounted future value $V_{km_k}^{p\xi}$ with bounded dynamic horizon $h_k = m_k - k + 1 \leq h_{max}$ converges i.m.s. to the true value $V_{km_k}^{p\mu}$, and the discounted future value $V_{k\gamma}^{p\xi}$ converges i.m. to $V_{k\gamma}^{p\mu}$ for any summable discount sequence γ_k . In detail:

$$\begin{aligned} i) \quad & |V_{km}^{p\xi} - V_{km}^{p\mu}| \leq (m - k + 1)r_{max}a_{k:m}, \quad |V_{k\gamma}^{p\xi} - V_{k\gamma}^{p\mu}| \leq r_{max}\sqrt{2d_{k:\infty}} \\ ii) \quad & \sum_{k=1}^{\infty} \mathbf{E}(V_{km_k}^{p\xi} - V_{km_k}^{p\mu})^2 \leq 2h_{max}^3 r_{max}^2 D_{\infty}, \quad \mathbf{E}(V_{k\gamma}^{p\xi} - V_{k\gamma}^{p\mu})^2 \leq 2r_{max}^2 (D_{\infty} - D_{k-1}) \rightarrow 0 \\ iii) \quad & V_{km_k}^{p\xi} \xrightarrow{k \rightarrow \infty} V_{km_k}^{p\mu} \quad \text{i.m.s. if } h_{max} < \infty, \quad V_{k\gamma}^{p\xi} \xrightarrow{k \rightarrow \infty} V_{k\gamma}^{p\mu} \quad \text{i.m. for any } \gamma \\ & a_{k:m} := \sum_{x_{k:m}} |\mu(y_{<k} \underline{y}_{k:m}) - \xi(y_{<k} \underline{y}_{k:m})|, \quad d_{k:m} := \sum_{x_{k:m}} \mu(y_{<k} \underline{y}_{k:m}) \ln \frac{\mu(y_{<k} \underline{y}_{k:m})}{\xi(y_{<k} \underline{y}_{k:m})} \end{aligned}$$

and $D_k := d_{1:k} \leq \ln w_{\mu}^{-1} < \infty$ are defined as in Section 3.7.1 with actions y as additional conditions.

Proof. (i)_{left} follows from $|V_{km}^{p\xi} - V_{km}^{p\mu}| = \left| \sum_{x_{k:m}} (r_k + \dots + r_m)[\xi() - \mu()] \right| \leq$

$$\sum_{x_{k:m}} (r_k + \dots + r_m) |\xi() - \mu()| \leq (m - k + 1)r_{max} \sum_{x_{k:m}} |\xi() - \mu()| = (m - k + 1)r_{max}a_{k:m},$$

where $\rho() = \rho(y_{<k} \underline{y}_{k:m})|_{y_{1:m}=p(x_{<m})}$. (i)_{right} is shown similarly: Let $V_{km\gamma}$ be the discounted future value $V_{k\gamma}$ but cut after cycle m . We have

$$|V_{km\gamma}^{p\xi} - V_{km\gamma}^{p\mu}| = \frac{1}{\Gamma_k} \left| \sum_{x_{k:m}} (\gamma_k r_k + \dots + \gamma_m r_m)[\xi() - \mu()] \right| \leq \dots \leq r_{max}a_{k:m} \leq r_{max}\sqrt{2d_{k:m}}.$$

In the last step we used $a_{k:m} \leq \sqrt{2d_{k:m}}$ (see Lemma 3.11 or Section 3.7.1). (i)_{right} follows by taking the limit $m \rightarrow \infty$, which exists since $V_{km\gamma}$ and $d_{k:m}$ are monotone

increasing in m and bounded. $(ii)_{left}$ follows from $(i)_{left}$, and $m_k - k + 1 \leq h_{max}$, and bound (3.71) with $t \rightsquigarrow k$, $n_t \rightsquigarrow m_k$, $h \rightsquigarrow h_{max}$, and y as additional conditions. $(ii)_{right}$ follows from $(i)_{right}$ and $\mathbf{E}[d_{k:n}] = D_n - D_{k-1}$ (see Section 3.7.1). (iii) follows from (ii) by Definition 3.8 of convergence i.m.(s.) \square

Convergence of the average values $\frac{1}{h_k} V_{km_k}^{p^\xi} \rightarrow h_k^{-1} V_{km_k}^{p^\mu}$ also holds, i.m.s. for bounded horizon, and i.m. for arbitrary horizon. Note also that if the history is generated by p^ξ , then (iii) implies $V_{k\gamma}^{*\xi} \rightarrow V_{k\gamma}^{p^\xi\mu}$, hence the universal value $V_{k\gamma}^{*\xi}$ can be used to estimate the true value $V_{k\gamma}^{p^\xi\mu}$, without any assumptions on \mathcal{M} and γ . Nevertheless, maximization of $V_{k\gamma}^{p^\xi}$ may asymptotically differ from maximization of $V_{k\gamma}^{p^\mu}$, since $V_{k\gamma}^{p^\xi} \not\rightarrow V_{k\gamma}^{p^\mu}$ for $p \neq p^\xi$ is possible (and also $V_{k\gamma}^{*\xi} \not\rightarrow V_{k\gamma}^{*\mu}$, see Section 5.3.2 and Problem 5.2).

5.6 Markov Decision Processes (MDP)

From all possible environments, Markov (Decision) Processes are probably the most intensively studied ones. To give an example, we apply Theorems 5.29 and 5.34 to ergodic Markov Decision processes (MDPs), but we will be very brief.

Definition 5.37 (Ergodic Markov Decision Processes) We call μ a (stationary) *Markov Decision Process (MDP)* if the probability of observing $x_k \in \mathcal{X}$, given history $y_{<k} y_k$ does only depend on the last action $y_k \in \mathcal{Y}$ and the last observation x_{k-1} , i.e. if $\mu(y_{<k} y_k \underline{x}_k) = \mu(x_{k-1} y_k \underline{x}_k)$. In this case x_k is called a *state*, \mathcal{X} the *state space*, and $\mu(x_{k-1} y_k \underline{x}_k)$ the *transition matrix*. An MDP μ is called *ergodic* if there exists a policy under which every state is visited infinitely often with probability 1. Let \mathcal{M}_{MDP} be the set of MDPs and \mathcal{M}_{MDP1} be the set of ergodic MDPs. If an MDP $\mu(x_{k-1} y_k \underline{x}_k)$ is independent of the action y_k it is a *Markov process*, if it is independent of x_{k-1} it is an *i.i.d.* process.

Stationary MDPs μ have stationary optimal policies p^μ mapping the same state/observation x_t always to the same action y_t . On the other hand a mixture ξ of MDPs is itself not an MDP, i.e. $\xi \notin \mathcal{M}_{MDP}$, which implies that p^ξ is, in general, not a stationary policy. The definition of ergodicity given here is least demanding, since it only demands the existence of a single policy under which the Markov process is ergodic. Often, stronger assumptions, e.g. that every policy is ergodic or that a stationary distribution exists, are made. We now show that there are self-optimizing policies for the class of ergodic MDPs in the following sense.

Theorem 5.38 (Self-optimizing policies for ergodic MDPs) There exist self-optimizing policies \tilde{p}_m for the class of ergodic MDPs in the sense that

$$i) \quad \exists \tilde{p}_m \forall \nu \in \mathcal{M}_{MDP1} : \frac{1}{m} V_{1m}^{*\nu} - \frac{1}{m} V_{1m}^{\tilde{p}_m \nu} = O(m^{-1/3})$$

In the discounted case, if the discount sequence γ_k has unbounded effective horizon $h_k^{eff} \xrightarrow{k \rightarrow \infty} \infty$, then there exist self-optimizing policies $\tilde{\pi}_k$ for the class of ergodic MDPs in the sense that

$$ii) \quad \exists \tilde{\pi}_k \forall \nu \in \mathcal{M}_{MDP1} : V_{k\gamma}^{\tilde{\pi}_k \nu} \xrightarrow{k \rightarrow \infty} V_{k\gamma}^{*\nu} \quad \text{for any history } y_{<k} \quad \text{if } \frac{\gamma_{k+1}}{\gamma_k} \rightarrow 1.$$

There is much literature on constructing and analyzing self-optimizing learning algorithms in MDP environments. The assumptions on the structure of the MDPs vary, all include some form of ergodicity, often stronger than Definition 5.37, demanding that the Markov process is ergodic under *every* policy. See, for instance, [KV86, BT96]. We will only briefly outline one algorithm satisfying Theorem 5.38 without trying to optimize performance.

Proof. For (i) one can choose a policy \tilde{p}_m which performs (uniformly) random actions in cycles $1 \dots k_0 - 1$ with $1 \ll k_0 \ll m$ and which follows thereafter the optimal policy based on an estimate of the transition matrix $T_{ss'}^a \equiv \nu(sas')$ from the initial $k_0 - 1$ cycles. The existence of an ergodic policy implies that for every pair of states $s_{start}, s \in \mathcal{X}$ there is a sequence of actions and transitions of length at most $|\mathcal{X}| - 1$ such that state s is reached from state s_{start} . The probability that the “right” transition occurs is at least T_{min} with T_{min} being the smallest non-zero transition probability in T . The probability that a random action is the “right” action is at least $|\mathcal{Y}|^{-1}$. So the probability of reaching a state s in $|\mathcal{X}| - 1$ cycles via a random policy is at least $(T_{min}/|\mathcal{Y}|)^{|\mathcal{X}|-1}$. In state s action a is taken with probability $|\mathcal{Y}|^{-1}$ and leads to state s' with probability $T_{ss'}^a \geq T_{min}$. Hence, the expected number of transitions $s \xrightarrow{a} s'$ to occur in the first k_0 cycles is $\geq \frac{k_0}{|\mathcal{X}|} (T_{min}/|\mathcal{Y}|)^{|\mathcal{X}|} \sim k_0$.⁸ The accuracy of the frequency estimate $\hat{T}_{ss'}^a$ of $T_{ss'}^a$ hence is $\sim k_0^{-1/2}$. Similar MDPs lead to “similar” optimal policies, which lead to similar values. More precisely, one can show (see Problem 5.12) that $\hat{T} - T \sim k_0^{-1/2}$ implies the same accuracy in the average value, i.e. $|\frac{1}{m} V_{k_0 m}^{\tilde{p}_m \nu} - \frac{1}{m} V_{k_0 m}^{*\nu}| \sim k_0^{-1/2}$, where \tilde{p}_m is the optimal policy based on \hat{T} and $*$ is the optimal policy based on $T (= \nu)$. Since $\frac{1}{m} V_{1k_0} \sim \frac{k_0}{m}$, (i) follows (with probability 1) by setting $k_0 \sim m^{2/3}$. The policy \tilde{p}_m can be derandomized, showing (i) for sure.

The discounted case (ii) can be proven similarly. The history $y_{<k}$ is simply ignored and the analogue to $m \rightarrow \infty$ is $h_k^{eff} \rightarrow \infty$ for $k \rightarrow \infty$, which is ensured by $\frac{\gamma_{k+1}}{\gamma_k} \rightarrow 1$. Let $\tilde{\pi}_k$ be the policy which performs (uniformly) random actions in cycles $k \dots k_0 - 1$ with $k \ll k_0 \ll h_k^{eff}$ and which follows thereafter the optimal policy⁹ based

⁸For $T_{ss'}^a = 0$ the estimate $\hat{T}_{ss'}^a = 0$ is exact.

⁹For non-geometric discounts as here, optimal policies are, in general, *not* stationary.

on an estimate \hat{T} of the transition matrix T from cycles $k \dots k_0 - 1$. The existence of an ergodic policy, again, ensures that the expected (after derandomization for sure) number of transitions $s \xrightarrow{a} s'$ occurring in cycles $k \dots k_0 - 1$ is proportional to $\Delta := k_0 - k$. The accuracy of the frequency estimate \hat{T} of T is $\sim \Delta^{-1/2}$ which implies by a strengthening of Theorem 5.35(iii) for ergodic MDPs similar to Problem 5.12 that

$$V_{k_0\gamma}^{\tilde{\pi}_k\nu} \rightarrow V_{k_0\gamma}^{*\nu} \quad \text{for } \Delta = k_0 - k \rightarrow \infty, \quad (5.39)$$

where $\tilde{\pi}_k$ is the optimal policy based on \hat{T} and $*$ is the optimal policy based on $T(=\nu)$. It remains to show that the achieved reward in the random phase $k \dots k_0 - 1$ gives a negligible contribution to $V_{k\gamma}$. The following implications for $k \rightarrow \infty$ are easy to show:

$$\frac{\gamma_{k+1}}{\gamma_k} \rightarrow 1 \Rightarrow \frac{\gamma_{k+\Delta}}{\gamma_k} \rightarrow 1 \Rightarrow \frac{\Gamma_{k+\Delta}}{\Gamma_k} \rightarrow 1 \Rightarrow \frac{1}{\Gamma_k} \sum_{i=k}^{k_0-1} \gamma_i r_i \leq \frac{r_{max}}{\Gamma_k} [\Gamma_{k+\Delta} - \Gamma_k] \rightarrow 0.$$

Since convergence to zero is true for all fixed finite Δ it is also true for sufficiently slowly increasing $\Delta(k) \rightarrow \infty$. This shows that the contribution of the first Δ rewards $r_k + \dots + r_{k_0-1}$ to $V_{k\gamma}$ is negligible. Together with (5.39) this shows $V_{k\gamma}^{\tilde{\pi}_k\nu} \rightarrow V_{k\gamma}^{*\nu}$ for $k_0 := k + \Delta(k)$. \square

The rate of convergence $m^{-1/3}$ rather than $m^{-1/2}$ in the undiscounted case may be a bit surprising. If we would explore for $k_0 = m$ steps and ask for the accuracy of the value function estimate afterwards, i.e. for $k > m$, we would get $\sim k_0^{-1/2} = m^{-1/2}$, of course, but since we are considering $\frac{1}{m} V_{1m}$ including the history $1..k_0$ we must get a worse result, namely $m^{-1/3}$. Although in the discounted case, we consider the future value $V_{k\gamma}$, the situation is nevertheless similar. Exploration takes place in cycles $k \dots k_0 - 1$, history $y_{x < k}$ is not exploited.

The conditions $\Gamma_k < \infty$ and $\frac{\gamma_{k+1}}{\gamma_k} \rightarrow 1$ on the discount sequence are, for instance, satisfied for $\gamma_k = 1/k^2$, so the theorem is not vacuous. The popular geometric discount $\gamma_k = \gamma^k$ fails the latter condition; it has finite effective horizon. Section 5.7 will give a detailed account on the discount and horizon issues.

Together with Theorems 5.29 and 5.34, Theorem 5.38 immediately implies that AIXI is self-optimizing for the class of ergodic MDPs.

Corollary 5.40 (AIXI is self-optimizing for ergodic MDPs) If \mathcal{M} is a countable class of ergodic MDPs, and $\xi := \sum_{\nu \in \mathcal{M}} w_\nu \nu$, then AIXI alias p_m^ξ maximizing $V_{1m}^{p_m^\xi}$ and p^ξ maximizing $V_{k\gamma}^{\pi^\xi}$ are self-optimizing in the sense that

$$\forall \nu \in \mathcal{M} : \frac{1}{m} V_{1m}^{p_m^\xi \nu} \xrightarrow{m \rightarrow \infty} \frac{1}{m} V_{1m}^{*\nu} \quad \text{and} \quad V_{k\gamma}^{p^\xi \nu} \xrightarrow{k \rightarrow \infty} V_{k\gamma}^{*\nu} \quad \text{if } \frac{\gamma_{k+1}}{\gamma_k} \rightarrow 1.$$

If \mathcal{M} is finite, then the speed of the first convergence is at least $O(m^{-1/3})$.

Continuous classes \mathcal{M} . There are uncountably many ergodic MDPs. Since we have restricted our development to countable classes \mathcal{M} we had to give the Corollary for a

countable subset of \mathcal{M}_{MDP1} . We may choose \mathcal{M} as the set of all ergodic MDPs with rational (or computable) transition probabilities. In this case \mathcal{M} is a dense subset of \mathcal{M}_{MDP1} which is, from a practical point of view, sufficiently rich. On the other hand, it is possible to extend the theory to continuously parameterized families of environments μ_θ and $\xi = \int d\theta w_\theta \mu_\theta$. Under some mild (differentiability and existence) conditions, most results of this work remain valid in some form, especially Corollary 5.40 for *all* ergodic MDPs \mathcal{M}_{MDP1} .

Bayesian self-optimizing policy. $AI_{\xi_{MDP1}}$ with unbounded horizon is the first purely Bayesian self-optimizing consistent policy for ergodic MDPs. The policies of all previous approaches were either hand crafted, like those in the proof of Theorem 5.38, or were Bayesian with a pre-chosen horizon m or with geometric discounting γ with finite effective horizon [KV86, BT96]. The combined conditions $\Gamma_k < \infty$ and $\frac{\gamma_{k+1}}{\gamma_k} \rightarrow 1$ allow a consistent self-optimizing Bayesian policy based on mixtures.

Bandits. For instance, consider the popular class of bandits B. In a two-armed bandit problem you pull repeatedly one out of two levers resulting in a gain of 1\$ with probability p_i for arm number i . The game can be described as an MDP with parameters p_i . If the p_i are unknown, Corollary 5.40 shows that AI_{ξ_B} yields asymptotically optimal payoff. The discounted unbounded horizon approach and result for Bandits is, to the best of our knowledge, also new.

Other environmental classes. Bandits, i.i.d. processes, classification tasks, and many more are all special (degenerate) cases of ergodic MDPs, for which Corollary 5.40 shows that p^ξ is self-optimizing. But the existence of self-optimizing policies is not limited to (subclasses of ergodic) MDPs. Certain classes of POMDPs, k^{th} order ergodic MDPs, factorizable environments, repeated games, and prediction problems are not MDPs, but nevertheless admit self-optimizing policies (to be shown elsewhere), and hence the corresponding Bayes-optimal mixture policy p^ξ is self-optimizing by Theorems 5.29 and 5.34.

Restricted policy classes. The development in this and the last paragraphs can be scaled down to restricted classes of policies \mathcal{P} . If one defines $V^* = \arg\max_{p \in \mathcal{P}} V^p$ all theorems remain valid, more or less unchanged. For instance, consider a finite class of quickly computable policies. For MDPs, ξ is quickly computable and V_ξ^p can be (efficiently) computed by Monte-Carlo sampling. Maximizing over the finitely many policies $p \in \mathcal{P}$ selects the asymptotically best policy p^ξ from \mathcal{P} for all ergodic MDPs.

Outlook. Future research could be the derivation of non-asymptotic bounds, possibly along the lines of [Hut01b]. To get good bounds one may have to exploit extra properties of the environments, like the mixing rate of MDPs [KS98]. Finally, instead of convergence of the expected reward sum, convergence with high probability of the actual reward sum, would be interesting to study (cf. Problem 3.4).

5.7 The Choice of the Horizon

The only significant arbitrariness in the AI ξ model lies in the choice of the horizon function $h_k \equiv m_k - k + 1$. We discuss some choices which seem to be natural and give preliminary conclusions at the end. We will not discuss ad hoc choices of h_k for specific problems (like the discussion in Section 6.3 in the context of finite strategic games). We are interested in universal choices of m_k .

Fixed horizon. If the lifetime of the agent is known to be m , which is in practice always large but finite, then the choice $m_k = m$ maximizes correctly the expected future reward. m is usually not known in advance, as in many cases the time we are willing to run an agent depends on the quality of its outputs. For this reason, it is often desirable that good outputs are not delayed too much, if this results in a marginal reward increase only. This can be incorporated by damping the future rewards. If, for instance, we assume that the survival of the agent in each cycle is proportional to the past reward an exponential damping (geometric discounting) $r_k := r'_k \cdot e^{-\lambda k}$ is appropriate, where r'_k are bounded, e.g. $r'_k \in [0,1]$. The expression (5.3) converges for $m_k \rightarrow \infty$ in this case¹⁰. But this does not solve the problem, as we introduced a new arbitrary time-scale $1/\lambda$. Every damping introduces a time-scale. Taking $\lambda \rightarrow 0$ is prone to the same problems as $m_k \rightarrow \infty$ in the undiscounted case.

General discounting. Geometric discounting does not solve the horizon problem, but the idea of discounting is fruitful. Let $r_k := \gamma_k r'_k$ with $\gamma_k > 0$ and $r'_k \in [0,1]$. If $\Gamma_k := \sum_{i=k}^{\infty} \gamma_i < \infty$, then $V_{k\gamma}^{pp} := \frac{1}{\Gamma_k} \lim_{m \rightarrow \infty} V_{km}^{pp}$ exists. Rewards r_{k+h} give only a small contribution to $V_{k\gamma}^{pp}$ for large h , since $\gamma_{k+h} \xrightarrow{h \rightarrow \infty} 0$. The instantaneous effective horizon may be defined as the \hat{h} for which $\gamma_{k+\hat{h}}$ is only half (or more generally a fraction $\beta < 1$) of γ_k . Formally we may define $\hat{h}_k^\beta := \min\{h \geq 0 : \gamma_{k+h} \leq \beta \gamma_k\}$. For any discount γ_k we have $\hat{h}_k^\beta \leq c \cdot k$ for some constant c independent of k . A better definition for the β -effective horizon is the h for which the cumulative discount $\Gamma_{k+h} \approx \beta \Gamma_k$, or more formally, $h_k^\beta := \min\{h \geq 0 : \Gamma_{k+h} \leq \beta \Gamma_k\}$. Approximating the infinite reward sum in $V_{k\gamma}$ by the first h_k^β terms introduces an error of at most βr_{max} . We define *the effective horizon* by $h_k^{eff} := h_k^{\beta=1/2}$. Table 5.41 shows effective horizons for various types of discounts γ_k .

Dynamic horizon (universal & harmonic discounting). The largest horizon with guaranteed finite and enumerable reward sum can be obtained by the universal discount $\gamma_k = 2^{-K(k)}$ (or the monotone variant $\gamma_k = \min_{i \leq k} 2^{-K(i)}$). This discount results in a truly farsighted agent with effective horizon which grows faster than any computable function. It is somewhat similar to a near-harmonic discount $\gamma_k = [k \log^2 k]^{-1}$, since $2^{-K(k)} \leq 1/k$ for most k and $2^{-K(k)} \geq c/(k \log^2 k)$ (see Theorem 2.10(ii)), but leads to $h_k^{eff} \sim k^2$. Similarly, the time-scale invariant power damping $\gamma_k = k^{-1-\varepsilon}$ introduces a dynamic time-scale. In cycle k the contribution of cycle

¹⁰More precisely $\dot{y}_k = \arg\max_{y_k} \lim_{m_k \rightarrow \infty} V_{km_k}^{*\xi}(\dot{y}_{x < k} y_k)$ exists.

Table 5.41 (Effective horizons) The table shows the effective horizons $\hat{h}_k^\beta := \min\{h \geq 0 : \gamma_{k+h} \leq \beta\gamma_k\}$ and $h_k^\beta := \min\{h \geq 0 : \Gamma_{k+h} \leq \beta\Gamma_k\}$ for various types of discounts γ_k .

Horizons	γ_k	\hat{h}_k^β	$\Gamma_k = \sum_{i=k}^\infty \gamma_i$	h_k^β	$h_k^{\text{eff}} = h_k^{\beta=1/2}$
Finite	$\begin{cases} 1 & \text{for } k \leq m \\ 0 & \text{for } k > m \end{cases}$	$m - k + 1$	$m - k + 1$	$\lceil (1-\beta)(m-k+1) \rceil$	$\lceil \frac{1}{2}(m-k+1) \rceil$
Geometric	$\gamma_k, 0 \leq \gamma < 1$	$\lceil \frac{\ln \beta}{\ln \gamma} \rceil$	$\frac{\gamma^k}{1-\gamma}$	$\lceil \frac{\ln \beta}{\ln \gamma} \rceil$	$\approx \frac{\ln 2}{1-\gamma}$ for $\gamma \approx 1$
Power	$k^{-1-\varepsilon}, \varepsilon > 0$	$\sim (\beta^{-\frac{1}{1+\varepsilon}} - 1)k$	$\sim \frac{1}{\varepsilon} k^{-\varepsilon}$	$\sim (\beta^{-1/\varepsilon} - 1)k$	$\propto k$
Near-Harmonic	$\frac{1}{k \ln^{1+\varepsilon} k}, \varepsilon > 0$	$\sim (\beta^{-1} - 1)k$	$\sim \frac{1}{\varepsilon} (\ln k)^{-\varepsilon}$	$\sim k^{\beta^{-1/\varepsilon}}$	$\sim k^{2^{1/\varepsilon}}$
Universal	$2^{-K(k)}$	$\approx k$ on average	decreases slower than any computable function	increases faster than any computable function	

$2^{\frac{1}{1+\varepsilon}} \cdot k$ is damped by a factor $\frac{1}{2}$. The instantaneous effective horizon \hat{h}_k in this case is $\sim k$, the maximum possible. The choice $h_k = \alpha \cdot k$ with $\alpha \sim 2^{\frac{1}{1+\varepsilon}}$ qualitatively models the same behavior. We have not introduced an arbitrary time-scale m , but limited the farsightedness to some multiple (or fraction) of the length of the current history. This avoids the pre-selection of a global time-scale m or $1/\lambda$. This choice has some appeal, as it seems that humans of age k years usually do not plan their lives for more than, perhaps, the next k years ($\alpha_{\text{human}} \approx 1$). From a practical point of view this model might serve all needs, but from a theoretical point we feel uncomfortable with such a limitation in the horizon from the very beginning. Note, that we have to choose $\alpha = O(1)$ because otherwise we would again introduce a number α , which has to be justified. We favor the universal discount $\gamma_k = 2^{-K(k)}$, since it allows us, if desired, to “mimic” all other more greedy behaviors based on other discounts γ_k by choosing $r_k \in [0, c \cdot \gamma_k] \subseteq [0, 2^{-K(k)}]$.

Infinite horizon. The naive limit $m_k \rightarrow \infty$ in (5.3) may turn out to be well defined and the previous discussion superfluous. In the following, we suggest a limit which is always well defined (for finite \mathcal{Y}). Let $\dot{y}_k^{(m_k)}$ be defined as in (5.3) with dependence on m_k made explicit. Further, let $\dot{\mathcal{Y}}_k^{(m)} := \{\dot{y}_k^{(m_k)} : m_k \geq m\}$ be the set of outputs in cycle k for the choices $m_k = m, m+1, m+2, \dots$. Because $\dot{\mathcal{Y}}_k^{(m)} \supseteq \dot{\mathcal{Y}}_k^{(m+1)} \neq \{\}$, we have $\dot{\mathcal{Y}}_k^{(\infty)} := \bigcap_{m=k}^\infty \dot{\mathcal{Y}}_k^{(m)} \neq \{\}$. We define the $m_k = \infty$ model to output any $\dot{y}_k^{(\infty)} \in \dot{\mathcal{Y}}_k^{(\infty)}$. This is the best output consistent with some arbitrary large choice of m_k . Choosing the lexicographically smallest $\dot{y}_k^{(\infty)} \in \dot{\mathcal{Y}}_k^{(\infty)}$ would correspond to the limes inferior $\lim_{m \rightarrow \infty} \dot{y}_k^{(m)}$, which always exists (for finite \mathcal{Y}). Generally $\dot{y}_k^{(\infty)} \in \dot{\mathcal{Y}}_k^{(\infty)}$ is unique, i.e. $|\dot{\mathcal{Y}}_k^{(\infty)}| = 1$ iff the naive limit $\lim_{m \rightarrow \infty} \dot{y}_k^{(m)}$ exists. Note, that the limit $\lim_{m \rightarrow \infty} V_{km}^*(y_{<k})$ needs not to exist for this construction.

Average reward and differential gain. Taking the raw average reward $(r_k + \dots + r_m)/(m - k + 1)$ and $m \rightarrow \infty$ also does not help: take an arbitrary policy for the first k time steps and the/an optimal policy for the remaining steps $k+1 \dots \infty$. All these policies give the same average. In MDP environments with a single recurrent class one

can define the relative or differential gain [Ber95b]. In more general environments (we are interested in) the differential gain can be infinite, which is acceptable, since differential gains can still be totally ordered. The major problem is the *existence* of the differential gain, i.e. whether it converges for $m \rightarrow \infty$ in $\mathbb{R} \cup \{\infty\}$ at all (and does not oscillate). This is just the old convergence problem in slightly different form.

Immortal agents are lazy. The construction above leads to a mathematically elegant, no-parameter $\text{AI}\xi$ model. Unfortunately this is not the end of the story. The limit $m_k \rightarrow \infty$ can cause undesirable results in the $\text{AI}\mu$ model for special μ , which might also happen in the $\text{AI}\xi$ model whatever we define $m_k \rightarrow \infty$. Consider $\mathcal{Y} = \mathcal{X} = \mathcal{R} = \{0,1\}$. Output $y_k = 0$ shall give reward $r_k = 0$ and output $y_k = 1$ shall give $r_k = 1$ iff $\dot{y}_{k-l-\sqrt{l}} \dots \dot{y}_{k-l} = 0 \dots 0$ for some l . I.e. the agent can achieve l consecutive positive rewards if there was a sequence of length at least \sqrt{l} with $y_k = r_k = 0$. If the lifetime of the $\text{AI}\mu$ agent is m , it outputs $\dot{y}_k = 0$ in the first r cycles and then $\dot{y}_k = 1$ for the remaining r^2 cycles with r such that $r + r^2 = m$. This will lead to the highest possible total reward $V_{1m} = r^2 = m + \frac{1}{2} - \sqrt{m + \frac{1}{4}}$. Any fragmentation of the 0 and 1 sequences would reduce V_{1m} . For $m \rightarrow \infty$ the $\text{AI}\mu$ agent can and will delay the point r of switching to $\dot{y}_k = 1$ indefinitely and always output 0 leading to total reward 0, obviously the worst possible behavior. The $\text{AI}\xi$ agent will explore the above rule after a while of trying $y_k = 0/1$ and then applies the same behavior as the $\text{AI}\mu$ agent, since the simplest rules covering past data dominate ξ . For finite m this is exactly what we want, but for infinite m the $\text{AI}\xi$ model (probably) fails just as the $\text{AI}\mu$ model does. The good point is, that this is not a weakness of the $\text{AI}\xi$ model in particular, as $\text{AI}\mu$ fails too. The bad point is that $m_k \rightarrow \infty$ has far reaching consequences, even when starting from an already very large $m_k = m$. The reason being that the μ of this example is highly non-local in time, i.e. it may violate one of our weak separability conditions.

Conclusions. We are not sure whether the choice of m_k is of marginal importance, as long as m_k is chosen sufficiently large and of low complexity, $m_k = 2^{2^{16}}$ for instance, or whether the choice of m_k will turn out to be a central topic for the $\text{AI}\xi$ model or for the planning aspect of any AI system in general. We suppose that the limit $m_k \rightarrow \infty$ for the $\text{AI}\xi$ model results in correct behavior for weakly separable μ . A proof of this conjecture, if true, would probably give interesting insights.

5.8 Outlook

Expert advice approach. We considered expected performance bounds for predictions based on Solomonoff's prior. The other, dual, currently very popular approach, is "prediction with expert advice" (PEA) invented by Littlestone and Warmuth (1989), Vovk (1992). Whereas PEA performs well in any environment, but only relative to a given set of experts, our Λ_ξ predictor competes with *any* other predictor, but only in expectation for environments with computable distribution. It

seems philosophically less compromising to make assumptions on prediction strategies than on the environment, however weak. One could investigate whether PEA can be generalized to the case of active agents, which would result in a model dual to AIXI. We believe the answer to be negative, which on the positive side would show the necessity of Occam's razor assumption, and the distinguishedness of AIXI.

Actions as random variables. The uniqueness for the choice of the generalized ξ (2.24) in the AIXI model could be explored. From the originally many alternatives, which could all be ruled out, there is one alternative which still seems possible. Instead of defining ξ as in (5.2) one could treat the agent's actions y also as universally distributed random variables and then conditionalize ξ on y by the chain rule (see Problem 5.1).

Structure of AIXI. The algebraic properties and the structure of AIXI could be investigated in more depth (we already saw that the value V_μ^p is a linear function in μ and V_μ^* is a convex function in μ). This would extract the essentials from AIXI which finally could lead to an axiomatic characterization of AIXI. The benefit is as in any axiomatic approach. It would clearly exhibit the assumptions, separate the essentials from technicalities, and simplify understanding and, most important, guide in finding proofs.

Posterization. Many properties of Kolmogorov complexity, Solomonoff's prior, and (policies based on) Bayes-mixtures remain valid after "posterization". With posterization we mean replacing V_{1m} , w_ν , $K(\nu)$, $\nu(\underline{y}_{1:m})$, etc. by the posteriors V_{km} , w_k^ν , $K(\nu|\underline{y}_{<k})$, $\nu(\underline{y}_{<k}\underline{y}_{k:m})$, etc. Strangely enough for w_ν chosen as $2^{-K(\nu)}$ it is not true that $w_k^\nu \sim 2^{-K(\nu|\underline{y}_{<k})}$. If this property were true, weak bounds as the one proven in Section 6.2 (which is too weak to be of practical importance) could be boosted to practical bounds of order 1. Hence, it is of high impact to rescue the posterization property in some way. It may be valid when grouping together essentially equal distributions ν .

5.9 Conclusions

All tasks which require intelligence to be solved can naturally be formulated as a maximization of some expected utility in the framework of agents. We gave an explicit expression (4.17) of such a decision theoretic agent. The main remaining problem is the unknown prior probability distribution μ^{AI} of the environment(s). Conventional learning algorithms are unsuitable, because they can neither handle large (unstructured) state spaces, nor do they converge in the theoretically minimal number of cycles, nor can they handle non-stationary environments appropriately. On the other hand, the universal semimeasure ξ (2.24), based on ideas from algorithmic information theory, solves the problem of the unknown prior distribution for induction problems. No explicit learning procedure is necessary, as ξ automatically converges to μ . We unified the theory of universal sequence prediction with the

decision theoretic agent by replacing the unknown true prior μ^{AI} by an appropriately generalized universal semimeasure ξ^{AI} . We gave strong arguments that the resulting AI ξ model is universally optimal. Furthermore, possible solutions to the horizon problem have been discussed. In Chapter 6 we outline for a number of problem classes, how the AI ξ model can solve them. They include sequence prediction, strategic games, function minimization and, especially, how AI ξ learns to learn supervised. In Chapter 7 we develop a more elegant but equivalent functional form of the AI ξ model. It is used to define a universal intelligence order relation, to discuss in which sense AI ξ is the most intelligent agent, and to construct a modified time-bounded (computable) AI ξ^{tl} version.

5.10 Converting Functions into Chronological Semimeasures

To complete the proof of the universality (5.6) of ξ we need to convert enumerable functions $\psi: \mathcal{B}^* \rightarrow \mathbb{R}^+$ into enumerable chronological semi-measures $\rho: (\mathcal{Y} \times \mathcal{X})^* \rightarrow \mathbb{R}^+$ with certain additional properties. The proof given here follows [LV97, p274], but is slightly more formal and compact. Every enumerable function like ψ and ρ can be approximated from below by definition¹¹ by primitive recursive functions $\varphi: \mathcal{B}^* \times \mathbb{N} \rightarrow \mathbb{Q}^+$ and $\phi: (\mathcal{Y} \times \mathcal{X})^* \times \mathbb{N} \rightarrow \mathbb{Q}^+$ with $\psi(s) = \sup_t \varphi(s, t)$ and $\rho(s) = \sup_t \phi(s, t)$ and recursion parameter t . For arguments of the form $s = yx_{1:n}$ we recursively (in n) construct ϕ from φ as follows:

$$\varphi'(yx_{1:n}, t) := \begin{cases} \varphi(yx_{1:n}, t) & \text{for } x_n < t \\ 0 & \text{for } x_n \geq t \end{cases}, \quad \varphi'(\epsilon, t) := \varphi(\epsilon, t) \quad (5.42)$$

$$\phi(\epsilon, t) := \max_{0 \leq i \leq t} \{ \varphi'(\epsilon, i) : \varphi'(\epsilon, i) \leq 1 \} \quad (5.43)$$

$$\phi(yx_{1:n}, t) := \max_{0 \leq i \leq t} \{ \varphi'(yx_{1:n}, i) : \sum_{x_n} \varphi'(yx_{1:n}, i) \leq \phi(yx_{<n}, t) \} \quad (5.44)$$

With $x_n < t$ we mean that the natural number associated with string x_n is smaller than t . According to (5.42) with φ also φ' as well as $\sum_{x_n} \varphi'$ are primitive recursive functions. Further, if we allow $t = 0$ we have $\varphi'(s, 0) = 0$. This ensures that ϕ is a total function.

In the following we prove by induction over n that ϕ is a primitive recursive chronological semimeasure monotone increasing in t . All necessary properties hold for $n = 0$ ($yx_{1:0} = \epsilon$) according to (5.43). For general n assume that the induction hypothesis is true for $\phi(yx_{<n}, t)$. We can see from (5.44) that $\phi(yx_{1:n}, t)$ is monotone increasing in t . ϕ is total as $\varphi'(yx_{1:n}, i = 0) = 0$ satisfies the inequality. By assumption

¹¹Defining enumerability as the supremum of total primitive recursive functions is more suitable for our purpose than the equivalent definition as a limit of monotone increasing partial recursive functions. In terms of Turing machines, the recursion parameter is the time after which a computation is terminated.

$\phi(\underline{y}_{<n}, t)$ is primitive recursive, hence with $\sum_{x_n} \varphi'$ also the order relation $\sum \varphi' \leq \phi$ is primitive recursive. This ensures that the non-empty finite set $\{\varphi' : \sum \varphi' \leq \phi\}_i$ and its maximum $\phi(\underline{y}_{1:n}, t)$ are primitive recursive. Further, $\phi(\underline{y}_{1:n}, t) = \varphi'(\underline{y}_{1:n}, i)$ for some i with $i \leq t$ independent of x_n . Thus, $\sum_{x_n} \phi(\underline{y}_{1:n}, t) = \sum_{x_n} \varphi'(\underline{y}_{1:n}, i) \leq \phi(\underline{y}_{<n}, t)$ which is the condition for ϕ being a chronological semimeasure. Inductively we have proved that ϕ is indeed a primitive recursive chronological semimeasure monotone increasing in t .

In the following we show that every (total)¹² enumerable chronological semimeasure ρ can be enumerated by some ϕ . By definition of enumerability there exist primitive recursive functions $\tilde{\varphi}$ with $\rho(s) = \sup_t \tilde{\varphi}(s, t)$. The function $\varphi(s, t) := (1 - 1/t) \cdot \max_{i < t} \tilde{\varphi}(s, i)$ also enumerates ρ but has the additional advantage of being strictly monotone increasing in t .

$\varphi'(\underline{y}_{1:n}, \infty) = \varphi(\underline{y}_{1:n}, \infty) = \rho(\underline{y}_{1:n})$ by definition (5.42). $\phi(\epsilon, t) = \varphi'(\epsilon, t)$ by (5.43) and the fact that $\varphi'(\epsilon, i-1) < \varphi'(\epsilon, i) \leq \varphi(\epsilon, i) \leq \rho(\epsilon) \leq 1$, hence $\phi(\epsilon, \infty) = \rho(\epsilon)$. $\phi(\underline{y}_{1:n}, t) \leq \varphi'(\underline{y}_{1:n}, t)$ by (5.44), hence $\phi(\underline{y}_{1:n}, \infty) \leq \rho(\underline{y}_{1:n})$. We prove the opposite direction $\phi(\underline{y}_{1:n}, \infty) \geq \rho(\underline{y}_{1:n})$ by induction over n . We have

$$\sum_{x_n} \varphi'(\underline{y}_{1:n}, i) \leq \sum_{x_n} \varphi(\underline{y}_{1:n}, i) < \sum_{x_n} \varphi(\underline{y}_{1:n}, \infty) = \sum_{x_n} \rho(\underline{y}_{1:n}) \leq \rho(\underline{y}_{<n}) \quad (5.45)$$

The strict monotony of φ and the semimeasure property of ρ have been used. By induction hypothesis $\lim_{t \rightarrow \infty} \phi(\underline{y}_{<n}, t) \geq \rho(\underline{y}_{<n})$ and (5.45) for sufficiently large t we have $\phi(\underline{y}_{<n}, t) > \sum_{x_n} \varphi'(\underline{y}_{1:n}, i)$. The condition in (5.44) is, hence, satisfied and therefore $\phi(\underline{y}_{1:n}, t) \geq \varphi'(\underline{y}_{1:n}, i)$ for sufficiently large t , especially $\phi(\underline{y}_{1:n}, \infty) \geq \varphi'(\underline{y}_{1:n}, i)$ for all i . Taking the limit $i \rightarrow \infty$ we get $\phi(\underline{y}_{1:n}, \infty) \geq \varphi'(\underline{y}_{1:n}, \infty) = \rho(\underline{y}_{1:n})$.

Combining all results, we have shown that the constructed $\phi(\cdot, t)$ are primitive recursive chronological semimeasures monotone increasing in t , which converge to the enumerable chronological semimeasure ρ . This finally proves the enumerability of the set of enumerable chronological semimeasures.

5.11 Proof of the Entropy Inequality

We show¹³ that

$$\sum_{i=1}^n (y_i - x_i)^2 \leq \sum_{i=1}^n y_i \ln \frac{y_i}{x_i} \quad \text{with} \quad y_i \geq 0, \quad x_i > 0, \quad \sum_{i=1}^n y_i = 1, \quad \sum_{i=1}^n x_i = \alpha \leq 1$$

with $0 \ln 0 := 0$ or equivalently that

$$\sum_{i=1}^n f(x_i, y_i) \geq 0 \quad \text{with} \quad f(x, y) := y \ln \frac{y}{x} - (y - x)^2, \quad f : (0, 1] \times [0, 1] \rightarrow \mathbb{R} \quad (5.46)$$

¹²Semimeasures are, by definition, total functions.

¹³We will not explicate every subtlety and only sketch the proof.

The proof of the case $N=2$ will not be repeated here, as it is elementary and well known. We will reduce the general case $n > 2$ to the case $N=2$. It is enough to show that $\sum f \geq 0$ at all extremal points and “at” the boundary.

The boundary is the set of all (\mathbf{x}, \mathbf{y}) where one x_i or one y_i is / tends to zero. If one $y_i = 0$ we can reduce (5.46) to $n-1$ (with a different α) since $f(x_i, 0) = 0$. If one $x_i \rightarrow 0$ then $f(x_i, y_i) \rightarrow \infty$. As f is bounded from below ($f > -2$), $\sum f$ tends to infinity and (5.46) is satisfied. Hence (5.46) is satisfied “at” the boundary.

The extrema in the interior are found by differentiation. To include the boundary conditions we add Lagrange multipliers λ and μ

$$L(\mathbf{x}, \mathbf{y}) := \sum_{i=1}^n f(x_i, y_i) + \lambda \cdot \left(\alpha - \sum_{i=1}^n x_i \right) + \mu \cdot \left(1 - \sum_{i=1}^n y_i \right) \quad (5.47)$$

The extrema are at $\partial L / \partial x_i = \partial L / \partial y_i = 0$ i.e. at

$$\lambda = \partial_{x_i} f(x_i, y_i) \quad , \quad \mu = \partial_{y_i} f(x_i, y_i) \quad (5.48)$$

Assume the $(\mathbf{x}^*, \mathbf{y}^*)$ is a solution of (5.48). We can determine (λ^*, μ^*) for this solution by inserting e.g. the first component (x_1^*, y_1^*) into (5.48). But all other components of $(\mathbf{x}^*, \mathbf{y}^*)$ must be consistent with (5.48) too. Let us assume that for given (λ^*, μ^*) there are $m < \infty$ different solutions of (5.48), i.e. $(\mathbf{x}^*, \mathbf{y}^*)$ consists only of m different components $(\tilde{x}_k, \tilde{y}_k)$ with $1 \leq k \leq m$ where each component has multiplicity $n_k \geq 1$. Define $\bar{x}_k := n_k \tilde{x}_k$ and $\bar{y}_k := n_k \tilde{y}_k$. We have

$$\sum_{k=1}^m n_k = n \quad , \quad \sum_{k=1}^m \bar{y}_k = \sum_{k=1}^m n_k \tilde{y}_k = \sum_{i=1}^n y_i = 1 \quad , \quad \sum_{k=1}^m \bar{x}_k = \alpha$$

Equal components in (5.46) can be grouped together

$$\begin{aligned} \sum_{i=1}^n f(x_i^*, y_i^*) &= \sum_{k=1}^m n_k \tilde{y}_k \left[\ln \frac{\tilde{y}_k}{\tilde{x}_k} - 2(\tilde{y}_k - \tilde{x}_k)^2 \right] \geq \\ &\geq \sum_{k=1}^m n_k \tilde{y}_k \left[\ln \frac{n_k \tilde{y}_k}{n_k \tilde{x}_k} - 2n_k^2 (\tilde{y}_k - \tilde{x}_k)^2 \right] = \sum_{k=1}^m \bar{y}_k \left[\ln \frac{\bar{y}_k}{\bar{x}_k} - 2(\bar{y}_k - \bar{x}_k)^2 \right] = \\ &= \sum_{k=1}^m f(\bar{x}_k, \bar{y}_k) \stackrel{??}{\geq} 0 \end{aligned}$$

So we have reduced the case $\sum_{i=1}^n$ in a somewhat unconventional way to the case $\sum_{k=1}^m$ with m the number of solutions of (5.48). We will show that $R(x, y) := f(x, y) - \lambda x - \mu y$ has at most two (non-degenerate) extrema which will in turn proof that $\partial_x f = \lambda$, $\partial_y f = \mu$ has at most two solutions. Let us consider R on a curve connecting two extrema $g(t) := R(x(t), y(t))$, $x(0) = \tilde{x}_k$, $x(1) = \tilde{x}_l$, $y(0) = \tilde{y}_k$, $y(1) = \tilde{y}_l$, $0 \leq t \leq 1$. From $g'(0) = g'(1) = 0$ we know that there is a t_0 with $g''(t_0) = 0$. I.e. every connecting curve between two extrema contains a point in which R has curvature zero in one

direction and hence zero Gauss curvature G . The support of R is divided by the curve(s) $G(x,y)=0$ into zones. Each zone can contain at most one extremum.

$$\begin{aligned} G(x,y) &:= \det \begin{pmatrix} \partial_x^2 R & \partial_x \partial_y R \\ \partial_y \partial_x R & \partial_y^2 R \end{pmatrix} = (\partial_x^2 f)(\partial_y^2 f) - (\partial_x \partial_y f)^2 = \\ &= \left(\frac{y}{x^2} - 2\right)\left(\frac{1}{y} - 2\right) - \left(-\frac{1}{x} + 2\right)^2 = -\frac{2}{x^2 y}(x-y)^2 \end{aligned}$$

This is zero for $x=y$ only. The support of f is divided into two zones ($x < y$ and $x > y$). The (infinitely many) degenerate extrema for $\tilde{x}_k = \tilde{y}_k$ give no contribution to (5.46) ($f(x,x)=0$) and are hence irrelevant. So there are at most 2 non-trivial solutions of $\partial_x f = \lambda$, $\partial_y f = \mu$, hence $m \leq 2$.

In summary we have reduced (5.46) for general n to $m=2$ which is true. \square

5.12 History & References

Most references relevant to this chapter have already been given in previous chapters or in the main text of this chapter. Below we only remark on and give references to two further, in the context of AI ξ interesting, topics: protocols in probability theory and Bandit problems.

Paradoxes, sample spaces, protocols, and incompleteness. There are many paradoxes in probability theory, like the Petersburg, Monty Hall, Simpson, Newcomb, rich Uncle, 3 prisoners, etc. paradoxes [Szé86, EF98, Res01, Mos65]. Some of them are not particularly related to probability theory, but just to the improper use of math in general. Probably the most interesting paradoxes directly related to probability theory concern the awareness and choice of sample spaces and protocols (see [GH02] and references therein, especially [Sha85]). If one phrases these paradoxes within the AI μ model, one automatically has to be aware of and choose a suitable sample space and protocol. If this procedure uniquely determines μ , the paradox is solved. If not, the problem description was not complete, i.e. the description is consistent with a whole set \mathcal{M} of possible environments. This incompleteness can sometimes be overcome by symmetry or maximum entropy arguments (see [GH02] and Section 2.3.4). In general, a universal prior $\xi_{\mathcal{M}}$ and the predictions/actions of the SP $\xi_{\mathcal{M}}$ /AI $\xi_{\mathcal{M}}$ model represent a satisfactory solution to the paradox, solving the sample space, protocol, *and* incompleteness problem.

Bandit Problems. Bandit problems arose historically from the desire to optimally assign treatments to patients. They are prototypical problems for the so-called exploration versus exploitation dilemma. They were originally introduced by Robbins [Rob52]. One out of several arms (treatments) can be chosen leading to a possible reward (success). The goal is to maximize ones reward in repeated trials. The simplest model is to assume that arm i leads to reward 1 (0) with unknown probability p_i ($1-p_i$), where the probabilities are unknown. The traditional Bayesian solution

to the uncertainty about p_i is to assume a (second order Dirichlet) prior over p_i . The goal is to maximize the (exponentially) discounted reward sum. A closed solution can be given in terms of Gittins indices [GJ74, Git89]. For the regular discount sequences these strategies are not self-optimizing [BF85, KV86]. Many efficient heuristic self-optimizing approaches exist. In a minimax approach one tries to find strategies which led to highest expected reward in the worst case over unknown chances p_i [Vog60]. A complete worst case approach without any probabilistic assumption on the environment can be found in [ACBFS95]. The default textbook on Bandits is [BF85] and on Gittins indices [Git89].

5.13 Problems

5.1 (Actions as random variables) [C35oi] Instead of defining $\xi^{AI}(\underline{y}_{1:n})$ as a universal distribution over observations $x_{1:n}$ conditioned under actions $y_{1:n}$ as in (5.2) one may think of the following alternative definition: We use a universal distribution over observations *and* actions and then conditionalize to the actions, i.e. $\xi_{alt}^{AI}(\underline{y}_{1:n}) := M(\underline{y}_{1:n}) / \sum_{x_{1:n}} M(\underline{y}_{1:n})$, where M is Solomonoff's prior (2.19) (we could use ξ_U as well). One motivation for doing so is to regard M as a prior belief in the whole arrangement of agent+environment. The major problem with this approach is that ξ_{alt}^{AI} is not enumerable. More precisely, the presented definition does not lead to an enumeration procedure for ξ_{alt}^{AI} . This does not necessarily imply non-enumerability of ξ_{alt}^{AI} . Whether $\xi_{alt}^{AI} \stackrel{\times}{=} \xi^{AI}$ is also an open problem (cf. Problem 2.6). If true it would imply universality of ξ_{alt}^{AI} and convergence to computable μ^{AI} . This alternative approach also allows conditionalization w.r.t. the observations and to determine $M(\underline{y}_{<k}\underline{y}_k)$, which may be interpreted as the agent's own belief in selecting action y_k . But an actual action selection based on this probability would lead to a poorly performing agent, which differs from the "optimal" action y_k^ξ and $y_k^{\xi_{alt}}$ via the expectimax expression. Could $M(\underline{y}_{<k}\underline{y}_k)$ nevertheless be close to the action of p^ξ and/or $p^{\xi_{alt}}$ for large k , justifying the above interpretation of M ? (cf. Section 8.5.2 on multi-agent systems). Prove or disprove the stated open questions, conjectures, and assertions.

5.2 (Absorbing two-state environment) [C15ui] The HeavenHell example of Section 5.3.2 was a 3-state MDP which did not allow for self-optimizing policies. Here, a similar two-state MDP shall be analyzed in more detail. Let $\mathcal{M} = \{\mu_0, \mu_1, \dots\}$ with $w_0 = w_1 = \frac{1}{2}(1-\beta) > 0$, $\sum_{i \geq 2} w_i = \beta \geq 0$, $r_k \in \mathcal{R} = \{0, \frac{2}{3}, 1\}$, and $\mathcal{Y} = \{a, b\}$. Environments μ_1 and μ_2 are deterministic MDPs defined as $\mu_i = [\overset{y=a}{r=\frac{2}{3}} \text{---} \textcircled{s} \xrightarrow{y=b} \textcircled{e} \xrightarrow{y=*}{r=i}]$, i.e. initially being in state s , action b irrevocably leads to state e . The reward in state s is $\frac{2}{3}$. Environments μ_0 and μ_1 differ only in the reward in state e , which is $r=0$ in environment μ_0 , and $r=1$ in environment μ_1 . Show that there is no policy, self-optimizing in μ_1 and μ_2 . Now, consider the case $\beta=0$: Determine $V_{\mu_i}^p$ for $i \in \{0, 1\}$ and V_ξ^p for all policies p . The results only depend on the first time action

b is taken (if at all) and on the farsightedness m . Now determine $V_{\mu_i}^*$, p^ξ , and $V_{\mu_i}^{p^\xi}$. The results show that p^ξ is not self-optimizing ($V_{\mu_2}^{p^\xi} \not\rightarrow V_{\mu_2}^*$). Generalize the latter result to $0 < \beta < \frac{1}{6}(1 - \frac{1}{m})$ with environments μ_i , $i \geq 2$ defined arbitrarily, to V_{km}^* , and to $V_{k\gamma}^*$.

5.3 (Computing ρ) [C35u/C25u] Show that for every enumerable chronological semimeasure ρ there exists a Turing machine T of length $K(\rho)$ which computes it, i.e.

$$\rho(\underline{y}_{1:n}) = \sum_{q: T(qy_{1:n})=x_{1:n}} 2^{-l(q)} \quad \text{and} \quad l(T) \pm K(\rho).$$

(see (5.7) for context). The easy way is to adapt Lemma 4.3.4 of [LV97, p255].

5.4 (Pareto-Optimality) [C30u] We have shown Pareto-optimality of $\text{AI}\xi$ with ξ given in the form $\xi = \sum_\nu w_\nu \nu$. Show that $\text{AI}\xi$ with ξ defined as in (5.2) is also Pareto-optimal. Compare with or use the results of Problem 2.2.

5.5 (Pareto-Optimality) [C30oi] We define policy p to be equivalent to policy p' if both policies lead to the same value V_ν in all environments ν , i.e. if $V_\nu^p = V_\nu^{p'} \forall \nu \in \mathcal{M}$. Are all Pareto-optimal policies equivalent to some (Pareto-optimal) mixture policy p^ξ for certain weights w_ν ? A positive answer to this question implies that one can restrict the search of optimal policies to mixture policies. Try to find necessary and/or sufficient conditions which make the above question true.

5.6 (Convergence of averages) [C20u] $\delta_\nu(m) = O(f(m)) \forall \nu$ does not necessarily imply $\delta(m) = O(f(m))$ if \mathcal{M} is infinite, where $\delta(m) := \sum_{\nu \in \mathcal{M}} w_\nu \delta_\nu(m)$ and $\sum_{\nu \in \mathcal{M}} w_\nu \leq 1$ (see Lemma 5.28(iv)). For instance, for $\mathcal{M} \cong \mathbb{N}$, $\delta_\nu(m) := e^{-m/\nu} \leq 1$ decays exponentially in m for every ν , but for $w_\nu := \frac{1}{\nu(\nu+1)}$, $\delta(m)$ decays only harmonically. Show that $\delta(m) \geq \frac{1}{2m}(1 - e^{-m}) \geq \frac{1}{4m}$ for $m \geq 1$ (easy) and even $\delta(m) \geq \frac{1}{2m}$ (harder). Furthermore, show that the boundedness assumption in Lemma 5.28(ii) is necessary.

5.7 (Domain of definitions) [C20u] Several subtleties concerning the domain of definition and existence have been ignored. First, Definition 5.19 defined p^ρ only on histories produced by p^ρ itself. Given a history $y_{<k} \neq \dot{y}_{<k}$ one has to generalize the definition similarly to 5.30. Even in this generalized form p^ρ is only defined for histories which occur with non-zero ρ probability. Show that p^μ and p^ξ are defined for all histories which have non-zero μ probability. Use this to verify the soundness of the definitions and theorems in this chapter.

5.8 (Self-optimizing policies for geometric discounting) [C30ui] Show that there are self-optimizing policies in ergodic MDPs even for geometric discounting ($\frac{\gamma_{k+1}}{\gamma_k} \not\rightarrow 1$). On the other hand, the Bayes-mix policy p^ξ is not self-optimizing for Bandit problems with geometric discounting. Since Bandits are special ergodic MDPs these results seem to contradict Theorem 5.34. Clarify this paradox and discuss the implications. Hint: Histories $y_{<k}$ are policy dependent.

5.9 (Relevant and non-computable environments μ) [C30oi] Assume feedback x consists of three parts $x = x'r'x''$, future I/O is completely independent of x'' , $x'r'$ is sampled from a computable distribution, x'' from a (possibly) non-computable distribution. Show that ξ multiplicatively dominates μ' , where μ' is the true distribution μ modified in a way such that $V_\mu = V_{\mu'}$ but μ' is computable. This shows that the computability assumption on μ can be weakened to the (for AI ξ) *relevant* parts of the environment. Formulate and proof all this rigorously and generalize it to less trivial cases, where the relevant computable and the irrelevant non-computable information in x cannot be factored so easily.

5.10 (Self-optimizing environments) [C35u] Ergodic MDPs admit self-optimizing policies, which implies that p_{MDP1}^ξ is self-optimizing (see Section 5.6). Show that Bandits, i.i.d. processes, and classification tasks, are special (degenerate) cases of ergodic MDPs. The existence of self-optimizing policies is not limited to (subclasses of ergodic) MDPs. Suitably define ergodic partially observable MDPs (ergodic POMDPs) and k^{th} order ergodic MDPs and show that these classes also admit self-optimizing policies. Furthermore, show that factorizable environments, defined in Section 4.3.1, admit self-optimizing policies.

5.11 (Belief contamination) [C30ui] Consider an environmental class \mathcal{M} which admits self-optimizing policies. Theorem 5.34 shows that p^ξ is self-optimizing in the sense of $\lim_{k \rightarrow \infty} [V_{k\gamma}^{*\nu} - V_{k\gamma}^{p^\xi \nu}] = 0$. The Bayes-mixture $\xi := \sum_{\nu \in \mathcal{M}} w_\nu \nu$ expresses the degree of belief w_ν in environment $\nu \in \mathcal{M}$. We want to study the effect of additionally believing in some $\rho \notin \mathcal{M}$ with some small probability α . The new belief prior is $\xi' := (1-\alpha)\xi + \alpha\rho$. Show that a belief α in ρ much smaller than the belief w_μ in the true environment $\mu \in \mathcal{M}$ causes only a small corruption of the self-optimizing property. More precisely, $\limsup_{k \rightarrow \infty} [V_{k\gamma}^{*\mu} - V_{k\gamma}^{p^{\xi'} \mu}] \leq \frac{\alpha \cdot r_{max}}{(1-\alpha)w_\mu}$. Construct examples for which $V_{k\gamma}^{p^\xi \mu} - V_{k\gamma}^{p^{\xi'} \mu} = \frac{\alpha \cdot r_{max}}{(1-\alpha)w_\mu}$. This shows that the upper bound cannot be improved in general and that a belief contamination α of magnitude comparable to w_μ can completely degrade performance.

5.12 (Continuity of Value for ergodic MDPs) [C40usm] Let μ and $\hat{\mu}$ be MDPs, which are “close” to each other in the sense that $\mu(as_{<k}as_k) = T_{s_{k-1}s_k}^{a_k}$ and $\hat{\mu}(as_{<k}as_k) = \hat{T}_{s_{k-1}s_k}^{a_k}$ and $\varepsilon := \max_{ss'a} |T_{ss'}^a - \hat{T}_{ss'}^a|$ is “small”. Furthermore let p be a stationary policy, i.e. $p(s_{<k}) = p(s_{k-1})$. Show properties (i)–(vii) of the value function(s)

- o) Condition: T is ergodic –or– $T_{ss'}^a = 0$ implies $\hat{T}_{ss'}^a = 0$.
- i) $v_T^p := \lim_{m \rightarrow \infty} V_{1m}^{pT}$ exists.
- ii) $|v_T^p - v_{\hat{T}}^p| = O(\varepsilon)$ if (o).
- iii) $p^T := \arg\max_p v_T^p$ can be chosen stationary. $v_T^* := \max_p v_T^p = v_T^{p^T}$.
- iv) $|v_T^* - v_{\hat{T}}^*| = O(\varepsilon)$.
- v) $|v_T^p - \frac{1}{m} V_{1m}^{pT}| = O(\frac{1}{m})$.
- vi) $|\frac{1}{m} V_{1m}^{p^T T} - \frac{1}{m} V_{1m}^{*T}| = O(\frac{1}{m})$ if (o).

vii) $|\frac{1}{m}V_{1m}^{p^T T} - \frac{1}{m}V_{1m}^{*T}| = 3 \cdot O(\frac{1}{m}) + 2 \cdot O(\varepsilon)$ if (o).

The “constant” factor hidden in $O()$ depends on T , but is independent of m , ε , and \hat{T} . Note that $\operatorname{argmax}_p V_{1m}^{*T}$ may be non-stationary. (vii) with $k_0 m$ instead of $1m$ and $\varepsilon \sim k_0^{-1/2} \sim m^{-1/3}$ has been used in the proof of Theorem 5.38(i).

Hints: (i) follows from the existence of $\bar{T} := \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} (T)^k$ [Ber95b, p187]. (ii) For ergodic T all rows of \bar{T} coincide with the stationary distribution, which is proportional to some column of the adjoint matrix of T , which itself is a polynomial in (the components) of T . (iii) similar to [Ber95b, p191]. (iv) from (ii) and (iii). (v) similar (i). (vi) similar (iii). For (vii) chain $(v) + (ii) + (iv) + (v) + (vi)$ and use the triangle inequality.

5.13 (Ergodic versus Forgetful Environments) [C20s/C10u] Forgetful environments have been defined in Section 5.3.6 as being asymptotically independent of the history. (Ergodic) MDPs have been defined in Section 5.6 as visiting every state infinitely often. An environment is called acyclic if the probability of infinitely repeating cycles is zero. Show that every acyclic ergodic MDP is forgetful, but not every forgetful MDP is ergodic. Note also that forgetfulness is a broader concept than (k^{th} order) MDPs.

5.14 (Uniform mixture of MDPs) [C30usi] In the following you are asked to derive explicit expressions for ξ^{MDP} for uniform prior belief w . Let $\mu_T \in \mathcal{M}_{MDP}$ be a (completely observable) MDP with transition matrix T . $T_{ss'}^a$ is the probability of going to state $s' \in \mathcal{X}$ under action $a \in \mathcal{Y}$ if currently in state $s \in \mathcal{X}$. Given a policy which determines the actions a_t , the probability of action-observation history $a_1 s_1 \dots a_n s_n$ after cycle n is $\mu_T(a_1 s_1 \dots a_n s_n) = T_{s_0 s_1}^{a_1} \dots T_{s_{n-1} s_n}^{a_n}$, where s_0 is some initial state (randomization over the initial state may be performed). Reward is a (given) function of state ($r_k = r(s_k)$). Optimal policies/actions follow from the recursive Bellman equations (4.29) or the explicit expectimax expression of the AI μ model (4.17). Assume now that we only know that the true environment is an MDP, but nothing more, i.e. $T \in \mathcal{T} := \{T : T_{ss'}^a \geq 0, \sum_{s'} T_{ss'}^a = 1\}$ is unknown. Since \mathcal{T} is continuous, the Bayes-mixture ξ has the form $\xi(\underline{as}_{1:n}) := \int_{\mathcal{T}} w_T \mu_T(\underline{as}_{1:n}) dT$.

i) Assume a uniform prior belief over T , i.e. $w_T \propto 1$ and the measure dT is the uniform measure on the polytope \mathcal{T} . Compute the integral and show that the ratio $\xi(as_{<n} \underline{as}_n) = \xi(\underline{as}_{1:n}) / \xi(\underline{as}_{<n}) = N_{s_{n-1} s_n}^{a_n} / (\sum_{s'} N_{s_{n-1} s'}^{a_n} + S - 1)$, where $S = |\mathcal{X}|$ is the number of states and $N_{ss'}^a$ is the historical number of transitions from s to s' under action a , (including the transition from s_0 if $s = s_0$ and to s_n if $s' = s_n$). This is just Laplace’s law of succession [Lap1814], one for each (ass') -tupel. For instance, initially all transitions are equally plausible $\xi(a_1 s_1) = \frac{1}{S}$.

ii) Show that, although the class \mathcal{T} is continuous and contains non-ergodic environments, the Bayes-optimal policy p^ξ is self-optimizing for ergodic environments $\mu_T \in \mathcal{M}_{MDP1}$. The intuitive reason is that \mathcal{T} is compact and the non-ergodic environments have measure zero.

iii) Model based reinforcement learning algorithms try to estimate T from past experience. Give an expression for the posterior believe $w_T(as_{1:n}) \propto \mu_T(as_{1:n})$ in transition T . Note that this is a (complex) distribution over T , while most reinforcement learning algorithms only estimate a single (e.g. a most likely) T . Show that the expected transition probability $\mathbf{E}[T_{ss'}^a | as_{1:n}] := \int_{\mathcal{T}} T_{ss'}^a w_T(as_{1:n}) dT = (N_{ss'}^a + 1) / (\sum_{s'} N_{ss'}^a + S)$ coincides with the relative historic occurrence of (ass') . Show that policy p^ξ based on (4.17) appropriately explores the environment, while the popular policy based on $\mathbf{E}[T]$ or other point estimates like Maximum Likelihood lack exploration.

iv) Assume we know that the environment is a *deterministic* MDP, i.e. $\mathcal{T} = \{T : T_{ss'}^a \in \{0,1\}, \sum_{s'} T_{ss'}^a = 1\}$. Repeat (i)–(iii) with this (now discrete) \mathcal{T} . Is the corresponding p^ξ self-optimizing?

v) Assume now that for every action a there exists a mirror “undo” action \bar{a} in the sense that $T_{ss'}^a = T_{s's}^{\bar{a}}$. Repeat (i)–(iii) for the set of all deterministic “symmetric” MDPs. Is the corresponding p^ξ self-optimizing? An example is a robot moving in a (noiseless) environment, like a maze. Non-symmetric MDPs contain one-way streets or doors, which are missing in symmetric MDPs.

vi) Incorporate further knowledge of the form $T_{ss'}^a = 0/1$ for some (ass') and repeat (i)–(iii). For example, if we know that the environment is an $l \times l$ grid-maze, and transitions are a priori only possible between neighboring cells, we know that $T_{ss'}^a = 0$ if s and s' are not neighboring grid cells.

vii) Explore the difficulties when extending the considerations in (i)–(iii) to POMDPs, potentially with variable state space size S , e.g. with prior $w_S \propto S^{-2}$.

5.15 (Effective horizons) [C25u] Derive the expressions for effective horizons presented in Table 5.41.



Richard Sutton

*“Ideas matter”
 “Approximate the solution, not the problem”
 (Richard Sutton)*

Chapter 6

Important Environmental Classes

6.1	Repetition of the AI_{μ}/ξ Models	602
6.2	Sequence Prediction (SP)	603
6.2.1	Using the AI_{μ} Model for Sequence Prediction	604
6.2.2	Using the AI_{ξ} Model for Sequence Prediction	606
6.3	Strategic Games (SG)	608
6.3.1	Introduction	608
6.3.2	Strictly Competitive Strategic Games	608
6.3.3	Using the AI_{μ} Model for Game Playing	609
6.3.4	Games of Variable Length	610
6.3.5	Using the AI_{ξ} Model for Game Playing	611
6.4	Function Minimization (FM)	612
6.4.1	Applications/Examples	612
6.4.2	The Greedy Model FMG_{μ}	613
6.4.3	The General FM_{μ}/ξ Model	615
6.4.4	Is the General Model Inventive?	616
6.4.5	Using the AI models for Function Minimization	617
6.4.6	Remark	618
6.5	Supervised Learning from Examples (EX)	619
6.5.1	Applications/Examples	619
6.5.2	Supervised Learning with the AI_{μ}/ξ Model	620
6.6	Other Aspects of Intelligence	621
6.7	Problems	622

In this and the following chapter we define $\xi = \xi_U \stackrel{\times}{=} M$ to be Solomonoff’s prior, i.e. $AI_{\xi} = AIXI$. In order to give further support for the universality and optimality of the AI_{ξ} theory, we apply AI_{ξ} in this chapter to a number of problem classes. They include sequence prediction, strategic games, function minimization and, especially,

how $\text{AI}\xi$ learns to learn supervised. For some classes we give concrete examples to illuminate the scope of the problem class. We first formulate each problem class in its natural way (when μ^{problem} is known) and then construct a formulation within the $\text{AI}\mu$ model and prove its equivalence. We then consider the consequences of replacing μ by ξ . The main goal is to understand why and how the problems are solved by $\text{AI}\xi$. We only highlight special aspects of each problem class. Sections 6.2–6.6 together should give a better picture of the $\text{AI}\xi$ model. We do not study every aspect for every problem class. The sections might be read selectively. They are not essential to understand the remaining chapters.

6.1 Repetition of the $\text{AI}\mu/\xi$ Models

In the last chapter we unified sequential decision theory with the theory of universal induction to a model of Artificial Intelligence, which we claimed to be universal and superior to any other model in various senses. All tasks which require intelligence to be solved can naturally be formulated as a maximization of some expected utility in the framework of agents. The main remaining problem is the unknown prior probability distribution μ^{AI} of the environment(s). Conventional learning algorithms are restricted in the sense that they can neither handle large (unstructured) state spaces, nor do they converge in the theoretically minimal number of cycles, nor can they handle non-stationary environments appropriately. On the other hand, the universal semimeasure ξ (2.24), based on ideas from algorithmic information theory, solves the problem of the unknown prior distribution for induction problems. No explicit learning procedure is necessary, as ξ automatically converges to μ . We unified the theory of universal sequence prediction with the decision theoretic agent by replacing the unknown true prior μ^{AI} by an appropriately generalized universal semimeasure ξ^{AI} . For convenience we repeat some definitions and results from previous chapters which we need in this chapter.

Let $\mu(y_{<k}\underline{y}_k)$ be the true probability of input x_k in cycle k , given the history $y_{<k}y_k$. $\mu(\underline{y}_{1:k})$ is the true chronological prior probability that the environment reacts with $x_{1:k}$ if provided with actions $y_{1:k}$ from the agent. We define $V_{k+1,m}^{*\mu}(y_{1:k})$ to be the μ -expected reward sum in cycles $k+1$ to m with outputs y_i generated by agent p^* and responses x_i from the environment. Adding reward $r_k \equiv r(x_k)$ to $V_{k+1,m}^{*\mu}$ we get the reward including cycle k . The probability of x_k , given $y_{<k}y_k$, is given by the conditional probability $\mu(y_{<k}\underline{y}_k)$. p^* chooses y_k as to maximize the future reward. So the expected reward sum in cycles k to m given $y_{<k}$ and y_k chosen by p^* is

$$V_{km}^{*\mu}(y_{<k}) = \max_{y_k} \sum_{x_k} [r_k + V_{k+1,m}^{*\mu}(y_{1:k})] \cdot \mu(y_{<k}\underline{y}_k) \quad (6.1)$$

Together with the induction start

$$V_{m+1,m}^{*\mu}(y_{1:m}) := 0 \quad (6.2)$$

$V_{km}^{*\mu}$ is completely defined. If m_k is our horizon function of p^* and $\dot{y}_{<k}$ is the actual history in cycle k , the output \dot{y}_k of the agent is given by

$$\dot{y}_k = \arg \max_{y_k} \sum_{x_k} [r_k + V_{k+1, m_k}^{*\mu}(\dot{y}_{<k} \underline{y}_k)] \cdot \mu(\dot{y}_{<k} \underline{y}_k) \quad (6.3)$$

which in turn defines the policy p^* . Then the environment responds \dot{x}_k with probability $\mu(\dot{y}_{<k} \dot{x}_k)$. Then cycle $k+1$ starts. We might unfold the recursion (6.1) further and give \dot{y}_k explicitly as

$$\dot{y}_k = \arg \max_{y_k} \sum_{x_k} \max_{y_{k+1}} \sum_{x_{k+1}} \dots \max_{y_{m_k}} \sum_{x_{m_k}} (r_k + \dots + r_{m_k}) \cdot \mu(\dot{y}_{<k} \underline{y}_{k:m_k}). \quad (6.4)$$

This expression has a direct interpretation: The probability of inputs $x_{k:m_k}$ in cycle k when the agent outputs $y_{k:m_k}$ with actual history $\dot{y}_{<k}$ is $\mu(\dot{y}_{<k} \underline{y}_{k:m_k})$. The future reward in this case is $r_k + \dots + r_{m_k}$. The best expected reward is obtained by averaging over the x_i ($\sum x_i$) and maximizing over the y_i . This has to be done in chronological order to correctly incorporate the dependency of x_i and y_i on the history. This is essentially the expectimax algorithm/sequence/tree (see Figure 4.13). The $\text{AI}\mu$ model is *optimal* in the sense that no other policy leads to higher expected reward. Unfortunately, in AI, the environment μ is often unknown. The $\text{AI}\xi$ model is defined similarly to the $\text{AI}\mu$ model, but with the unknown μ replaced by the (known) universal prior ξ :

$$\dot{y}_k = \arg \max_{y_k} \sum_{x_k} \max_{y_{k+1}} \sum_{x_{k+1}} \dots \max_{y_{m_k}} \sum_{x_{m_k}} (r_k + \dots + r_{m_k}) \cdot \xi(\dot{y}_{<k} \underline{y}_{k:m_k}) \quad (6.5)$$

$$\text{with } \xi(\underline{y}_{1:k}) := \sum_{q: q(y_{1:k}) = x_{1:k}} 2^{-l(q)}. \quad (6.6)$$

Motivations for $\text{AI}\xi$ being a good substitute for $\text{AI}\mu$ were the convergence of ξ to μ

$$\xi(\underline{y}_{<k} \underline{y}_{k:m_k}) \xrightarrow{k \rightarrow \infty} \mu(\underline{y}_{<k} \underline{y}_{k:m_k}) \quad \text{if } h_k \equiv m_k - k + 1 \leq h_{\max} < \infty \quad (6.7)$$

and the tight error and loss bounds in the case a passive sequence prediction.

6.2 Sequence Prediction (SP)

We have introduced the $\text{AI}\xi$ model as a unification of ideas of sequential decision theory and universal probability distribution. We might expect $\text{AI}\xi$ to behave identically to $\text{SP}\xi$, when faced with a sequence prediction problem, but things are not that simple, as we will see. Let us repeat the definition of the total number of expected erroneous predictions the $\text{SP}\rho$ agent makes for the first n observations:

$$E_n^{\Theta_\rho} := \sum_{k=1}^n \sum_{x_{<k}} \mu(x_{<k}) [1 - \mu(x_{<k} \underline{x}_k^{\Theta_\rho})] \quad \text{with } x_k^{\Theta_\rho} := \arg \max_{x_k} \rho(x_{<k} \underline{x}_k). \quad (6.8)$$

The $\text{SP}\mu$ agent is best in the sense that $E_n^{\Theta_\mu} \leq E_n^{\Theta_\rho}$ for *any* ρ . We have shown that the universal predictor $\text{SP}\xi$ is not much worse

$$E_n^{\Theta_\xi} - E_n^{\Theta_\rho} \leq 2D + 2\sqrt{E_n^{\Theta_\rho} D} = O(\sqrt{E_n^{\Theta_\rho}}), \quad D \stackrel{+}{\leq} \ln 2 \cdot K(\mu). \quad (6.9)$$

6.2.1 Using the AI μ Model for Sequence Prediction

We have seen in Chapter 3 how to predict sequences for known and unknown prior distribution μ^{SP} . Here we consider binary sequences¹ $z_1 z_2 z_3 \dots \in \mathcal{B}^\infty$ with known prior probability $\mu^{SP}(\underline{z_1 z_2 z_3 \dots})$.

We want to show how the AI μ model can be used for sequence prediction. We will see that it makes the same prediction as the SP μ agent. For simplicity we only discuss the special error-loss $\ell_{xy} = 1 - \delta_{xy}$, where δ is the Kronecker symbol, defined as $\delta_{ab} = 1$ for $a = b$ and 0 otherwise. First, we have to specify *how* the AI μ model should be used for sequence prediction. The following choice is natural:

The system's output y_k is interpreted as a prediction for the k^{th} bit z_k of the string under consideration. This means that y_k is binary ($y_k \in \mathcal{B} =: \mathcal{Y}$). As a reaction of the environment, the agent receives reward $r_k = 1$ if the prediction was correct ($y_k = z_k$), or $r_k = 0$ if the prediction was erroneous ($y_k \neq z_k$). The question is what the input x'_k in the next cycle should be. One choice would be to inform the agent about the correct k^{th} bit of the string and set $x'_k = z_k$. But as from the reward r_k in conjunction with the prediction y_k , the true bit $z_k = \delta_{y_k r_k}$ can be inferred, this information is redundant. There is no need for this additional feedback. So we set $x'_k = \epsilon \in \mathcal{X} = \{\epsilon\}$ thus having $x_k \equiv r_k$. The agent's performance does not change when we include this redundant information, it merely complicates the notation. The prior probability μ^{AI} of the AI μ model is

$$\mu^{AI}(y_1 \underline{x_1} \dots y_k \underline{x_k}) = \mu^{AI}(y_1 \underline{r_1} \dots y_k \underline{r_k}) = \mu^{SP}(\underline{\delta_{y_1 r_1} \dots \delta_{y_k r_k}}) = \mu^{SP}(\underline{z_1 \dots z_k}) \quad (6.10)$$

In the following, we will drop the superscripts of μ because they are clear from the arguments of μ and the μ equal in any case. The formula (6.1) for the expected reward reduces to

$$V_{km}^{*\mu}(y_{<k}) = \max_{y_k} \sum_{r_k} [r_k + V_{k+1,m}^{*\mu}(y_{1:k})] \cdot \mu(\delta_{y_1 r_1} \dots \delta_{y_{k-1} r_{k-1}} \underline{\delta_{y_k r_k}}) \quad (6.11)$$

The first observation we can make, is that for this special μ , $V_{km}^{*\mu}$ only depends on $\delta_{y_i r_i}$, i.e. replacing y_i and r_i simultaneously with their complements does not change the value of $V_{km}^{*\mu}$. We have a symmetry in $y_i r_i$. For $k = m + 1$ this is definitely true as $V_{m+1,m}^{*\mu} = 0$ in this case (see (6.2)). For $k \leq m$ we prove it by induction. The r.h.s. of (6.11) is symmetric in $y_i r_i$ for $i < k$ because μ possesses this symmetry and $V_{k+1,m}^{*\mu}$ possesses it by induction hypothesis, so the symmetry holds for the l.h.s., which completes the proof. The prediction \dot{y}_k is

$$\begin{aligned} \dot{y}_k &= \arg \max_{y_k} \sum_{r_k} [r_k + V_{k+1,m}^{*\mu}(\dot{y}_{<k} y_k)] \cdot \mu(\delta_{y_1 r_1} \dots \delta_{y_{k-1} r_{k-1}} \underline{\delta_{y_k r_k}}) = \\ &= \arg \max_{y_k} \sum_{r_k} r_k \cdot \mu(\delta_{\dot{y}_1 r_1} \dots \underline{\delta_{y_k r_k}}) = \arg \max_{y_k} \mu(\dot{z}_1 \dots \dot{z}_{k-1} \underline{y_k}) = \arg \max_{z_k} \mu(\dot{z}_1 \dots \dot{z}_{k-1} \underline{z_k}) \end{aligned} \quad (6.12)$$

¹We use z_k to avoid notational conflicts with the agent's inputs x_k .

The first equation is the definition of the agent's action (6.3), where we have used (6.10), which gives the r.h.s. of (6.11) with \max_{y_k} replaced by $\operatorname{argmax}_{y_k}$. $\sum_r f(\dots \delta_{yr} \dots)$ is independent of y for any function, depending on the combination δ_{yr} only. Therefore, the $\sum_r V^* \mu$ term is independent of y_k because $V_{k+1,m}^{*\mu}$ as well as μ depend on $\delta_{y_k r_k}$ only. In the second equation, we can therefore drop this term, as adding a constant to the argument of $\operatorname{argmax}_{y_k}$ does not change the location of the maximum. In the second last equation we evaluated the \sum_{r_k} . Further, if the true reward to y_i is r_i the true i^{th} bit of the string must be $z_i = \delta_{y_i r_i}$. The last equation is just a renaming.

So, the $\text{AI}\mu$ model predicts that z_k that has maximal μ probability, given $z_1 \dots z_{k-1}$. This prediction is independent of the choice of m_k . It is exactly the prediction scheme of the sequence predictor $\text{SP}\mu$ with known prior described in Section 3.3. As this model was optimal, $\text{AI}\mu$ is optimal too, i.e. has minimal number of expected errors (maximal μ -expected reward) as compared to any other sequence prediction scheme.

From this, it is already clear that the value $V_{km}^{*\mu}$ must be closely related to the expected sequence prediction error $E_m^{\Theta\mu}$ (6.8). In the following we prove that $V_{1m}^{*\mu} = m - E_m^{\Theta\mu}$. We rewrite $V_{km}^{*\mu}$ in (6.11) as a function of z_i instead of $y_i r_i$ as it is symmetric in $y_i r_i$. Further, we can pull $V_{k+1,m}^{*\mu}$ out of the maximization, as it is independent of y_k similarly as in (6.12). Renaming the bounded variables y_k and r_k we get

$$V_{km}^{*\mu}(z_{<k}) = \max_{z_k} \mu(z_{<k} z_k) + \sum_{z_k} V_{k+1,m}^{*\mu}(z_{1:k}) \cdot \mu(z_{<k} z_k) \quad (6.13)$$

Recursively inserting the l.h.s. into the r.h.s. we get

$$V_{km}^{*\mu}(z_{<k}) = \sum_{i=k}^m \sum_{z_{k:i-1}} \max_{z_i} \mu(z_{<k} z_{k:i}) \quad (6.14)$$

This is most easily proven by induction. For $k = m$ we have $V_{mm}^{*\mu}(z_{<m}) = \max_{z_m} \mu(z_{<m} z_m)$ from (6.13) and (6.2), which equals (6.14). By induction hypothesis, we assume that (6.14) is true for $k+1$. Inserting this into (6.13) we get

$$\begin{aligned} V_{km}^{*\mu}(z_{<k}) &= \max_{z_k} \mu(z_{<k} z_k) + \sum_{z_k} \left[\sum_{i=k+1}^m \sum_{z_{k+1:i-1}} \max_{z_i} \mu(z_{1:k} z_{k+1:i}) \right] \mu(z_{<k} z_k) = \\ &= \max_{z_k} \mu(z_{<k} z_k) + \sum_{i=k+1}^m \sum_{z_{k:i-1}} \max_{z_i} \mu(z_{<k} z_{k:i}) \end{aligned}$$

which equals (6.14). This was the induction step and hence (6.14) is proven. By setting $k=1$ and slightly reformulating (6.14), we get the total expected reward in the first m cycles

$$V_{1m}^{*\mu}(\epsilon) = \sum_{i=1}^m \sum_{z_{<i}} \mu(z_{<i}) \max\{\mu(z_{<i} \underline{0}), \mu(z_{<i} \underline{1})\} = m - E_m^{\Theta\mu}$$

with $E_m^{\Theta\mu}$ defined in (6.8).

6.2.2 Using the AI ξ Model for Sequence Prediction

Now we want to use the universal AI ξ model instead of AI μ for sequence prediction and try to derive error bounds analogous to (6.9). Like in the AI μ case, the agent's output y_k in cycle k is interpreted as a prediction for the k^{th} bit z_k of the string under consideration. The reward is $r_k = \delta_{y_k z_k}$ and there are no other inputs $x_k = \epsilon$. What makes the analysis more difficult is that ξ is not symmetric in $y_i r_i \leftrightarrow (1 - y_i)(1 - r_i)$ and (6.10) does not hold for ξ . On the other hand, ξ^{AI} converges to μ^{AI} in the limit (6.7), and (6.10) should hold asymptotically for ξ in some sense. So we expect that everything proven for AI μ holds approximately for AI ξ . The AI ξ model should behave similarly to Solomonoff prediction SP ξ . Especially we expect error bounds similar to (6.9). Making this rigorous seems difficult. Some general remarks have been made in the last chapter. Note that bounds like (5.15) can't hold in general, but could be valid for AI ξ in (pseudo) passive environments.

Here we concentrate on the special case of a deterministic computable environment, i.e. the environment is a sequence $\dot{z} = \dot{z}_1 \dot{z}_2 \dots$, $Km(\dot{z}_1 \dots \dot{z}_n) \leq Km(\dot{z}) < \infty$. $Km(\dot{z}_{1:n})$ is the length of the shortest (possibly non-halting) program printing a string starting with $z_{1:n}$. Furthermore, we only consider the simplest horizon model $m_k = k$, i.e. greedily maximize only the next reward. This is sufficient for sequence prediction, as the reward of cycle k only depends on output y_k and not on earlier decisions. This choice is in no way sufficient and satisfactory for the full AI ξ model, as *one* single choice of m_k should serve for *all* AI problem classes. So AI ξ should allow good sequence prediction for some universal choice of m_k and not only for $m_k = k$, which definitely does not suffice for more complicated AI problems. The analysis of this general case is a challenge for the future. For $m_k = k$ the AI ξ model (6.5) with $x'_i = \epsilon$ reduces to

$$\dot{y}_k = \arg \max_{y_k} \sum_{r_k} r_k \cdot \xi(\dot{y}_{<k} y_k) = \arg \max_{y_k} \xi(\dot{y}_{<k} y_k \underline{1}) = \arg \max_{y_k} \xi(\dot{y}_{<k} y_k \underline{1}) \quad (6.15)$$

The environmental response \dot{r}_k is given by $\delta_{\dot{y}_k \dot{z}_k}$; it is 1 for a correct prediction ($\dot{y}_k = \dot{z}_k$) and 0 otherwise. In the following, we want to bound the number of errors this prediction scheme makes. We need the following inequality

$$\xi(\underline{y}_1 \dots \underline{y}_k) > 2^{-Km(\delta_{y_1 r_1} \dots \delta_{y_k r_k}) - O(1)} \quad (6.16)$$

We have to find a short program in the sum (6.6) calculating $r_1 \dots r_k$ from $y_1 \dots y_k$. If we knew $z_i := \delta_{y_i r_i}$ for $1 \leq i \leq k$ a program of size $O(1)$ could calculate $r_1 \dots r_k = \delta_{y_1 z_1} \dots \delta_{y_k z_k}$. So combining this program with a shortest coding of $z_1 \dots z_k$ leads to a program q of size $l(q) = Km(z_1 \dots z_k) + O(1)$ with $q(y_{1:k}) = r_{1:k}$, which proves (6.16).

Let us now assume that we make a wrong prediction in cycle k , i.e. $\dot{r}_k = 0$, $\dot{y}_k \neq \dot{z}_k$. The goal is to show that $\dot{\xi}$ defined by

$$\dot{\xi}_k := \xi(\dot{y}_{1:k}) = \xi(\dot{y}_{<k} \dot{y}_k \underline{0}) \leq \xi(\dot{y}_{<k}) - \xi(\dot{y}_{<k} \dot{y}_k \underline{1}) < \dot{\xi}_{k-1} - \alpha$$

decreases for every wrong prediction, at least by some α . The \leq arose from the fact that ξ is only a semimeasure.

$$\begin{aligned} \xi(\dot{y}_{<k}\dot{y}_k\mathbf{1}) &> \xi(\dot{y}_{<k}(1-\dot{y}_k)\mathbf{1}) \stackrel{\times}{>} 2^{-Km(\delta_{\dot{y}_1\dot{r}_1}\dots\delta_{(1-\dot{y}_k)\mathbf{1}})} = \\ &= 2^{-Km(\dot{z}_1\dots\dot{z}_k)} > 2^{-Km(\dot{z})-O(1)} =: \alpha \end{aligned}$$

In the first inequality we have used the fact that \dot{y}_k maximizes by definition (6.15) the argument, i.e. $1-\dot{y}_k$ has lower ξ probability than \dot{y}_k . Bound (6.16) has been applied in the second inequality. The equality holds, because $\dot{z}_i = \delta_{\dot{y}_i\dot{r}_i}$ and $\delta_{(1-\dot{y}_k)\mathbf{1}} = \delta_{\dot{y}_k\mathbf{0}} = \delta_{\dot{y}_k\dot{r}_k} = \dot{z}_k$. The last inequality follows from the definition of \dot{z} .

We have shown that each erroneous prediction reduces $\dot{\xi}$ by at least the α defined above. Together with $\dot{\xi}_0 = 1$ and $\dot{\xi}_k > 0$ for all k this shows that the agent can make at most $1/\alpha$ errors, since otherwise $\dot{\xi}_k$ would become negative. So the number of wrong predictions $E_{n\dot{\xi}}^{AI}$ of agent (6.15) is bounded by

$$E_{n\dot{\xi}}^{AI} < \frac{1}{\alpha} = 2^{Km(\dot{z})+O(1)} < \infty \quad (6.17)$$

for a computable deterministic environment string $\dot{z}_1\dot{z}_2\dots$. The intuitive interpretation is that each wrong prediction eliminates at least one program p of size $l(p) \stackrel{+}{\leq} Km(\dot{z})$. The size is smaller than $Km(\dot{z})$, as larger policies could not mislead the agent to a wrong prediction, since there is a program of size $Km(\dot{z})$ making a correct prediction. There are at most $2^{Km(\dot{z})+O(1)}$ such policies, which bounds the total number of errors.

We have derived a finite bound for $E_{n\dot{\xi}}^{AI}$, but unfortunately, a rather weak one as compared to (6.9). The reason for the strong bound in the SP case was that every error at least halves $\dot{\xi}$ because the sum of the argmax_{x_k} arguments was 1. Here we have

$$\begin{aligned} \xi(\dot{y}_1\dot{r}_1\dots\dot{y}_{k-1}\dot{r}_{k-1}\mathbf{00}) + \xi(\dot{y}_1\dot{r}_1\dots\dot{y}_{k-1}\dot{r}_{k-1}\mathbf{01}) &= 1 \\ \xi(\dot{y}_1\dot{r}_1\dots\dot{y}_{k-1}\dot{r}_{k-1}\mathbf{10}) + \xi(\dot{y}_1\dot{r}_1\dots\dot{y}_{k-1}\dot{r}_{k-1}\mathbf{11}) &= 1 \end{aligned}$$

but argmax_{y_k} runs over the right top and right bottom ξ , for which no sum criterion holds.

The AI ξ model would not be sufficient for realistic applications if the bound (6.17) were sharp, but we have the strong feeling (but only weak arguments) that better bounds proportional to $Km(\dot{z})$ analogous to (6.9) exist. The technique used above may not be appropriate for achieving this. One argument for a better bound is the formal similarity between $\text{argmax}_{z_k}\xi(\dot{z}_{<k}\dot{z}_k)$ and (6.15), the other is that we were unable to construct an example sequence for which (6.15) makes more than $O(Km(\dot{z}))$ errors (see Problem 6.2).

6.3 Strategic Games (SG)

6.3.1 Introduction

A very important class of problems are strategic games, like chess. In fact, what is subsumed under game theory, is so general, that it includes not only a huge variety of games, from simple games of chance like roulette, combined with strategy like backgammon, up to purely strategic games like chess or checkers or go. Game theory can also describe political and economic competitions and coalitions, Darwinism and many more. It seems that nearly every AI problem could be brought into the form of a game. Nevertheless, the intention of a game is that several players perform actions with (partial) observable consequences. The goal of each player is to maximize some utility function (e.g. to win the game). The players are assumed to be rational, taking into account all information they possess. The different goals of the players are usually in conflict. For an introduction into game theory, see [FT91, OR94, RN95, NM44].

If we interpret the $AI\mu$ model as one player and the environment models the other rational player *and* the environment provides the reinforcement feedback r_k , we see that the agent-environment configuration satisfies all criteria of a game. On the other hand, the AI models can handle more general situations, since it interacts optimally with an environment, even if the environment is not a rational player with conflicting goals.

6.3.2 Strictly Competitive Strategic Games

In the following, we restrict ourselves to deterministic, strictly competitive strategic² games with alternating moves. Player 1 makes move y'_k in round k , followed by the move x'_k of player 2. So a game with n rounds consists of a sequence of alternating moves $y'_1 x'_1 y'_2 x'_2 \dots y'_n x'_n$. At the end of the game in cycle n the game or final board situation is evaluated with $V(y'_1 x'_1 \dots y'_n x'_n)$. Player 1 tries to maximize V , whereas player 2 tries to minimize V . In the simplest case, V is 1 if player 1 won the game, $V = -1$ if player 2 won and $V = 0$ for a draw. We assume a fixed game length n independent of the actual move sequence. For games with variable length but maximal possible number of moves n , we could add dummy moves and pad the length to n . The optimal strategy (Nash equilibrium) of both players is a minimax strategy

$$\dot{x}'_k = \arg \min_{x'_k} \max_{y'_{k+1}} \min_{x'_{k+1}} \dots \max_{y'_n} \min_{x'_n} V(\dot{y}'_1 \dot{x}'_1 \dots \dot{y}'_k \dot{x}'_k \dots y'_n x'_n) \quad (6.18)$$

$$\dot{y}'_k = \arg \max_{y'_k} \min_{x'_k} \dots \max_{y'_n} \min_{x'_n} V(\dot{y}'_1 \dot{x}'_1 \dots \dot{y}'_{k-1} \dot{x}'_{k-1} \dot{y}'_k \dot{x}'_k \dots y'_n x'_n) \quad (6.19)$$

But note, that the minimax strategy is only optimal if both players behave rationally. If, for instance, player 2 has limited capabilities or makes errors and player 1 is able to discover these (through past moves) he could exploit these and improve his

²In game theory, games like chess are often called ‘extensive’, whereas ‘strategic’ is reserved for a different kind of game.

performance by deviating from the minimax strategy. At least, the classical game theory of Nash equilibria does not take into account limited rationality, whereas the AI ξ agent should.

6.3.3 Using the AI μ Model for Game Playing

In the following, we demonstrate the applicability of the AI μ model to games. The AI μ model takes the position of player 1. The environment provides the evaluation V . For a symmetric situation we could take a second AI μ model as player 2, but for simplicity we take the environment as the second player and assume that this environmental player behaves according to the minimax strategy (6.18). The environment serves as a perfect player *and* as a teacher, albeit a very crude one as it tells the agent at the end of the game, only whether it won or lost.

The minimax behavior of player 2 can be expressed by a (deterministic) probability distribution μ^{SG} as the following

$$\mu^{SG}(y'_1 \underline{x}'_1 \dots y'_n \underline{x}'_n) := \begin{cases} 1 & \text{if } x'_k = \arg \min_{x''_k} \dots \max_{y''_n} \min_{x''_n} V(y'_1 x'_1 \dots y'_k x''_k \dots y''_n x''_n) \quad \forall k \\ 0 & \text{otherwise} \end{cases} \quad (6.20)$$

The probability that player 2 makes move x'_k is $\mu^{SG}(y'_1 \underline{x}'_1 \dots y'_k \underline{x}'_k)$ which is 1 for $x'_k = \dot{x}'_k$ as defined in (6.18) and 0 otherwise.

Clearly, the AI μ system receives no feedback, i.e. $r_1 = \dots = r_{n-1} = 0$, until the end of the game, where it should receive positive/negative/neutral feedback on a win/loss/draw, i.e. $r_n = V(\dots)$. The environmental prior probability is therefore

$$\mu^{AI}(y_1 \underline{x}_1 \dots y_n \underline{x}_n) = \begin{cases} \mu^{SG}(y'_1 \underline{x}'_1 \dots y'_n \underline{x}'_n) & \text{if } r_1 \dots r_{n-1} = 0 \text{ and } r_n = V(y'_1 x'_1 \dots y'_n x'_n) \\ 0 & \text{otherwise} \end{cases} \quad (6.21)$$

where $y_i = y'_i$ and $x_i = r_i x'_i$. If the environment is a minimax player (6.18) plus a crude teacher V , i.e. if μ^{AI} is the true prior probability, the question now is, what is the behavior \dot{y}_k^{AI} of the AI μ agent. It turns out that if we set $m_k = n$ the AI μ agent is also a minimax player (6.19) and hence optimal

$$\begin{aligned} \dot{y}_k^{AI} &= \arg \max_{y_k} \sum_{x'_k} \dots \max_{y_n} \sum_{x'_n} V(\dot{y}'_{<k} y'_{k:n}) \cdot \mu^{SG}(\dot{y}'_{<k} y'_{k:n}) = \\ &= \arg \max_{y_k} \sum_{x'_k} \dots \max_{y_{n-1}} \sum_{x'_{n-1}} \max_{y_n} \min_{x'_n} V(\dot{y}'_{<k} y'_{k:n}) \cdot \mu^{SG}(\dot{y}'_{<k} y'_{k:n-1}) = \\ &= \dots = \arg \max_{y_k} \min_{x'_{k+1}} \dots \max_{y_n} \min_{x'_n} V(\dot{y}'_{<k} y'_{k:n}) = \dot{y}_k^{SG} \end{aligned} \quad (6.22)$$

In the first line we inserted $m_k = n$ and (6.21) into the definition (6.4) of \dot{y}_k^{AI} . This removes all sums over the r_i . Further, the sum over x'_n gives only a contribution for $x'_n = \arg \min_{x'_n} V(y'_1 \underline{x}'_1 \dots y'_n x'_n)$ by definition (6.20) of μ^{SG} . Inserting this x'_n gives

the second line. Effectively, μ^{SG} is reduced to a lower number of arguments and the sum over x'_n replaced by $\min_{x'_n}$. Repeating this procedure for $x'_{n-1}, \dots, x'_{k+1}$ leads to the last line, which is just the minimax strategy of player 1 defined in (6.19).

Let us now assume that the game under consideration is played s times. The prior probability then is

$$\mu^{AI}(\underline{y}_1 \dots \underline{y}_{sn}) = \prod_{r=0}^{s-1} \mu_1^{AI}(\underline{y}_{rn+1} \dots \underline{y}_{(r+1)n}) \quad (6.23)$$

where we have renamed the prior probability (6.21) for one game to μ_1^{AI} . (6.23) is a special case of a factorizable μ (defined in Section 4.3.1) with identical factors $\mu_r = \mu_1^{AI}$ for all r and equal episode lengths $n_{r+1} - n_r = n$. The AI μ agent (6.23) for repeated game playing also implements the minimax strategy,

$$\dot{y}_k^{AI} = \arg \max_{y_k} \min_{x'_k} \dots \max_{y_{(r+1)n}} \min_{x'_{(r+1)n}} V(\dot{y}'_{rn+1:k-1} \dots \dot{y}'_{k:(r+1)n}) \quad (6.24)$$

with r such that $rn < k \leq (r+1)n$ and for any choice of m_k as long as the horizon $h_k \geq n$. This can be proven by using (4.27) and (6.22).

6.3.4 Games of Variable Length

We have argued that a single game of variable but bounded length can be padded to a fixed length without effect. We now analyze in a sequence of games the effect of replacing the games with fixed length by games of variable length. The sequence $y'_1 x'_1 \dots y'_n x'_n$ can still be grouped into episodes corresponding to the moves of separated consecutive games, but now the length and total number of games that fit into the n moves depend on the actual moves taken³. $V(y'_1 x'_1 \dots y'_n x'_n)$ equals the number of games where the agent wins, minus the number of games where the environment wins. Whenever a loss, win or draw has been achieved by the agent or the environment, a new game starts. The player whose turn it would next be, begins the next game. The games are still separated in the sense that the behavior and reward of the current game does not influence the next game. On the other hand, they are slightly entangled, because the length of the current game determines the time of start of the next. As the rules of the game are time invariant, this does not influence the next game directly. If we play a fixed number of games, the games are completely independent, but if we play a fixed number of total moves n , the number of games depends on their lengths. This has the following consequences: the better player tries to keep the games short, to win more games in the given time n . The poorer player tries to draw the games out, in order to lose less games. The better player might further prefer a quick draw, rather than to win a long game. Formally, this entanglement is represented by the fact that the prior probability μ does no longer factorize. The reduced form (6.24) of \dot{y}_k^{AI} to one episode is no longer valid. Also,

³If the sum of game lengths do not fit exactly into n moves, we pad the last game appropriately.

the behavior y_k^{AI} of the agent depends on m_k , even if the horizon h_k is chosen larger than the longest possible game. The important point is that the agent realizes that keeping games short/long can lead to increased reward. In practice, a horizon much larger than the average game length should be sufficient to incorporate this effect. The details of games in the distant future do not affect the current game and can, therefore, be ignored. A more quantitative analysis could be interesting, but would lead us too far astray.

6.3.5 Using the AI ξ Model for Game Playing

When going from the specific AI μ model, where the rules of the game have been explicitly modeled into the prior probability μ^{AI} , to the universal model AI ξ we have to ask whether these rules can be learned from the assigned rewards r_k . Here, another (actually the main) reason for studying the case of repeated games, rather than just one game arises. For a single game there is only one cycle of non-trivial feedback namely the end of the game – too late to be useful except when there are further games following.

Even in the case of repeated games, there is only very limited feedback, at most $\log_2 3$ bits of information per game if the 3 outcomes win/loss/draw have the same frequency. So there are at least $O(K(game))$ number of games necessary to learn a game of complexity $K(game)$. Apart from extremely simple games, even this estimate is far too optimistic. As the AI ξ agent has no information about the game to begin with, its moves will be more or less random and it can win the first few games merely by pure luck. So the probability that the agent loses is near to one and hence the information content I in the feedback r_k at the end of the game is much less than $\log_2 3$. This situation remains for a very large number of games. But in principle, every game should be learnable after a very long sequence of games even with this minimal feedback only, as long as $I \neq 0$.

The important point is that no other learning scheme with no extra information can learn the game more quickly than AI ξ . We expect this to be true as μ^{AI} factorizes in the case of games of fixed length, i.e. μ^{AI} satisfies a strong separability condition. In the case of variable game length the entanglement is also low. μ^{AI} should still be sufficiently separable allowing us to formulate and prove good reward bounds for AI ξ .

To learn realistic games like tic-tac-toe (noughts and crosses) in realistic time one has to provide more feedback. This could be achieved by intermediate help during the game. The environment could give positive(negative) feedback for every good(bad) move the agent makes. The demand on whether a move is to be valued as good should be adapted to the gained experience of the agent in such a way that approximately half of the moves are valued as good and the other half as bad, in order to maximize the information content of the feedback.

For more complicated games like chess, even more feedback is necessary from a practical point of view. One way to increase the feedback far beyond a few bits

per cycle is to train the agent by teaching it good moves. This is called supervised learning. Despite the fact that the $\text{AI}\mu$ model has only a reward feedback r_k , it is able to learn supervised, as will be shown in Section 6.5. Another way would be to start with more simple games containing certain aspects of the true game and to switch to the true game when the agent has learned the simple game.

No other difficulties are expected when going from μ to ξ . Eventually ξ^{AI} will converge to the minimax strategy μ^{AI} . In the more realistic case, where the environment is not a perfect minimax player, $\text{AI}\xi$ can detect and exploit the weakness of the opponent.

Finally, we want to comment on the input/output space \mathcal{X}/\mathcal{Y} of the AI models. In practical applications, \mathcal{Y} will possibly include also illegal moves. If \mathcal{Y} is the set of moves of e.g. a robotic arm, the agent could move a wrong figure or even knock over the figures. A simple way to handle illegal moves y_k is by interpreting them as losing moves, which terminate the game. Further, if e.g. the input x_k is the image of a video camera which makes one shot per move, \mathcal{X} is not the set of moves by the environment but includes the set of states of the game board. The discussion in this section handles this case as well. There is no need to explicitly design the systems I/O space \mathcal{X}/\mathcal{Y} for a specific game.

The discussion above on the $\text{AI}\xi$ agent was rather informal for the following reason: game playing (the $\text{SG}\xi$ agent) has (nearly) the same complexity as fully general AI, and quantitative results for the $\text{AI}\xi$ agent are difficult (but not impossible) to obtain.

6.4 Function Minimization (FM)

6.4.1 Applications/Examples

There are many problems that can be reduced to a function minimization problem (FM). The minimum of a (real valued) function $f: \mathcal{Y} \rightarrow \mathbb{R}$ over some domain \mathcal{Y} or a good approximate to the minimum has to be found, usually with some limited resources.

One popular example is the traveling salesman problem (TSP). \mathcal{Y} is the set of different routes between towns and $f(y)$ the length of route $y \in \mathcal{Y}$. The task is to find a route of minimal length visiting all cities. This problem is NP hard. Getting good approximations in limited time is of great importance in various applications. Another example is the minimization of production costs (MPC), e.g. of a car, under several constraints. \mathcal{Y} is the set of all alternative car designs and production methods compatible with the specifications and $f(y)$ the overall cost of alternative $y \in \mathcal{Y}$. A related example is finding materials or (bio)molecules with certain properties (MAT). E.g. solids with minimal electrical resistance or maximally efficient chlorophyll modifications or aromatic molecules that taste as close as possible to strawberry. We can also ask for nice paintings (NPT). \mathcal{Y} is the set of all existing

or imaginable paintings and $f(y)$ characterizes how much person A likes painting y . The agent should present paintings, which A likes.

For now, these are enough examples. The TSP is very rigorous from a mathematical point of view, as f , i.e. an algorithm of f , is usually known. In principle, the minimum could be found by exhaustive search, were it not for computational resource limitations. For MPC, f can often be modeled in a reliable and sufficiently accurate way. For MAT you need very accurate physical models, which might be unavailable or too difficult to solve or implement. For NPT all we have is the judgement of person A on every presented painting. The evaluation function f cannot be implemented without scanning A 's brain, which is not possible with today's technology.

So there are different limitations, some depending on the application we have in mind. An implementation of f might not be available, f can only be tested at some arguments y and $f(y)$ is determined by the environment. We want to (approximately) minimize f with as few function calls as possible or, conversely, find an as close as possible approximation for the minimum within a fixed number of function evaluations. If f is available or can quickly be inferred by the agent and evaluation is quick, it is more important to minimize the total time needed to imagine new trial minimum candidates plus the evaluation time for f . As we do not consider computational aspects of AI ξ till Section 7.2 we concentrate on the first case, where f is not available or dominates the computational requirements.

6.4.2 The Greedy Model FMG μ

The FM model consists of a sequence $\dot{y}_1 \dot{z}_1 \dot{y}_2 \dot{z}_2 \dots$ where \dot{y}_k is a trial of the FM agent for a minimum of f and $\dot{z}_k = f(\dot{y}_k)$ is the true function value returned by the environment. We randomize the model by assuming a probability distribution $\mu(f)$ over the functions. There are several reasons for doing this. We might really not know the exact function f , as in the NPT example, and model our uncertainty by the probability distribution μ . More importantly, we want to parallel the other AI classes, like in the SP μ model, where we always started with a probability distribution μ that was finally replaced by ξ to get the universal Solomonoff prediction SP ξ . We want to do the same thing here. Further, the probabilistic case includes the deterministic case by choosing $\mu(f) = \delta_{ff_0}$, where f_0 is the true function. A final reason is that the deterministic case is trivial when μ and hence f_0 is known, as the agent can internally (virtually) check all function arguments and output the correct minimum from the very beginning.

We will assume that \mathcal{Y} is countable or finite and that μ is a discrete measure, e.g. by taking only computable functions. The probability that the function values of y_1, \dots, y_n are z_1, \dots, z_n is then given by

$$\mu^{FM}(y_1 \dot{z}_1 \dots y_n \dot{z}_n) := \sum_{f: f(y_i) = \dot{z}_i \ \forall 1 \leq i \leq n} \mu(f) \quad (6.25)$$

We start with a model that minimizes the expectation z_k of the function value f for the next output y_k , taking into account previous information:

$$\dot{y}_k := \arg \min_{y_k} \sum_{z_k} z_k \cdot \mu(\dot{y}_1 \dot{z}_1 \dots \dot{y}_{k-1} \dot{z}_{k-1} y_k \dot{z}_k)$$

This type of greedy algorithm, just minimizing the next feedback, was sufficient for sequence prediction (SP) and is also sufficient for classification (CF). It is, however, not sufficient for function minimization as the following example demonstrates.

Take $f : \{0,1\} \rightarrow \{1,2,3,4\}$. There are 16 different functions which shall be equiprobable, $\mu(f) = \frac{1}{16}$. The function expectation in the first cycle

$$\langle z_1 \rangle := \sum_{z_1} z_1 \cdot \mu(y_1 \dot{z}_1) = \frac{1}{4} \sum_{z_1} z_1 = \frac{1}{4}(1+2+3+4) = 2.5$$

is just the arithmetic average of the possible function values and is independent of y_1 . Therefore, $\dot{y}_1 = 0$, as argmin is defined to take the lexicographically first minimum in an ambiguous case. Let us assume that $f_0(0) = 2$, where f_0 is the true environment function, i.e. $\dot{z}_1 = 2$. The expectation of z_2 is then

$$\langle z_2 \rangle := \sum_{z_2} z_2 \cdot \mu(0 \dot{y}_2 \dot{z}_2) = \begin{cases} 2 & \text{for } y_2 = 0 \\ 2.5 & \text{for } y_2 = 1 \end{cases}$$

For $y_2 = 0$ the agent already knows $f(0) = 2$, for $y_2 = 1$ the expectation is, again, the arithmetic average. The agent will again output $\dot{y}_2 = 0$ with feedback $\dot{z}_2 = 2$. This will continue forever. The agent is not motivated to explore other y 's as $f(0)$ is already smaller than the expectation of $f(1)$. This is obviously not what we want. The greedy model fails. The agent ought to be inventive and try other outputs when given enough time.

The general reason for the failure of the greedy approach is that the information contained in the feedback z_k depends on the output y_k . A FM agent can actively influence the knowledge it receives from the environment by the choice in y_k . It may be more advantageous to first collect certain knowledge about f by an (in greedy sense) non-optimal choice for y_k , rather than to minimize the z_k expectation immediately. The non-minimality of z_k might be over-compensated in the long run by exploiting this knowledge. In SP, the received information is always the current bit of the sequence, independent of what SP predicts for this bit. This is the reason why a greedy strategy in the SP case is already optimal.

6.4.3 The General FM μ/ξ Model

To get a useful model we have to think more carefully about what we really want. Should the FM agent output a good minimum in the last output in a limited number of cycles m , or should the average of the z_1, \dots, z_m values be minimal, or does it suffice that just one of the z is as small as possible? Let us define the FM μ model as to

minimize the μ averaged weighted sum $\alpha_1 z_1 + \dots + \alpha_m z_m$ for some given $\alpha_k \geq 0$. Building the μ average by summation over the z_i and minimizing w.r.t. the y_i has to be performed in the correct chronological order. With a similar reasoning as in (6.1) to (6.4) we get

$$\dot{y}_k^{FM} = \arg \min_{y_k} \sum_{z_k} \dots \min_{y_m} \sum_{z_m} (\alpha_1 z_1 + \dots + \alpha_m z_m) \cdot \mu(\dot{y}_1 \dot{z}_1 \dots \dot{y}_{k-1} \dot{z}_{k-1} y_k z_k \dots y_m z_m) \quad (6.26)$$

If we want the final output \dot{y}_m to be optimal we should choose $\alpha_k = 0$ for $k < m$ and $\alpha_m = 1$ (final model FMF μ). If we want to already have a good approximation during intermediate cycles, we should demand that the output of all cycles together are optimal in some average sense, so we should choose $\alpha_k = 1$ for all k (sum model FMS μ). If we want to have something in between, for instance, increase the pressure to produce good outputs, we could choose the $\alpha_k = e^{\gamma(k-m)}$ exponentially increasing for some $\gamma > 0$ (exponential model FME μ). For $\gamma \rightarrow \infty$ we get the FMF μ , for $\gamma \rightarrow 0$ the FMS μ model. If we want to demand that the best of the outputs $y_1 \dots y_k$ is optimal, we must replace the α weighted z -sum by $\min\{z_1, \dots, z_m\}$ (minimum Model FMM μ). We expect the behavior to be very similar to the FMF μ model, and do not consider it further.

By construction, the FM μ models guarantee optimal results in the usual sense that no other model knowing only μ can be expected to produce better results. The variety of FM variants is not a fault of the theory. They just reflect the fact that there is some interpretational freedom of what is meant by minimization within m function calls. In most applications, probably FMF is appropriate. In the NPT application one might prefer the FMS model.

The interesting case (in AI) is when μ is unknown. We define for this case, the FM ξ model by replacing $\mu(f)$ with some $\xi(f)$, which should assign high probability to functions f of low complexity. So we might define $\xi(f) = \sum_{q: \forall x [U(qx)=f(x)]} 2^{-l(q)}$. The problem with this definition is that it is, in general, undecidable whether a TM q is an implementation of a function f . $\xi(f)$ defined in this way is uncomputable, not even approximable. As we only need a ξ analogous to the l.h.s. of (6.25), the following definition is natural

$$\xi^{FM}(y_1 \dot{z}_1 \dots y_n \dot{z}_n) := \sum_{q: q(y_i)=z_i \ \forall 1 \leq i \leq n} 2^{-l(q)} \quad (6.27)$$

ξ^{FM} is actually equivalent to inserting the uncomputable $\xi(f)$ into (6.25). ξ^{FM} is an enumerable semi-measure and dominates all enumerable probability distributions of the form (6.25). We will not prove this here.

Alternatively, we could have constrained the sum in (6.27) by $q(y_1 \dots y_n) = z_1 \dots z_n$ analogous to (6.6), but these two definitions are not equivalent. Definition (6.27) ensures the symmetry⁴ in its arguments and $\xi^{FM}(\dots y \dot{z} \dots y \dot{z}' \dots) = 0$ for $z \neq z'$. It incorporates all general knowledge we have about function minimization, whereas (6.6)

⁴See [Sol99] for a discussion on symmetric universal distributions on unordered data.

does not. But this extra knowledge has only low information content (complexity of $O(1)$), so we do not expect FM ξ to perform much worse when using (6.6) instead of (6.27). But there is no reason to deviate from (6.27) at this point.

We can now define an “error” measure $E_{m\mu}^{FM}$ as (6.26) with $k=1$ and argmin_{y_1} replaced by \min_{y_1} and, additionally, μ replaced by ξ for $E_{m\xi}^{FM}$. We expect $|E_{m\xi}^{FM} - E_{m\mu}^{FM}|$ to be bounded in a way that justifies the use of ξ instead of μ for computable μ , i.e. computable f_0 in the deterministic case. The arguments are the same as for the AI ξ model.

6.4.4 Is the General Model Inventive?

In the following we will show that FM ξ will never cease searching for minima, but will test an infinite set of different y' s for $m \rightarrow \infty$.

Let us assume that the agent tests only a finite number of $y_i \in \mathcal{A} \subset \mathcal{Y}$, $|\mathcal{A}| < \infty$. Let $t-1$ be the cycle in which the last new $y \in \mathcal{A}$ is selected (or some later cycle). Selecting y' s in cycles $k \geq t$ a second time, the feedback z does not provide any new information, i.e. does not modify the probability ξ^{FM} . The agent can minimize $E_{m\xi}^{FM}$ by outputting in cycles $k \geq t$ the best $y \in \mathcal{A}$ found so far (in the case $\alpha_k = 0$, the output does not matter). Let us fix f for a moment. Then we have

$$E^a := \alpha_1 z_1 + \dots + \alpha_m z_m = \sum_{k=1}^{t-1} \alpha_k f(y_k) + f_1 \cdot \sum_{k=t}^m \alpha_k, \quad f_1 := \min_{1 \leq k < t} f(y_k)$$

Let us now assume that the agent tests one additional $y_t \notin \mathcal{A}$ in cycle t , but no other $y \notin \mathcal{A}$. Again, it will keep to the best output for $k > t$, which is either the one of the previous agent or y_t .

$$E^b = \sum_{k=1}^t \alpha_k f(y_k) + \min\{f_1, f(y_t)\} \cdot \sum_{k=t+1}^m \alpha_k$$

The difference can be represented in the form

$$E^a - E^b = \left(\sum_{k=t}^m \alpha_k \right) \cdot f^+ - \alpha_t \cdot f^- \quad , \quad f^\pm := \max\{0, \pm(f_1 - f(y_t))\} \geq 0.$$

As the true FM strategy is the one which minimizes E , assumption a is ruled out if $E^a > E^b$. We will say that b is favored over a , which does not mean that b is the correct strategy, only that a is not the true one. For probability distributed f , b is favored over a when

$$E^a - E^b = \left(\sum_{k=t}^m \alpha_k \right) \cdot \langle f^+ \rangle - \alpha_t \cdot \langle f^- \rangle > 0 \quad \Leftrightarrow \quad \sum_{k=t}^m \alpha_k > \alpha_t \frac{\langle f^- \rangle}{\langle f^+ \rangle}$$

where $\langle f^\pm \rangle$ is the ξ expectation of $\pm(f_1 - f(y_t))$ under the condition that $\pm f_1 \geq \pm f(y_t)$ and under the constraints imposed in cycles $1 \dots t-1$. As ξ assigns a strictly

positive probability to every non-empty event, $\langle f^+ \rangle \neq 0$. Inserting $\alpha_k = e^{\gamma(k-m)}$, assumption *a* is ruled out in model FME ξ if

$$m - t > \frac{1}{\gamma} \ln \left[1 + \frac{\langle f^- \rangle}{\langle f^+ \rangle} (e^\gamma - 1) \right] - 1 \rightarrow \begin{cases} 0 & \text{for } \gamma \rightarrow \infty \text{ (FMF}\xi\text{)} \\ \langle f^- \rangle / \langle f^+ \rangle - 1 & \text{for } \gamma \rightarrow 0 \text{ (FMS}\xi\text{)} \end{cases}$$

We see that if the condition is not satisfied for some t , it will remain wrong for all $t' > t$. So the FME ξ agent will test each y only once up to a point from which on it always outputs the best found y . Further, for $m \rightarrow \infty$ the condition always gets satisfied. As this is true for any finite \mathcal{A} , the assumption of a finite \mathcal{A} is wrong. For $m \rightarrow \infty$ the agent tests an increasing number of different y 's, provided \mathcal{Y} is infinite. The FMF ξ model will never repeat any y except in the last cycle m where it chooses the best found y . The FMS ξ model will test a new y_t for fixed m , only if the expected value of $f(y_t)$ is not too large.

The above does not necessarily hold for other choices of α_k . The above also holds for the FMF μ agent if $\langle f^+ \rangle \neq 0$. $\langle f^+ \rangle = 0$ if the agent can already exclude that y_t is a better guess, so there is no reason to test it explicitly.

Nothing has been said about the quality of the guesses, but for the FM μ agent they are optimal by definition. If $K(\mu)$ for the true distribution μ is finite, we expect the FM ξ agent to solve the ‘exploration versus exploitation’ problem in a universally optimal way, as ξ converges to μ .

6.4.5 Using the AI models for Function Minimization

The AI μ model can be used for function minimization in the following way. The output y_k of cycle k is a guess for a minimum of f , like in the FM model. The reward r_k should be high for small function values $z_k = f(y_k)$. The reward should also be weighted with α_k to reflect the same strategy as in the FM case. The choice of $r_k = -\alpha_k z_k$ is natural. Here, the feedback is not binary but $r_k \in \mathcal{R} \subset \mathbb{R}$, with \mathcal{R} being a countable subset of \mathbb{R} , e.g. the computable reals or all rational numbers. The feedback x'_k should be the function value $f(y_k)$. So we set $x'_k = z_k$. Note, that there is a redundancy if $\alpha_{()}$ is a computable function with no zeros, as $r_k = -\alpha_k x'_k$. So, for small $K(\alpha_{()})$ like in the FMS model, one might set $x_k \equiv \epsilon$. If we keep x'_k the AI prior probability is

$$\mu^{AI}(y_1 \underline{x}_1 \dots y_n \underline{x}_n) = \begin{cases} \mu^{FM}(y_1 \underline{z}_1 \dots y_n \underline{z}_n) & \text{for } r_k = -\alpha_k z_k, \ x'_k = z_k, \ x_k = r_k x'_k \\ 0 & \text{else.} \end{cases} \quad (6.28)$$

Inserting this into (6.4) with $m_k = m$ we get

$$\begin{aligned} \dot{y}_k^{AI} &= \arg \max_{y_k} \sum_{x_k} \dots \max_{y_m} \sum_{x_m} (r_k + \dots + r_m) \cdot \mu^{AI}(\dot{y}_1 \dot{x}_1 \dots y_k \underline{x}_k \dots y_m \underline{x}_m) = \\ &= \arg \min_{y_k} \sum_{z_k} \dots \min_{y_m} \sum_{z_m} (\alpha_k z_k + \dots + \alpha_T z_m) \cdot \mu^{FM}(\dot{y}_1 \dot{z}_1 \dots y_k \underline{z}_k \dots y_m \underline{z}_m) = \dot{y}_k^{FM} \end{aligned}$$

where \dot{y}_k^{FM} has been defined in (6.26). The proof of equivalence was so simple because the FM model has already a rather general structure, which is similar to the full $AI\mu$ model.

One might expect no problems when going from the already very general FM ξ model to the universal AI ξ model (with $m_k=m$), but there is a pitfall in the case of the FMF model. All rewards r_k are zero in this case, except for the last one being r_m . Although there is a feedback z_k in every cycle, the AI ξ agent cannot learn from this feedback as it is not told that in the final cycle r_m will equal to $-z_m$. There is no problem in the FM ξ model because in this case this knowledge is hardcoded into ξ^{FM} . The AI ξ model must first learn that it has to minimize a function but it can only learn if there is a non-trivial credit assignment r_k . FMF works for repeated minimization of (different) functions, such as minimizing N functions in $N \cdot m$ cycles. In this case there are N non-trivial feedbacks and AI ξ has time to learn that there is a relation between $r_{k \cdot m}$ and $x'_{k \cdot m}$ every m^{th} cycle. This situation is similar to the case of strategic games discussed in Section 6.3.

There is no problem in applying AI ξ to FMS because the r feedback provides enough information in this case. The only thing the AI ξ model has to learn, is to ignore the x' feedbacks as all information is already contained in r . Interestingly the same argument holds for the FME model if $K(\gamma)$ and $K(m)$ are small⁵. The AI ξ model has additionally only to learn the relation $r_k = -e^{-\gamma(k-m)}x'_k$. This task is simple as every cycle provides one data point for a simple function to learn. This argument is no longer valid for $\gamma \rightarrow \infty$ as $K(\gamma) \rightarrow \infty$ in this case.

6.4.6 Remark

TSP seems to be trivial in the $AI\mu$ model but non-trivial in the AI ξ model. The reason being that (6.26) just implements an internal complete search as $\mu(f) = \delta_{f^{TSP}}$ contains all necessary information. $AI\mu$ outputs from the very beginning, the exact minimum of f^{TSP} . This "solution" is, of course, unacceptable from performance perspective. As long as we give no efficient approximation ξ^c of ξ , we have not contributed anything to a solution of the TSP by using AI ξ^c . The same is true for any other problem where f is computable and easily accessible. Therefore, TSP is not (yet) a good example because all we have done is to replace a NP complete problem with the uncomputable AI ξ model or by a computable AI ξ^c model, for which we have said nothing about computation time yet. It is simply an overkill to reduce simple problems to AI ξ . TSP is a simple problem in this respect, until we consider the AI ξ^c model seriously. For the other examples, where f is inaccessible or complicated, an AI ξ^c model would provide a true solution to the minimization problem as an explicit definition of f is not needed for AI ξ and AI ξ^c . A computable version of AI ξ will be defined in Section 7.2.

⁵If we set $\alpha_k = e^{\gamma k}$ the condition on $K(m)$ can be dropped.

6.5 Supervised Learning from Examples (EX)

The developed AI models provide a frame for reinforcement learning. The environment provides feedback r , informing the agent about the quality of its last (or earlier) output y ; it assigns reward r to output y . In this sense, reinforcement learning is explicitly integrated into the $AI\rho$ model. $AI\mu$ maximizes the true expected reward, whereas the $AI\xi$ model is a universal, environment independent, reinforcement learning algorithm.

There is another type of learning method: Supervised learning by presentation of examples (EX). Many problems learned by this method are association problems of the following type. Given some examples $x \in R \subset \mathcal{X}$, the agent should reconstruct, from a partially given x' , the missing or corrupted parts, i.e. complete x' to x such that relation R contains x . In many cases, \mathcal{X} consists of pairs (z, v) , where v is the possibly missing part.

6.5.1 Applications/Examples

Learning functions by presenting $(z, f(z))$ pairs and asking for the function value of z by presenting $(z, ?)$ falls into this category.

A basic example is learning properties of geometrical objects coded in some way. E.g. if there are 18 different objects characterized by their size (small or big), their colors (red, green or blue) and their shapes (square, triangle, circle), then $(object, property) \in R$ if the *object* possesses the *property*. Here, R is a relation which is not the graph of a single valued function.

When teaching a child, by pointing to objects and saying “this is a tree” or “look how green” or “how beautiful”, one establishes a relation of $(object, property)$ pairs in R . Pointing to a (possibly different) tree later and asking “what is this ?” corresponds to a partially given pair $(object, ?)$, where the missing part “?” should be completed by the child saying “tree”.

A final example we want to give is chess. We have seen that, in principle, chess can be learned by reinforcement learning. In the extreme case the environment only provides reward $r = 1$ when the agent wins. The learning rate is completely unacceptable from a practical point of view. The reason is the very low amount of information feedback. A more practical method of teaching chess is to present example games in the form of sensible $(board-state, move)$ sequences. They contain information about legal and good moves (but without any explanation). After several games have been presented, the teacher could ask the agent to make its own move by presenting $(board-state, ?)$ and then evaluate the answer of the agent.

6.5.2 Supervised Learning with the $AI\mu/\xi$ Model

Let us define the EX model as follows: The environment presents inputs $x'_{k-1} = z_k v_k \equiv (z_k, v_k) \in R \cup (Z \times \{?\}) \subset Z \times (\mathcal{Y} \cup \{?\}) = \mathcal{X}'$ to the agent in cycle $k-1$. The agent is

expected to output y_k in the next cycle, which is evaluated with $r_k=1$ if $(z_k, y_k) \in R$ and 0 otherwise. To simplify the discussion, an output y_k is expected and evaluated even when $v_k(\neq?)$ is given. To complete the description of the environment, the probability distribution $\mu_R(\underline{x}'_1 \dots \underline{x}'_n)$ of the examples and questions x'_i (depending on R) has to be given. Wrong examples should not occur, i.e. μ_R should be 0 if $x'_i \notin R \cup (Z \times \{?\})$ for some $1 \leq i \leq n$. The relations R might also be probability distributed with $\sigma(\underline{R})$. The example prior probability in this case is

$$\mu(\underline{x}'_1 \dots \underline{x}'_n) = \sum_R \mu_R(\underline{x}'_1 \dots \underline{x}'_n) \cdot \sigma(\underline{R}) \quad (6.29)$$

The knowledge of the valuation r_k on output y_k restricts the possible relations R , consistent with $R(z_k, y_k) = r_k$, where $R(z, y) := 1$ if $(z, y) \in R$ and 0 otherwise. The prior probability for the input sequence $x_1 \dots x_n$ if the output sequence is $y_1 \dots y_n$, is therefore

$$\mu^{AI}(y_1 \underline{x}_1 \dots y_n \underline{x}_n) = \sum_{R: \forall 1 < i \leq n [R(z_i, y_i) = r_i]} \mu_R(\underline{x}'_1 \dots \underline{x}'_n) \cdot \sigma(\underline{R})$$

where $x_i = r_i x'_i$ and $x'_{i-1} = z_i v_i$ with $v_i \in \mathcal{Y} \cup \{?\}$. In the I/O sequence $y_1 x_1 y_2 x_2 \dots = y_1 r_1 z_2 v_2 y_2 r_2 z_3 v_3 \dots$ the $r_1 y_1$ are dummies, after which regular behavior starts with example (z_2, v_2) .

The $AI\mu$ model is optimal by construction of μ^{AI} . For computable prior μ_R and σ , we expect a near optimal behavior of the universal $AI\xi$ model if μ_R additionally satisfies some separability property. In the following, we give some motivation why the $AI\xi$ model takes into account the supervisor information contained in the examples and why it learns faster than by reinforcement.

We keep R fixed and assume $\mu_R(x'_1 \dots x'_n) = \mu_R(x'_1) \cdot \dots \cdot \mu_R(x'_n) \neq 0 \Leftrightarrow x'_i \in R \cup (Z \times \{?\}) \forall i$ to simplify the discussion. Short codes q contribute most to $\xi^{AI}(y_1 \underline{x}_1 \dots y_n \underline{x}_n)$. As $x'_1 \dots x'_n$ is distributed according to the computable probability distribution μ_R , a short code of $x'_1 \dots x'_n$ for large enough n is a Huffman code w.r.t. the distribution μ_R . So we expect μ_R and hence R to be coded in the dominant contributions to ξ^{AI} in some way, where the plausible assumption was made that the y on the input tape do not matter. Much more than one bit per cycle will usually be learned, i.e. relation R will be learned in $n \ll K(R)$ cycles by appropriate examples. This coding of R in q evolves independently of the feedbacks r . To maximize the feedback r_k , the agent has to learn to output a y_k with $(z_k, y_k) \in R$. The agent has to invent a program extension q' to q , which extracts z_k from $x_{k-1} = (z_k, ?)$ and searches for and outputs a y_k with $(z_k, y_k) \in R$. As R is already coded in q , q' can re-use this coding of R in q . The size of the extension q' is, therefore, of order 1. To learn this q' , the agent requires feedback r with information content $O(1) = K(q')$ only.

Let us compare this with reinforcement learning, where only $x'_k = (z_k, ?)$ pairs are presented. A coding of R in a short code q for $x'_1 \dots x'_n$ is of no use and will therefore be absent. Only the rewards r force the agent to learn R . q' is therefore expected to be of size $K(R)$. The information content in the r 's must be of the order $K(R)$. In practice, there are often only very few $r_k=1$ at the beginning of the learning phase

and the information content in $r_1 \dots r_n$ is much less than n bits. The required number of cycles to learn R by reinforcement is, therefore, at least but in many cases much larger than $K(R)$.

Although $AI\xi$ was never designed or told to learn supervised, it learns how to take advantage of the examples from the supervisor. μ_R and R are learned from the examples, the rewards r are not necessary for this process. The remaining task of learning how to learn supervised is then a simple task of complexity $O(1)$, for which the rewards r are necessary.

6.6 Other Aspects of Intelligence

In AI, a variety of general ideas and methods have been developed. In the last sections, we have seen how several problem classes can be formulated within $AI\xi$. As we claim universality of the $AI\xi$ model, we want to enlight which of, and how the other AI methods are incorporated in the $AI\xi$ model, by looking at its structure. Some methods are directly included, others are or should be emergent. We do not claim the following list to be complete.

Probability theory and *utility theory* are the heart of the $AI\mu/\xi$ models. The probabilities are the true/universal behaviors of the environment. The utility function is what we called total reward, which should be maximized. Maximization of an expected utility function in a probabilistic environment is usually called *sequential decision theory*, and is explicitly integrated in full generality in our model. In a sense this includes probabilistic (a generalization of deterministic) *reasoning*, where the objects of reasoning are not true and false statements, but the prediction of the environmental behavior. *Reinforcement Learning* is explicitly built in, due to the rewards. Supervised learning is an emergent phenomenon (Section 6.5). *Algorithmic information theory* leads us to use ξ as a universal estimate for the prior probability μ .

For horizon > 1 , the expectimax series in (6.4) and the process of selecting maximal values may be interpreted as abstract *planning*. The expectimax series is a form of *informed search*, in the case of $AI\mu$, and *heuristic search*, for $AI\xi$, where ξ could be interpreted as a heuristic for μ . The minimax strategy of *game playing* in case of $AI\mu$ is also subsumed. The $AI\xi$ model converges to the minimax strategy if the environment is a minimax player but it can also take advantage of environmental players with limited rationality. *Problem solving* occurs (only) in the form of how to maximize the expected future reward.

Knowledge is accumulated by $AI\xi$ and is stored in some form not specified further on the work tape. Any kind of information in any representation on the inputs y is exploited. The problem of *knowledge engineering* and *representation* appears in the form of how to train the $AI\xi$ model. More practical aspects, like *language* or *image processing* have to be learned by $AI\xi$ from scratch.

Other theories, like *fuzzy logic*, *possibility theory*, *Dempster-Shafer theory*, ... are

partly outdated and partly reducible to Bayesian probability theory [Che85, Che88]. The interpretation and consequences of the evidence gap $g := 1 - \sum_{x_k} \xi(y_{<k} \underline{y}_k) > 0$ in ξ may be similar to those in Dempster-Shafer theory. Boolean logical reasoning about the external world plays, at best, an emergent role in the AI ξ model.

Other methods, which do not seem to be contained in the AI ξ model might also be emergent phenomena. The AI ξ model has to construct short codes of the environmental behavior, the AI ξ^{il} (see next section) has to construct short action programs. If we would analyze and interpret these programs for realistic environments, we might find some of the unmentioned or unused or new AI methods at work in these programs. This is, however, pure speculation at this point. More important: when trying to make AI ξ practically usable, some other AI methods, like genetic algorithms or neural nets, may be useful.

The main thing we wanted to point out is that the AI ξ model does not lack any important known property of intelligence or known AI methodology. What *is* missing, however, are computational aspects, which are addressed in the next chapter.

6.7 Problems

6.1 (Self-optimizingness) [C35uo] Formally define the environmental classes \mathcal{M}_{EC} for $EC \in \{\text{FMS}, \text{SGR}, \text{EX}\}$ similarly to $EC \in \{\text{SP}, \text{AI}\}$. \mathcal{M}_{EC} shall be the class of all (lower-semi)computable environments consistent with the problem setup EC, $\xi^{EC} := \sum_{\nu \in \mathcal{M}_{EC}} 2^{-K(\nu)} \nu$ the corresponding universal prior, and EC ξ alias p_{EC}^ξ the Bayes-optimal policy. Show that function minimization (FMS), repeated strategic games (SGR), and supervised learning (EX) admit self-optimizing policies (cf. Problem 5.10), hence FMS ξ , SGR ξ and EX ξ being self-optimizing. Interpret the results and compare them to the properties of SP ξ and AI ξ .

6.2 (Prediction loss bounds for AI ξ) [C40oi] In Section 6.2.2 we derived a bound (6.17) exponential in $K(\dot{z})$ on the number of prediction errors made by AI ξ with horizon $h_k = 1$ in deterministic passive environments. Try to generalize/improve this bound to (a) general loss functions, (b) bounds on $E_{n\xi}^{AI} - E_{n\mu}^{AI}$ for probabilistic passive environments, (c) the case $h_k > 1$, (d) bounds linear or “polynomial” in $K(\dot{z})$ as in the case of SP ξ – or – find examples demonstrating the impossibility of such generalizations/improvements.

6.3 (Posterization of prediction errors) [C20u/40o] Show $\xi^{AI}(\underline{y}_{1:n}) \stackrel{\times}{\geq} \xi^{SP}(\underline{z}_{1:n})$, where $z_k = \delta_{y_k x_k}$, and that the other direction $\stackrel{\times}{\leq}$ is wrong. Use this result to “improve” the bound (6.17) to $E_{n\xi}^{AI} \stackrel{\times}{\leq} [\xi^{SP}(\dot{z}_{1:n})]^{-1}$. Posterize this to a bound $E_{kn\xi}^{AI} \stackrel{\times}{\leq} \xi^{AI}(\underline{y}_{<k}) / \xi^{SP}(\dot{z}_{1:n})$ on the number of errors in cycles k to n . Is it possible to improve the numerator to $\xi^{SP}(\dot{z}_{<k})$ and to bound the expression by $\approx 2^{K(\dot{z}_{k:n} | \dot{z}_{<k})}$ (cf. Problem 3.13)?



John von Neumann
(1903-1957)

“Only math nerds would call 2^{500} finite.” (Leonid Levin)

“The biggest difference between time and space is that you can’t reuse time.” (Merrick Furst)

“The only reason for time is so that everything doesn’t happen at once.” (Albert Einstein / John Wheeler)

“You insist that there is something a machine cannot do. If you will tell me precisely what it is that a machine cannot do, then I can always make a machine which will do just that!” (John von Neumann)

Chapter 7

Computational Aspects

7.1	The Fastest & Shortest Algorithm for All Problems	702
7.1.1	Introduction & Main Result	702
7.1.2	Levin Search	704
7.1.3	Fast Matrix Multiplication	705
7.1.4	Applicability of the Fast Algorithm M_p^ε	706
7.1.5	The Fast Algorithm M_p^ε	707
7.1.6	Time Analysis	708
7.1.7	Assumptions on the Machine Model	710
7.1.8	Algorithmic Complexity and the Shortest Algorithm	710
7.1.9	Generalizations	712
7.1.10	Summary & Outlook	712
7.2	Time Bounded AIXI Model	713
7.2.1	Introduction	713
7.2.2	Time Limited Probability Distributions	714
7.2.3	The Idea of the Best Vote Algorithm	716
7.2.4	Extended Chronological Programs	716
7.2.5	Valid Approximations	717
7.2.6	Effective Intelligence Order Relation	717
7.2.7	The Universal Time Bounded AIXI \tilde{t} Agent	718
7.2.8	Limitations and Open Questions	719
7.2.9	Remarks	720

Up to now we have shown the universal character of the AI ξ model but have completely ignored computational aspects, which we make up for in this chapter.

We start in Section 7.1 by developing an algorithm M that is capable of solving any well-defined problem p as quickly as the fastest algorithm computing a solution

to p , save for a factor of $1+\varepsilon$ and lower-order additive terms. M optimally distributes resources between the execution of provably correct p -solving programs and an enumeration of all proofs, including relevant proofs of program correctness and of time bounds on program runtimes. The solution is somewhat involved from an implementational aspect. An implementation would include first order logic, the definition of a Universal Turing machine within it and proof theory. M avoids Blum's speed-up theorem by ignoring programs without correctness proof. M has broader applicability and can be faster than Levin's universal search, the fastest method for inverting functions save for a large multiplicative constant. An extension of Kolmogorov complexity and two novel natural measures of function complexity are used to show that the most efficient program computing some function f is also among the shortest programs provably computing f .

Based on a similar idea we construct in Section 7.2 a computable version of the AI ξ model. Let us assume that there exists some algorithm \tilde{p} of size \tilde{l} with computation time per interaction cycle \tilde{t} , which behaves in a sufficiently intelligent way (this assumption is the very basis of AI). The algorithm p^* should run all algorithms of length $\leq \tilde{l}$ for \tilde{t} time steps in every cycle and select the best output among them. So we have an algorithm which runs in time $\tilde{t} \cdot 2^{\tilde{l}}$ and is at least as good as \tilde{p} , i.e. it also serves our needs apart from the (very large but) constant multiplicative factor in computation time. This idea of the 'typing monkeys', one of them eventually producing 'Shakespeare', is well known and widely used in theoretical computer science. The difficult part here is the selection of the algorithm with the best output. A further complication is that the selection process itself must have only limited computation time. We present a suitable modification of the AI ξ model which solves these difficult problems. The assumptions behind this construction are discussed at the end.

7.1 The Fastest & Shortest Algorithm for All Well-Defined Problems

7.1.1 Introduction & Main Result

Searching for fast algorithms to solve certain problems is a central and difficult task in computer science. Positive results usually come from explicit constructions of efficient algorithms for specific problem classes. A wide class of problems can be phrased in the following way. Given a formal specification of a problem depending on some parameter $x \in X$, we are interested in a fast algorithm computing solution $y \in Y$. This means that we are interested in a fast algorithm computing $f: X \rightarrow Y$, where f is a formal (logical, mathematical, not necessarily algorithmic), specification of the problem. Ideally, we would like to have the fastest algorithm, maybe apart from some small constant factor in computation time. Unfortunately, Blum's Speed-up Theorem [Blu67, Blu71] shows that there are problems for which an (incomputable)

sequence of speed-improving algorithms (of increasing size) exists, but no fastest algorithm.

In the approach presented here, we consider only those algorithms which *provably* solve a given problem, and have a fast (i.e. quickly computable) time bound. Neither the programs themselves, nor the proofs need to be known in advance. Under these constraints we construct the asymptotically fastest algorithm save a factor of $1+\varepsilon$ that solves any well-defined problem f .

Theorem 7.1 (The fastest algorithm) Let p^* be a given algorithm computing $p^*(x)$ from x , or, more generally, a specification of a function. Let p be any algorithm, computing provably the same function as p^* with computation time provably bounded by the function $t_p(x)$ for all x . $time_{t_p}(x)$ is the time needed to compute the time bound $t_p(x)$. Fix some $\varepsilon \in (0, \frac{1}{2})$. Then the algorithm $M_{p^*}^\varepsilon$ constructed in Section 7.1.5 computes $p^*(x)$ in time

$$time_{M_{p^*}^\varepsilon}(x) \leq (1+\varepsilon) \cdot t_p(x) + \frac{d_p}{\varepsilon} \cdot time_{t_p}(x) + \frac{c_p}{\varepsilon}$$

with constants c_p and d_p depending on p but not on x . Neither p , t_p , nor the proofs need to be known in advance for the construction of $M_{p^*}^\varepsilon$.

Known time bounds for practical problems can often be computed quickly, i.e. $time_{t_p}(x)/time_p(x)$ often converges very quickly to zero. Furthermore, from a practical point of view, the provability restrictions are often rather weak. Hence, we have constructed for all those problems a solution, which is asymptotically only a factor $1+\varepsilon$ slower than the (provably) fastest algorithm! There is no large multiplicative factor and the problems are not restricted to inversion problems, as in Levin's algorithm (see Section 7.1.2). What somewhat spoils the practical applicability of $M_{p^*}^\varepsilon$ is the large additive constant c_p , which will be estimated in Section 7.1.6.

An interesting and counter-intuitive consequence of Theorem 7.1, derived in Section 7.1.8, is that the fastest program that computes a certain function is also among the shortest programs that provably computes this function. Looking for larger programs saves at most a finite number of computation steps, but cannot improve the time order.

In Section 7.1.2 we review Levin search and the universal search algorithms SIMPLE and SEARCH, described in [LV97]. We point out that SIMPLE has the same asymptotic time complexity as SEARCH not only w.r.t. the problem instance, but also w.r.t. to the problem class. In Section 7.1.3 we elucidate Theorem 7.1 and the applicability to an example problem (matrix multiplication) unsolvable by Levin search. Section 7.1.4 discusses the general applicability of $M_{p^*}^\varepsilon$. In Section 7.1.5 we give formal definitions of the expressions *time*, *proof*, *compute*, etc., which occur in Theorem 7.1, and define the fast algorithm $M_{p^*}^\varepsilon$. In Section 7.1.6 we analyze the algorithm $M_{p^*}^\varepsilon$, especially its computation time, prove Theorem 7.1, and give upper bounds for the constants c_p and d_p . Subtleties regarding the underlying machine

model are briefly discussed in Section 7.1.7. In Section 7.1.8 we show that the fastest program computing a certain function is also among the shortest programs provably computing this function. For this purpose, we extend the definition of the Kolmogorov complexity of a string and define two new natural measures for the complexity of functions and programs. Section 7.1.9 outlines generalizations of Theorem 7.1 to I/O streams and other time-measures. Conclusions are given in Section 7.1.10.

7.1.2 Levin Search

Levin search is one of the few rather general speed-up algorithms. Within a (typically large) factor, it is the fastest algorithm for inverting a function $g: Y \rightarrow X$, if g can be evaluated quickly [Lev73b, Lev84]. Given x , an inversion algorithm p tries to find a $y \in Y$, called g -witness for x , with $g(y) = x$. Levin search just runs and verifies the result of *all* algorithms p in parallel with relative computation time $2^{-l(p)}$; i.e. a time fraction $2^{-l(p)}$ is devoted to execute p , where $l(p)$ is the length of program p (coded in binary). Verification is necessary since the output of *any* program can be *anything*. This is the reason why Levin search is only effective if a fast implementation of g is available. Levin search halts if the first g -witness has been produced and verified. The total computation time to find a solution (if one exists) is bounded by $2^{l(p)} \cdot \text{time}_p^+(x)$. $\text{time}_p^+(x)$ is the runtime of $p(x)$ *plus* the time to verify the correctness of the result ($g(p(x)) = x$) by a *known* implementation for g .

Li and Vitányi [LV97, p503] propose a very simple variant, called SIMPLE, which runs all programs $p_1 p_2 p_3 \dots$ one step at a time according to the following scheme: p_1 is run every second step, p_2 every second step in the remaining unused steps, p_3 every second step in the remaining unused steps, and so forth, i.e. according to the sequence of indices 121312141213121512.... If p_k inverts g on x in $\text{time}_{p_k}(x)$ steps, then SIMPLE will do the same in *at most* $2^k \text{time}_{p_k}^+(x) + 2^{k-1}$ steps. In order to improve the factor 2^k , they define the algorithm SEARCH, which runs all p (of length less than i) for $2^i 2^{-l(p)}$ steps in phase $i = 1, 2, 3, \dots$, until it has inverted g on x . The computation time of SEARCH is bounded by $2^{K(k)+O(1)} \text{time}_{p_k}^+(x)$, where $K(k) \leq l(p_k) \leq 2 \log_2 k$ is the Kolmogorov complexity of k . They suggest that SIMPLE has worse asymptotic behavior w.r.t. k than SEARCH, but actually this is not the case.

In fact, SIMPLE and SEARCH have the same asymptotics also in k , because SEARCH itself is an algorithm with some index $k_{\text{SEARCH}} = O(1)$. Hence, SIMPLE executes SEARCH every $2^{k_{\text{SEARCH}}}$ -th step, and can at most be a constant (independent of k and x) factor $2^{k_{\text{SEARCH}}} = O(1)$ slower than SEARCH. However, in practice, SEARCH should be favored, because also constants matter, and $2^{k_{\text{SEARCH}}} \approx 2^{2^{l(p_{k_{\text{SEARCH}}})}}$ is rather large.

Levin search can be modified to handle time-limited optimization problems as well [Sol86]. Many, but not all problems, are of inversion or optimization type. The

matrix multiplication example (Section 7.1.4), the *decision* problem SAT [LV97, p503], and reinforcement learning [Hut01d], for instance, are not of this form. Furthermore, the large factor $2^{l(p)}$ somewhat limits the applicability of Levin search.

Levin search in program space cannot be used directly in $M_{p^*}^\varepsilon$ for computing p^* because we have to decide somehow whether a certain program solves our problem or computes something else. For this, we have to search through the space of proofs. In order to avoid the large time-factor $2^{l(p)}$, we also have to search through the space of time-bounds. Only *one* (fast) program should be executed for a significant time interval. The algorithm $M_{p^*}^\varepsilon$ essentially consists of 3 interwoven algorithms: *sequential* program execution, sequential search through proof space, and Levin search through time-bound space. A tricky scheduling prevents performance degradation from computing slow p 's before *the* p has been found.

7.1.3 Fast Matrix Multiplication

To illustrate Theorem 7.1, we consider the problem of multiplying two $n \times n$ matrices. If p^* is the standard algorithm for multiplying two matrices¹ $x \in R^{n \times n} \times R^{n \times n}$ of size $l(x) \sim n^2$, then $t_{p^*}(x) := 2n^3$ upper bounds the true computation time $time_{p^*}(x) = n^2(2n-1)$. We know there exists algorithms p' for matrix multiplication with $time_{p'}(x) \leq t_{p'}(x) := c \cdot n^\omega$ ($\omega = 2.81$ [Str69], $\omega = 2.50$ [CW82], $\omega = 2.38$ [CW90], ...). The time-bound function (cast to an integer) can, as in many cases, be computed very quickly, $time_{t_{p'}}(x) = O(\log^2 n)$. Hence, using Theorem 7.1, also $M_{p^*}^\varepsilon$ is fast, $time_{M_{p^*}^\varepsilon}(x) \leq (1+\varepsilon)c \cdot n^\omega + O(\log^2 n)$. Of course, $M_{p^*}^\varepsilon$ would be of no real use if p' is already the fastest program, since p' is known and could be used directly. We do not know however, whether there is an algorithm p'' with $time_{p''}(x) \leq d \cdot n^2 \log n$, for instance. But if it does exist, $time_{M_{p^*}^\varepsilon}(x) \leq (1+\varepsilon)d \cdot n^2 \log n + O(1)$ for all x is guaranteed.

There is no contradiction to [CW82] who show that there is no fastest bilinear λ -algorithm (b λ A) for matrix multiplication. For every b λ A p_i with computation time n^{ω_i} one can find another b λ A p_{i+1} with computation time $n^{\omega_{i+1}}$ and $\omega_{i+1} < \omega_i$, but there is no b λ A with computation time n^{ω_0} and $\omega_0 = \inf_i \{\omega_i\}$. On the other hand, this says nothing about the existence of a non-b λ A M with computation time of, for instance, $n^{\omega_0} \log n$, which is faster than all b λ A p_i . Indeed a formal construction of such an algorithm is easy. The sequence $\{p_1, p_2, p_3, \dots\}$ is enumerable, i.e. there is an algorithm which creates the programs p_1, p_2, p_3, \dots , say in time $\tau_1, \tau_2, \tau_3, \dots$. We enumerate p_1, p_2, p_3, \dots and start executing them in parallel as soon as they have been constructed and assign a fraction $\frac{1}{i(i+1)}$ of time to $p_i(x)$. The first p_i which halts outputs the result. The total computation time of this (meta)algorithm M is

$$time_M(x) = \min_i \{\tau_i + i(i+1)time_{p_i}(x)\} = O(time_{p_i}(x)) \forall i$$

¹Instead of interpreting R as the set of real numbers one might take the field $IF_2 = \{0,1\}$ to avoid subtleties arising from large numbers. Arithmetic operations are assumed to need one unit of time.

Hence M has better time complexity than any of the p_i . For instance, for $\omega_i = \omega_0 + O(i^{-2})$ and τ_i polynomial in i , it is easy to see that $\text{time}_M(x) = O(n^{\omega_0} \log n)$. The construction above works in general as long as the program sequence is enumerable. It fails for incomputable sequences, like in Blum's speed-up construction.

The matrix multiplication example has been chosen for specific reasons. First, it is not an inversion or optimization problem directly suitable for Levin search. The computation time of Levin search is lower-bounded by the time to verify the solution with a known algorithm (which is currently $c \cdot n^{2.376\dots}$) multiplied with the (large) number of necessary verifications. Second, although matrix multiplication is a very important and time-consuming issue, p' is not used in practice, since c is so large that for all practically occurring n , the cubic algorithm is faster. The same is true for c_p and d_p , but we must admit that although c is large, the bounds we obtain for c_p and d_p are tremendous. On the other hand, even Levin search, which has a tremendous multiplicative factor, can be successfully applied [Sch97, SZW97, Sch02b], when handled with care. The same should hold for Theorem 7.1, as will be discussed. We avoid the $O()$ notation as far as possible, as it can be severely misleading (e.g. $10^{42} = O(1)^{O(1)} = O(1)$). This chapter could be viewed as another $O()$ warning showing, how important factors, and even subdominant additive terms, are.

7.1.4 Applicability of the Fast Algorithm $M_{p^*}^\varepsilon$

An obvious time bound for p is the actual computation time itself. An obvious algorithm to compute $\text{time}_p(x)$ is to count the number of steps needed for computing $p(x)$. Hence, inserting $t_p = \text{time}_p$ into Theorem 7.1 and using $\text{time}_{\text{time}_p}(x) \leq \text{time}_p(x)$, we see that the computation time of $M_{p^*}^\varepsilon$ is optimal within a multiplicative constant $(d_p + 1 + \varepsilon)$ and an additive constant c_p . This result is weaker than the one in Theorem 7.1, but no assumption concerning the computability of time bounds has to be made.

When do we trust that a fast algorithm solves a given problem? At least for well specified problems, like satisfiability, solving a combinatoric puzzle, computing the digits of π , ..., we usually invent algorithms, prove that they solve the problem and in many cases also can prove good and quickly computable time bounds. In these cases, the provability assumptions in Theorem 7.1 are no real restriction. The same holds for approximate algorithms which guarantee a precision ε within a known time bound (many numerical algorithms are of this kind). For exact/approximate programs provably computing/converging to the right answer (e.g. traveling salesman problem, and also many numerical programs), but for which no good, and easy to compute time bound exists, M_{p^*} is only optimal apart from a huge constant factor $1 + \varepsilon + d_p$ in time, as discussed above. Universal reinforcement learning could be a problem of this kind. There is no known efficient algorithm for computing the optimal policy for sequential decision problems in non-Markovian environments. The algorithm $\text{AI}\xi^{tl}$ developed in Section 7.2 is based on a similar idea as the $M_{p^*}^\varepsilon$. It creates an incremental policy for an agent in an unknown non-Markovian environment, which is superior to any other time t and length l bounded agent. The computation time

of $\text{AI}\xi^{tl}$ is of the order $t \cdot 2^l$. For poorly specified problems, Theorem 7.1 does not help at all.

7.1.5 The Fast Algorithm $M_{p^*}^\varepsilon$

One ingredient of algorithm $M_{p^*}^\varepsilon$ is an enumeration of proofs of increasing length in some formal axiomatic system. If a proof actually proves that p and p^* are functionally equivalent and p has time bound t_p , (p, t_p) is added to a list L . The program p in L with the currently smallest time bound $t_p(x)$ is executed. By construction, the result $p(x)$ is identical to $p^*(x)$. The trick to achieve the time bound stated in Theorem 1 is to schedule everything in a proper way, in order not to lose too much performance by computing slow p 's and t_p 's before *the* p has been found.

To avoid confusion, we formally define p and t_p to be binary strings. That is, p is neither a program nor a function, but can be informally interpreted as such. A formal definition of the interpretations of p is given below. We say “ p computes function f ”, when a universal reference Turing machine U on input (p, x) computes $f(x)$ for all x . This is denoted by $U(p, x) = f(x)$. To be able to talk about proofs, we need a formal logic system $(\forall, \lambda, y_i, c_i, f_i, R_i, \rightarrow, \wedge, =, \dots)$, and axioms, and inference rules. A proof is a sequence of formulas, where each formula is either an axiom or inferred from previous formulas in the sequence by applying the inference rules. See [Fit96, Sho67] or any other textbook on logic or proof theory. We only need to know that *provability*, *Turing Machines*, and *computation time* can be formalized:

1. The set of (correct) proofs is enumerable.
2. A term u can be defined such that the formula $[\forall y : u(p, y) = u(p^*, y)]$ is true if, and only if $U(p, x) = U(p^*, x)$ for all x , i.e. if p and p^* describe the same function.
3. A term tm can be defined such that the formula $[tm(p, x) = n]$ is true if, and only if the computation time of U on (p, x) is n , i.e. if $n = \text{time}_p(x)$.

We say that p is provably equivalent to p^* if the formula $[\forall y : u(p, y) = u(p^*, y)]$ can be proven. $M_{p^*}^\varepsilon$ runs three algorithms A , B , and C in parallel:

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

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

Start algorithms A , B , and C in parallel with relative computational resources ε , ε , and $1 - 2\varepsilon$, respectively.

That is, C performs $\frac{1}{\varepsilon}$ steps when A and B perform 1 step each.

Algorithm A

```

for  $i:=1,2,3,\dots$  do
  pick the  $i^{th}$  proof in the list of all proofs and
  isolate the last formula in the proof.
  if this formula is equal to  $[\forall y: u(p,y)=u(p^*,y) \wedge u(t,y) \geq tm(p,y)]$ 
  for some strings  $p$  and  $t$ ,
  then add  $(p,t)$  to  $L$ .
next  $i$ 

```

Algorithm B

```

for all  $(p,t) \in L$ 
  run  $U$  on all  $(t,x)$  in parallel for all  $t$  with relative computational resources
   $2^{-l(p)-l(t)}$ .
  if  $U$  halts for some  $t$  and  $U(t,x) < t_{fast}$ ,
  then  $t_{fast} := U(t,x)$  and  $p_{fast} := p$  and restart algorithm  $C$ .
continue  $(p,t)$ 

```

Algorithm C

```

run  $U$  on  $(p_{fast},x)$ . For each executed step decrease  $t_{fast}$  by 1.
if  $U$  halts then print result  $U(p_{fast},x)$  and abort computation of  $A$ ,  $B$  and  $C$ .

```

Note that A and B only terminate when aborted by C . The discussion of the algorithm(s) in the following sections clarifies details and proves Theorem 7.1.

7.1.6 Time Analysis

Henceforth we return to the convenient abbreviations $p(x) := U(p,x)$ and $t_p(x) := U(t_p,x)$. Let p' be some fixed algorithm that is provably equivalent to p^* , with computation time $time_{p'}$ provably bounded by $t_{p'}$. Let $l(\text{proof}(p'))$ be the length of the binary coding of the, for instance, shortest proof. *Computation time* always refers to true overall computation time, whereas *computation steps* refer to instruction steps. $steps = \alpha \cdot time$, if a percentage α of computation time is assigned to an algorithm.

A) To write down (not to invent!) a proof requires $O(l(\text{proof}))$ steps. A time $O(N_{axiom} \cdot l(F_i))$ is needed to check whether a formula F_i in the proof $F_1 F_2 \dots F_n$ is an axiom, where N_{axiom} is the number of axioms or axiom-schemes, which is finite. Variable substitution (binding) can be performed in linear time. For a suitable set of axioms, the only necessary inference rule is modus ponens. If F_i is not an axiom, one searches for a formula F_j , $j < i$ of the form $F_k \rightarrow F_i$ and then for the formula F_k , $k < i$. This takes time $O(l(\text{proof}))$. There are $n \leq O(l(\text{proof}))$ formulas F_i to check in this way. Whether the sequence of formulas constitutes a valid proof

can, hence, be checked in $O(l(\text{proof})^2)$ steps. There are less than 2^{l+1} proofs of (binary) length $\leq l$. Algorithm A receives a fraction ε of relative computation time. Hence, for a proof of $(p', t_{p'})$ to occur, and for $(p', t_{p'})$ to be added to L , at most time $T_A \leq \frac{1}{\varepsilon} \cdot 2^{l(\text{proof}(p'))+1} \cdot O(l(\text{proof}(p'))^2)$ is needed. Note that the same program p can and will be accompanied by different time bounds t_p ; for instance (p, time_p) will occur.

B) The time assignment of algorithm B to the t_p 's only works if the Kraft inequality $\sum_{(p, t_p) \in L} 2^{-l(p)-l(t_p)} \leq 1$ is satisfied [Kra49]. This can be ensured by using prefix free (e.g. Shannon-Fano) codes [Sha48, LV97]. The number of steps to calculate $t_{p'}(x)$ is, by definition, $\text{time}_{t_{p'}}(x)$. The relative computation time available for computing $t_{p'}(x)$ is $\varepsilon \cdot 2^{-l(p')-l(t_{p'})}$. Hence, $t_{p'}(x)$ is computed and $t_{fast} \leq t_{p'}(x)$ is checked after time $T_B \leq T_A + \frac{1}{\varepsilon} \cdot 2^{l(p')+l(t_{p'})} \cdot \text{time}_{t_{p'}}(x)$. We have to add T_A , since B has to wait, in the worst case, time T_A before it can start executing $t_{p'}(x)$.

C) If algorithm C halts, its construction guarantees that the output is correct. In the following, we show that C always halts, and give a bound for the computation time.

- i) Assume that algorithm C stops before B performed the check $t_{p'}(x) < t_{fast}$, because a different p already computed $p(x)$. In this case $T_C \leq T_B$.
- ii) Assume that B performs the check $t_{p'}(x) < t_{fast}$ and the check succeeds. Running-time T_B has passed until this point. C is restarted computes $p_{fast}(x) = p'(x)$ in time $t_{fast} := t'_{p'}$, or faster, if during the computation, p_{fast} gets replaced by an even faster algorithm constructed by A and B (t_{fast} is a decreasing variable). Since a fraction $1 - 2\varepsilon$ of relative computation time is assigned to C it halts after time $T_C \leq T_B + \frac{1}{1-2\varepsilon} t_{p'}(x)$.
- iii) At any point in time the remaining time until C halts is bounded by $\frac{1}{1-2\varepsilon} t_{fast}$, since t_{fast} is never increasing. Hence, if the check $t_{p'}(x) < t_{fast}$ fails, $T_C \leq T_B + \frac{1}{1-2\varepsilon} t_{fast} \leq T_B + \frac{1}{1-2\varepsilon} t_{p'}(x)$.

The maximum of the cases (i) to (iii) bounds the computation time of C and, hence, of $M_{p^*}^\varepsilon$ by

$$\text{time}_{M_{p^*}^\varepsilon}(x) = T_C \leq T_B + \frac{1}{1-2\varepsilon} t_p(x) \leq (1 + 3\varepsilon) \cdot t_p(x) + \frac{d_p}{3\varepsilon} \cdot \text{time}_{t_p}(x) + \frac{c_p}{3\varepsilon},$$

$$d_p = 3 \cdot 2^{l(p)+l(t_p)}, \quad c_p = 3 \cdot 2^{l(\text{proof}(p))+1} \cdot O(l(\text{proof}(p))^2),$$

where we have dropped the prime from p and used $\frac{1}{1-2\varepsilon} \leq 1 + 3\varepsilon$ for $\varepsilon \leq \frac{1}{6}$. We have also suppressed the dependency of c_p and d_p on p^* ($\text{proof}(p)$ depends on p^* too), since we considered p^* to be a fixed given algorithm. Rescaling $\varepsilon \rightsquigarrow \varepsilon/3$ leads to the bound in Theorem 7.1.

7.1.7 Assumptions on the Machine Model

In the time analysis above we have assumed that program simulation with abort possibility and scheduling parallel algorithms can be performed in real-time, i.e. without loss of performance. Parallel computation can be avoided by sequentially performing time slices of N operations and then switching to the next task. Algorithms A and C , and every $(p, t) \in L$ in algorithm B constitute a task. If switching between time slices needs constant time s and we choose $N \sim \frac{1}{s\varepsilon}$, then time slicing increases computation time by a factor $1 + \varepsilon$. Also, in order to avoid a possible slowdown of p in algorithm C due to decrementing t_{fast} , one should decrement t_{fast} by N every N 'th time step, possibly synchronously to the task switching. Counting can be performed in time $O(1)$ [SV88].

A thorough construction of a real-time machine U goes beyond the scope of this paper. The above discussion should be a motivation that universal real-time machines U are something reasonable. Note that we use the same universal Turing machine U with the same underlying Turing machine model (number of heads, symbols, ...) for measuring computation time of all programs (strings) p , including M_p^ε . This prevents us from applying the linear speedup theorem (which is cheating somewhat anyway), but allows the possibility of designing a U which allows real-time simulation with abort possibility. Theorem 7.1 should also hold for suitable Kolmogorov-Uspenskii [KU63] and Pointer machines [Sch80].

7.1.8 Algorithmic Complexity and the Shortest Algorithm

Data compression is a very important issue in computer science. Saving space or channel capacity are obvious applications. In Chapter 2 we have seen that a less obvious (but not far fetched) application is that of inductive inference in various forms (hypothesis testing, forecasting, classification, ...). A free interpretation of Occam's razor is that the shortest theory consistent with past data is the most likely to be correct. This has been put into a rigorous scheme by [Sol64] and proved to be optimal in Chapter 3. Kolmogorov complexity is a universal notion of the information content of a string [Kol65, Cha66, ZL70]. It is defined as the length of the shortest program computing string x .

$$K_U(x) := \min_p \{l(p) : U(p) = x\} = K(x) + O(1)$$

where U is some universal Turing Machine. It can be shown that $K_U(x)$ varies, at most, by an additive constant independent of x by varying the machine U . Hence, *the* Kolmogorov Complexity $K(x)$ is universal in the sense that it is uniquely defined up to an additive constant. $K(x)$ can be approximated from above (is co-enumerable), but is not finitely computable. See [LV97] for an excellent introduction to Kolmogorov Complexity and [VL00] for a review of Kolmogorov inspired prediction schemes.

Recently, Schmidhuber [Sch00, Sch02a] has generalized Kolmogorov complexity in various ways to the limits of computability and beyond. In the following, we also need a generalization, but of a different kind. We need a short description of a function, rather than a string. The following definition of the complexity of a function f

$$K'(f) := \min_p \{l(p) : U(p, x) = f(x) \forall x\}$$

seems natural, but suffers from not even being approximable. There exists no algorithm converging to $K'(f)$, because it is in general undecidable whether a program p is equivalent to (some formal definition of) a function f . Even if we have a program p^* computing f , $K'(p^*)$ is not approximable. Using $K(p^*)$ is not a suitable alternative, since $K(p^*)$ might be considerably larger than $K'(p^*)$, as in the former case all information conveyed by p^* will be kept – even that which is functionally irrelevant (e.g. dead code). An alternative is to restrict ourselves to provably equivalent programs. The length of the shortest one is

$$K''(p^*) := \min_p \{l(p) : \text{a proof of } [\forall y: u(p, y) = u(p^*, y)] \text{ exists}\}$$

It can be approximated from above, since the set of all programs provably equivalent to p^* is enumerable.

Having obtained, after some time, a very short description p' of p^* for some purpose (e.g. for defining a prior probability for some inductive inference scheme), it is usually also necessary to obtain values for some arguments. We are now concerned with the computation time of p' . Could we get slower and slower algorithms by compressing p^* more and more? Interestingly this is not the case. Inventing complex (long) programs is *not* necessary to construct asymptotically fast algorithms, under the stated provability assumptions, in contrast to Blum's Theorem [Blu67, Blu71]. The following theorem roughly says that there is a *single* program, which is the fastest *and* the shortest program.

Theorem 7.2 (The fastest & shortest algorithm) Let p^* be a given algorithm or formal specification of a function. There exists a program \tilde{p} , equivalent to p^* , for which the following holds

$$i) \quad l(\tilde{p}) \leq K''(p^*) + O(1)$$

$$ii) \quad \text{time}_{\tilde{p}}(x) \leq (1+\varepsilon) \cdot t_p(x) + \frac{d_p}{\varepsilon} \cdot \text{time}_{t_p}(x) + \frac{c_p}{\varepsilon}$$

where p is any program provably equivalent to p^* with computation time provably less than $t_p(x)$. The constants c_p and d_p depend on p but not on x .

To prove the theorem, we just insert the shortest algorithm p' provably equivalent to p^* into M , that is $\tilde{p} := M_{p'}^\varepsilon$. As only $O(1)$ instructions are needed to build $M_{p'}^\varepsilon$ from p' , $M_{p'}^\varepsilon$ has size $l(p') + O(1) = K''(p^*) + O(1)$. The computation time of $M_{p'}^\varepsilon$ is the same as of $M_{p^*}^\varepsilon$ apart from “slightly” different constants c_p and d_p .

The following subtlety has been pointed out by Peter van Emde Boas. Neither $M_{p^*}^\varepsilon$, nor \tilde{p} is *provably* equivalent to p^* . The construction of $M_{p^*}^\varepsilon$ in Section 7.1.5 shows equivalence of $M_{p^*}^\varepsilon$ (and of \tilde{p}) to p^* , but it is a meta-proof which cannot be formalized within the considered proof system. A formal proof of the correctness of $M_{p^*}^\varepsilon$ would prove the consistency of the proof system, which is impossible by Gödel's second incompleteness theorem. See [Har79] for details in a related context.

7.1.9 Generalizations

If p^* has to be evaluated repeatedly, algorithm A can be modified to remember its current state and continue operation for the next input (A is independent of $x!$). The large offset time c_p is only needed on the first invocation.

$M_{p^*}^\varepsilon$ can be modified to handle I/O streams, definable by a Turing machine with monotone input and output tapes (and bidirectional work tapes) receiving an input stream and producing an output stream. The currently read prefix of the input stream is x . $\text{time}_p(x)$ is the time used for reading x . $M_{p^*}^\varepsilon$ caches the input and output streams, so that algorithm C can repeatedly read/write the streams for each new p . The true input/output tapes are used for requesting/producing a new symbol. Algorithm B is reset after 1,2,4,8,... steps (not after reading the next symbol of $x!$) to appropriately take into account increased prefixes x . Algorithm A just continues. The bound of Theorem 7.1 holds for this case too, with slightly increased d_p .

The construction above also works if time is defined as a function of the current output rather than the current input x . This measure is, for example, used for the time-complexity of calculating the n^{th} digit of a computable real (e.g. π), where there is no input, but only an output stream.

7.1.10 Summary & Outlook

We presented an algorithm $M_{p^*}^\varepsilon$ which accelerates the computation of a program p^* . $M_{p^*}^\varepsilon$ combines (A) sequential search through proof space, (B) Levin search through time-bound space, (C) and *sequential* program execution, using a somewhat tricky scheduling. Under certain provability constraints, $M_{p^*}^\varepsilon$ is the asymptotically fastest algorithm for computing p^* apart from a factor $1+\varepsilon$ in computation time. Blum's Theorem shows that the provability constraints are essential. We have shown that the conditions on Theorem 7.1 are often, but not always, satisfied for practical problems. For complex approximation problems, for instance, where no good and quickly computable time bound exists, $M_{p^*}^\varepsilon$ is still optimal, but in this case, only apart from a large multiplicative factor. We briefly outlined how $M_{p^*}^\varepsilon$ can be modified to handle I/O streams and other time-measures. An interesting and counter-intuitive consequence of Theorem 7.1 was that the fastest program computing a certain function is also among the shortest programs provably computing this function. Looking for larger programs saves at most a finite number of computation steps, but cannot

improve the time order. To quantify this statement, we extended the definition of Kolmogorov complexity and defined two novel natural measures for the complexity of a function. The large constants c_p and d_p seem to spoil a direct implementation of $M_{p^*}^\varepsilon$. On the other hand, Levin search has been successfully adapted/generalized and applied to solve rather difficult machine learning problems [Sch97, SZW97, Sch02b], even though it suffers from a large multiplicative factor of similar origin. The use of more elaborate theorem-provers, rather than brute force enumeration of all proofs, could lead to smaller constants and bring M_p^* closer to practical applications, possibly restricted to subclasses of problems [RV01]. A more fascinating (and more speculative) way may be the utilization of so called transparent or holographic proofs [Bab91]. The correctness of these proofs can be checked by only reading a logarithmic number of their bits. This would mean that exponentially many proofs are checked simultaneously, reducing the constants c_p and d_p to their logarithm. I would like to conclude with a general question. Will the ultimate search for asymptotically fastest programs typically lead to fast or slow programs for arguments of practical size? Levin search, matrix multiplication and the algorithm $M_{p^*}^\varepsilon$ seem to support the latter, but this might be due to our inability to do better.

7.2 Time Bounded AIXI Model

7.2.1 Introduction

Until now, we have not bothered with the non-computability of the universal probability distribution ξ . As all universal models in this paper are based on ξ , they are not effective in this form. In this section, we outline how the previous models and results can be modified/generalized to the time-bounded case. Indeed, the situation is not as bad as it could be. ξ is enumerable and \dot{y}_k is still approximable, i.e. there exists an algorithm that will produce a sequence of outputs eventually converging to the exact output \dot{y}_k , but we can never be sure whether we have already reached it. Besides this, the convergence is extremely slow, so this type of asymptotic computability is of no direct (practical) use, but will nevertheless, be important later.

Let \tilde{p} be a program which calculates within a reasonable time \tilde{t} per cycle, a reasonable intelligent output, i.e. $\tilde{p}(\dot{x}_{<k}) = \dot{y}_{1:k}$. This sort of computability assumption, that a general purpose computer of sufficient power is able to behave in an intelligent way, is the very basis of AI, justifying the hope to be able to construct agents which eventually reach and outperform human intelligence. For a contrary viewpoint see [Luc61, Pen89, Pen94]. It is not necessary to discuss here, what is meant by ‘reasonable time/intelligence’ and ‘sufficient power’. What we are interested in, in this section, is whether there is a computable version $\text{AI}\xi^{\tilde{t}}$ of the $\text{AI}\xi$ agent which is superior or equal to any p with computation time per cycle of at most \tilde{t} . By ‘superior’, we mean ‘more intelligent’, so what we need is an order relation (like) (5.14) for intelligence.

The best result we could think of would be an $\text{AI}\xi^{\tilde{t}}$ with computation time $\leq \tilde{t}$

at least as intelligent as any p with computation time $\leq \tilde{t}$. If AI is possible at all, we would have reached the final goal, the construction of the most intelligent algorithm with computation time $\leq \tilde{t}$. Just as there is no universal measure in the set of computable measures (within time \tilde{t}), such an $\text{AI}\xi^{\tilde{t}}$ may neither exist.

What we can realistically hope to construct, is an $\text{AI}\xi^{\tilde{t}}$ agent of computation time $c \cdot \tilde{t}$ per cycle for some constant c . The idea is to run all programs p of length $\leq \tilde{l} := l(\tilde{p})$ and time $\leq \tilde{t}$ per cycle and pick the best output. The total computation time is $c \cdot \tilde{t}$ with $c = 2^{\tilde{l}}$. This sort of idea of ‘typing monkeys’ with one of them eventually writing Shakespeare, has been applied in various forms and contexts in theoretical computer science. The realization of this *best vote* idea, in our case, is not straightforward and will be outlined in this section. An idea related to this, is that of basing the decision on the majority of algorithms. This ‘democratic vote’ idea has been used in [LW94, Vov92] for sequence prediction, and is referred to as ‘weighted majority’ there.

7.2.2 Time Limited Probability Distributions

In the literature one can find time limited versions of Kolmogorov complexity [Dal73, Dal77, Ko86] and time limited universal semimeasures [LV91, LV97, Sch02c]. In the following, we utilize and adapt the latter and see how far we get. One way to define a time-limited universal chronological semimeasure is as a sum over all enumerable chronological semimeasures similar to the unbounded case (5.5) but computable within time \tilde{t} and of size at most \tilde{l} .

$$\xi^{\tilde{t}\tilde{l}}(\underline{y}_{1:n}) := \sum_{\rho : l(\rho) \leq \tilde{l} \wedge t(\rho) \leq \tilde{t}} 2^{-l(\rho)} \rho(\underline{y}_{1:n}) \quad (7.3)$$

Let us assume that the true environmental prior probability μ^{AI} is equal to or sufficiently accurately approximated by a ρ with $l(\rho) \leq \tilde{l}$ and $t(\rho) \leq \tilde{t}$ with \tilde{t} and \tilde{l} of reasonable size. There are several AI problems that fall into this class. In function minimization of Section 6.4, the computation of f and μ^{FM} are often feasible. In many cases, the sequences of Section 6.2 which should be predicted, can be easily calculated when μ^{SP} is known. In a classifier problem, the probability distribution μ^{CF} , according to which examples are presented, is, in many cases, also elementary. But not all AI problems are of this ‘easy’ type. For the strategic games of Section 6.3, the environment is usually, itself, a highly complex strategic player with a μ^{SG} that is difficult to calculate, although one might argue that the environmental player may have limited capabilities too. But it is easy to think of a difficult to calculate physical (probabilistic) environment like the chemistry of biomolecules.

The number of interesting applications makes this restricted class of AI problems, with time and space bounded environment $\mu^{\tilde{t}\tilde{l}}$, worth being studied. Superscripts to a probability distribution except for $\xi^{\tilde{t}\tilde{l}}$ indicate their length and maximal computation time. $\xi^{\tilde{t}\tilde{l}}$ defined in (7.3), with a yet to be determined computation time, multiplicatively dominates all $\mu^{\tilde{t}\tilde{l}}$ of this type. Hence, an $\text{AI}\xi^{\tilde{t}\tilde{l}}$ model, where we use

$\xi^{\tilde{t}\tilde{l}}$ as prior probability, is universal, relative to all $\text{AI}\mu^{\tilde{t}\tilde{l}}$ models in the same way as $\text{AI}\xi$ is universal to $\text{AI}\mu$ for all enumerable chronological semimeasures μ . The argmax_{y_k} in (5.3) selects a y_k for which $\xi^{\tilde{t}\tilde{l}}$ has the highest expected utility V_{km_k} , where $\xi^{\tilde{t}\tilde{l}}$ is the weighted average over the $\rho^{\tilde{t}\tilde{l}}$. $\dot{y}_k^{\text{AI}\xi^{\tilde{t}\tilde{l}}}$ is determined by a weighted majority. We expect $\text{AI}\xi^{\tilde{t}\tilde{l}}$ to outperform all (bounded) $\text{AI}\rho^{\tilde{t}\tilde{l}}$, analogous to the unrestricted case.

In the following we analyze the computability properties of $\xi^{\tilde{t}\tilde{l}}$ and $\text{AI}\xi^{\tilde{t}\tilde{l}}$, i.e. of $\dot{y}_k^{\text{AI}\xi^{\tilde{t}\tilde{l}}}$. To compute $\xi^{\tilde{t}\tilde{l}}$ according to the definition (7.3) we have to enumerate all chronological enumerable semimeasures $\rho^{\tilde{t}\tilde{l}}$ of length $\leq \tilde{l}$ and computation time $\leq \tilde{t}$. This can be done similarly to the unbounded case (5.42-5.44). All $2^{\tilde{l}}$ enumerable functions of length $\leq \tilde{l}$, computable within time \tilde{t} have to be converted to chronological probability distributions. For this, one has to evaluate each function for $|\mathcal{X}| \cdot k$ different arguments. Hence, $\xi^{\tilde{t}\tilde{l}}$ is computable within time² $t(\xi^{\tilde{t}\tilde{l}}(\underline{y}_{1:k})) = O(|\mathcal{X}| \cdot k \cdot 2^{\tilde{l}} \cdot \tilde{t})$. The computation time of $\dot{y}_k^{\text{AI}\xi^{\tilde{t}\tilde{l}}}$ depends on the size of \mathcal{X} , \mathcal{Y} and m_k . $\xi^{\tilde{t}\tilde{l}}$ has to be evaluated $|\mathcal{Y}|^{h_k} |\mathcal{X}|^{h_k}$ times in (5.3). It is possible to optimize the algorithm and perform the computation within time

$$t(\dot{y}_k^{\text{AI}\xi^{\tilde{t}\tilde{l}}}) = O(|\mathcal{Y}|^{h_k} |\mathcal{X}|^{h_k} \cdot 2^{\tilde{l}} \cdot \tilde{t}) \quad (7.4)$$

per cycle. If we assume that the computation time of $\mu^{\tilde{t}\tilde{l}}$ is exactly \tilde{t} for all arguments, the brute force time \bar{t} for calculating the sums and maxs in (4.17) is $\bar{t}(\dot{y}_k^{\text{AI}\mu^{\tilde{t}\tilde{l}}}) \geq |\mathcal{Y}|^{h_k} |\mathcal{X}|^{h_k} \cdot \tilde{t}$. Combining this with (7.4), we get

$$t(\dot{y}_k^{\text{AI}\xi^{\tilde{t}\tilde{l}}}) = O(2^{\tilde{l}} \cdot \bar{t}(\dot{y}_k^{\text{AI}\mu^{\tilde{t}\tilde{l}}}))$$

This result has the proposed structure, that there is a universal $\text{AI}\xi^{\tilde{t}\tilde{l}}$ agent with computation time $2^{\tilde{l}}$ times the computation time of a special $\text{AI}\mu^{\tilde{t}\tilde{l}}$ agent.

Unfortunately, the class of $\text{AI}\mu^{\tilde{t}\tilde{l}}$ systems with brute force evaluation of \dot{y}_k , according to (4.17) is completely uninteresting from a practical point of view. E.g. in the context of chess, the above result says that the $\text{AI}\xi^{\tilde{t}\tilde{l}}$ is superior within time $2^{\tilde{l}} \cdot \tilde{t}$ to any brute force minimax strategy of computation time \tilde{t} . Even if the factor of $2^{\tilde{l}}$ in computation time would not matter, the $\text{AI}\xi^{\tilde{t}\tilde{l}}$ agent is, nevertheless practically useless, as a brute force minimax chess player with reasonable time \tilde{t} is a very poor player.

Note, that in the case of binary sequence prediction ($h_k = 1$, $|\mathcal{Y}| = |\mathcal{X}| = 2$) the computation time of ρ coincides with that of $\dot{y}_k^{\text{AI}\rho}$ within a factor of 2. The class $\text{AI}\rho^{\tilde{t}\tilde{l}}$ includes *all* non-incremental sequence prediction algorithms of size $\leq \tilde{l}$ and computation time $\leq \tilde{t}/2$. By non-incremental, we mean that no information of previous cycles is taken into account for speeding up the computation of \dot{y}_k of the current cycle.

²We assume that a TM can be simulated by another in linear time.

The shortcomings (mentioned and unmentioned ones) of this approach are cured in the next subsection, by deviating from the standard way of defining a time bounded ξ as a sum over functions or programs.

7.2.3 The Idea of the Best Vote Algorithm

A general agent is a chronological program $p(x_{<k}) = y_{1:k}$. This form, introduced in Section 4.1, is general enough to include any AI system (and also less intelligent systems). In the following, we are interested in programs p of length $\leq \tilde{l}$ and computation time $\leq \tilde{t}$ per cycle. One important point in the time-limited setting is that p should be incremental, i.e. when computing y_k in cycle k , the information of the previous cycles stored on the work tape can be re-used. Indeed, there is probably no practically interesting, non-incremental AI system at all.

In the following, we construct a policy p^* , or more precisely, policies p_k^* for every cycle k that outperform all time and length limited AI systems p . In cycle k , p_k^* runs all $2^{\tilde{l}}$ programs p and selects the one with the best output y_k . This is a ‘best vote’ type of algorithm, as compared to the ‘weighted majority’ like algorithm of the last subsection. The ideal measure for the quality of the output would be the ξ -expected future reward

$$V_{km}^{p\xi}(\dot{y}\dot{x}_{<k}) := \sum_{q \in \dot{Q}_k} 2^{-l(q)} V_{km}^{pq} \quad , \quad V_{km}^{pq} := r(x_k^{pq}) + \dots + r(x_m^{pq}) \quad (7.5)$$

The program p which maximizes $V_{km_k}^{p\xi}$ should be selected. We have dropped the normalization \mathcal{N} unlike in (5.13), as it is independent of p and does not change the order relation which we are solely interested in here. Furthermore, without normalization, $V_{km}^{*\xi}(\dot{y}\dot{x}_{<k}) := \max_{p \in \dot{P}} V_{km}^{p\xi}(\dot{y}\dot{x}_{<k})$ is enumerable, which will be important later.

7.2.4 Extended Chronological Programs

In the (functional form of the) AI ξ model it was convenient to maximize V_{km_k} over all $p \in \dot{P}_k$, i.e. all p consistent with the current history $\dot{y}\dot{x}_{<k}$. This was no restriction, because for every possibly inconsistent program p there exists a program $p' \in \dot{P}_k$ consistent with the current history and identical to p for all future cycles $\geq k$. For the time limited best vote algorithm p^* it would be too restrictive to demand $p \in \dot{P}_k$. To prove universality, one has to compare *all* $2^{\tilde{l}}$ algorithms in every cycle, not just the consistent ones. An inconsistent algorithm may become the best one in later cycles. For inconsistent programs we have to include the \dot{y}_k into the input, i.e. $p(\dot{y}\dot{x}_{<k}) = y_{1:k}^p$ with $\dot{y}_i \neq y_i^p$ possible. For $p \in \dot{P}_k$ this was not necessary, as p knows the output $\dot{y}_k \equiv y_k^p$ in this case. The r_i^{pq} in the definition of V_{km} are the valuations emerging in the I/O sequence, starting with $\dot{y}\dot{x}_{<k}$ (emerging from p^*) and then continued by applying p and q with $\dot{y}_i := y_i^p$ for $i \geq k$.

Another problem is that we need V_{km_k} to select the best policy, but unfortunately V_{km_k} is uncomputable. Indeed, the structure of the definition of V_{km_k} is very similar to that of \dot{y}_k , hence a brute force approach to approximate V_{km_k} requires too much computation time as for \dot{y}_k . We solve this problem in a similar way, by supplementing each p with a program that estimates V_{km_k} by w_k^p within time \tilde{t} . We combine the calculation of y_k^p and w_k^p and extend the notion of a chronological program once again to

$$p(\dot{y}_{<k}) = w_1^p y_1^p \dots w_k^p y_k^p \quad (7.6)$$

with chronological order $w_1^p y_1^p \dot{y}_1 x_1 w_2^p y_2^p \dot{y}_2 x_2 \dots$

7.2.5 Valid Approximations

p might suggest any output y_k^p but it is not allowed to rate it with an arbitrarily high w_k^p if we want w_k^p to be a reliable criterion for selecting the best p . We demand that no policy is allowed to claim that it is better than it actually is. We define a (logical) predicate $\text{VA}(p)$ called *valid approximation*, which is true if, and only if, p always satisfies $w_k^p \leq V_{km_k}^{p\xi}$, i.e. never overrates itself.

$$\text{VA}(p) \equiv \forall k \forall w_1^p y_1^p \dot{y}_1 x_1 \dots w_k^p y_k^p : p(\dot{y}_{<k}) = w_1^p y_1^p \dots w_k^p y_k^p \Rightarrow w_k^p \leq V_{km_k}^{p\xi}(\dot{y}_{<k}) \quad (7.7)$$

In the following, we restrict our attention to programs p , for which $\text{VA}(p)$ can be proven in some formal axiomatic system. A very important point is that $V_{km_k}^{*\xi}$ is enumerable. This ensures the existence of sequences of programs p_1, p_2, p_3, \dots for which $\text{VA}(p_i)$ can be proven and $\lim_{i \rightarrow \infty} w_k^{p_i} = V_{km_k}^{*\xi}$ for all k and all I/O sequences. p_i may be defined as the naive (non-halting) approximation scheme (by enumeration) of $V_{km_k}^{*\xi}$ terminated after i time steps and using the approximation obtained so far for $w_k^{p_i}$ together with the corresponding output $y_k^{p_i}$. The convergence $w_k^{p_i} \xrightarrow{i \rightarrow \infty} V_{km_k}^{*\xi}$ ensures that $V_{km_k}^{*\xi}$, which we claimed to be the universally optimal value, can be approximated by p with provable $\text{VA}(p)$ arbitrarily well, when given enough time. The approximation is not uniform in k , but this does not matter as the selected p is allowed to change from cycle to cycle.

Another possibility would be to consider only those p which check $w_k^p \leq V_{km_k}^{p\xi}$ online in every cycle, instead of the pre-check $\text{VA}(p)$, either by constructing a proof (on the work tape) for this special case, or $w_k^p \leq V_{km_k}^{p\xi}$ is already evident by the construction of w_k^p . In cases where p cannot guarantee $w_k^p \leq V_{km_k}^{p\xi}$ it sets $w_k = 0$ and, hence, trivially satisfies $w_k^p \leq V_{km_k}^{p\xi}$. On the other hand, for these p it is also no problem to prove $\text{VA}(p)$ as one has simply to analyze the internal structure of p and recognize that p shows the validity internally itself, cycle by cycle, which is easy by assumption on p . The cycle by cycle check is, therefore, a special case of the pre-proof of $\text{VA}(p)$.

7.2.6 Effective Intelligence Order Relation

In Section 5.1 we have introduced an intelligence order relation \succeq on AI systems, based on the expected reward $V_{km_k}^{p\xi}$. In the following we need an order relation \succeq^c based on the claimed reward w_k^p which might be interpreted as an approximation to \succeq .

Definition 7.8 (Effective intelligence order relation) We call p *effectively more or equally intelligent* than p' if

$$p \succeq^c p' :\Leftrightarrow \forall k \forall \dot{y}_{<k} \exists w_{1:n} w'_{1:n} : \\ p(\dot{y}_{<k}) = w_1 * \dots w_k * \wedge p'(\dot{y}_{<k}) = w'_1 * \dots w'_k * \wedge w_k \geq w'_k$$

i.e. if p always claims higher reward estimate w than p' .

\succeq^c is a co-enumerable partial order relation on extended chronological programs. Restricted to valid approximations it orders the policies w.r.t. the quality of their outputs *and* their ability to justify their outputs with high w_k .

7.2.7 The Universal Time Bounded AIXItl Agent

In the following, we describe the algorithm p^* underlying the universal time bounded AIXitl agent. It is essentially based on the selection of the best algorithms p_k^* out of the time \tilde{t} and length \tilde{l} bounded p , for which there exists a proof of $\text{VA}(p)$ with length $\leq l_P$.

1. Create all binary strings of length l_P and interpret each as a coding of a mathematical proof in the same formal logic agent in which $\text{VA}(\cdot)$ has been formulated. Take those strings which are proofs of $\text{VA}(p)$ for some p and keep the corresponding programs p .
2. Eliminate all p of length $> \tilde{l}$.
3. Modify all p in the following way: all output $w_k^p y_k^p$ of p is temporarily written on an auxiliary tape. If p stops in \tilde{t} steps the internal ‘output’ is copied to the output tape. If p does not stop after \tilde{t} steps a stop is forced and $w_k = 0$ and some arbitrary y_k is written on the output tape. Let P be the set of all those modified programs.
4. Start first cycle: $k := 1$.
5. Run every $p \in P$ on extended input $\dot{y}_{<k}$, where all outputs are redirected to some auxiliary tape: $p(\dot{y}_{<k}) = w_1^p y_1^p \dots w_k^p y_k^p$. This step is performed incrementally by adding y_{k-1} for $k > 1$ to the input tape and continuing the computation of the previous cycle.
6. Select the program p with highest claimed reward w_k^p : $p_k^* := \arg\max_p w_k^p$.
7. Write $\dot{y}_k := y_k^{p_k^*}$ to the output tape.
8. Receive input \dot{x}_k from the environment.

9. Begin next cycle: $k := k + 1$, goto step 5.

It is easy to see that the following theorem holds.

Theorem 7.9 (Optimality of AIXI $_{tl}$) Let p be any extended chronological (incremental) program like (7.6) of length $l(p) \leq \tilde{l}$ and computation time per cycle $t(p) \leq \tilde{t}$, for which there exists a proof of $\text{VA}(p)$ defined in (7.7) of length $\leq l_P$. The algorithm p^* constructed in the last subsection, depending on \tilde{l} , \tilde{t} and l_P but not on p , is effectively more or equally intelligent, according to \succeq^c (see Definition 7.8) than any such p . The size of p^* is $l(p^*) = O(\log(\tilde{l} \cdot \tilde{t} \cdot l_P))$, the setup-time is $t_{\text{setup}}(p^*) = O(l_P \cdot 2^{l_P})$ and the computation time per cycle is $t_{\text{cycle}}(p^*) = O(2^{\tilde{l}} \cdot \tilde{t})$.

Roughly speaking, the theorem says, that if there exists a computable solution to some (or all) AI problem(s) at all, the explicitly constructed algorithm p^* is such a solution. Although this theorem is quite general, there are some limitations and open questions which we discuss in the following.

The construction of the algorithm p^* needs the specification of a formal logic system $(\forall, \lambda, y_i, c_i, f_i, R_i, \rightarrow, \wedge, =, \dots)$, and axioms, and inference rules. A proof is a sequence of formulas, where each formula is either an axiom or inferred from previous formulas in the sequence by applying the inference rules. Details have been presented in Section 7.1.5. We only need to know that *provability* and *Turing Machines* can be formalized. The setup time in the main theorem is just the time needed to verify the 2^{l_P} proofs, each needing time $O(l_P^2)$.

7.2.8 Limitations and Open Questions

- Formally, the total computation time of p^* for cycles $1 \dots k$ increases linearly with k , i.e. is of order $O(k)$ with a coefficient $2^{\tilde{l}} \cdot \tilde{t}$. The unreasonably large factor $2^{\tilde{l}}$ is a well known drawback in best/democratic vote models and will be taken without further comments, whereas the factor \tilde{t} can be assumed to be of reasonable size. If we don't take the limit $k \rightarrow \infty$ but consider reasonable k , the practical usefulness of the time bound on p^* is somewhat limited, due to the additional additive constant $O(l_P \cdot 2^{l_P})$. It is much larger than $k \cdot 2^{\tilde{l}} \cdot \tilde{t}$ as typically $l_P \gg l(\text{VA}(p)) \geq l(p) \equiv \tilde{l}$.
- p^* is superior only to those p which justify their outputs (by large w_k^p). It might be possible that there are p which produce good outputs y_k^p within reasonable time, but it takes an unreasonably long time to justify their outputs by sufficiently high w_k^p . We do not think that (from a certain complexity level onwards) there are policies where the process of constructing a good output is completely separated from some sort of justification process. But this justification might not be translatable (at least within reasonable time) into a reasonable estimate of $V_{km_k}^{p\xi}$.

- The (inconsistent) programs p must be able to continue strategies started by other policies. It might happen that a policy p steers the environment to a direction for which p is specialized. A ‘foreign’ policy might be able to displace p only between loosely connected episodes. There is probably no problem for factorizable μ . Think of a chess game, where it is usually very difficult to continue the game/strategy of a different player. When the game is over, it is usually advantageous to replace a player by a better one for the next game. There might also be no problem for sufficiently separable μ .
- There might be (efficient) valid approximations p for which $\text{VA}(p)$ is true but not provable, or for which only a very long ($> l_P$) proof exists.

7.2.9 Remarks

- The idea of suggesting outputs and justifying them by proving reward bounds implements one aspect of human thinking. There are several possible reactions to an input. Each reaction possibly has far reaching consequences. Within a limited time one tries to estimate the consequences as well as possible. Finally, each reaction is valued and the best one is selected. What is inferior to human thinking is, that the estimates w_k^p must be rigorously proved and the proofs are constructed by blind exhaustive search, further, that *all* behaviors p of length $\leq \tilde{l}$ are checked. It is inferior ‘only’ in the sense of necessary computation time but not in the sense of the quality of the outputs.
- In practical applications there are often cases with short and slow programs p_s performing some task T , e.g. the computation of the digits of π , for which there exist long but quick programs p_l too. If it is not too difficult to prove that this long program is equivalent to the short one, then it is possible to prove $K^{t(p_l)}(T) \stackrel{+}{\leq} l(p_s)$ with K^t being the time bounded Kolmogorov complexity. Similarly, the method of proving bounds w_k for V_{km_k} can give high lower bounds without explicitly executing these short and slow programs, which mainly contribute to V_{km_k} .
- Dovetailing all length and time-limited programs is a well known elementary idea (typing monkeys). The crucial part which has been developed here, is the selection criterion for the most intelligent agent.
- By construction of $\text{AI}\xi^{\tilde{t}\tilde{l}}$ and due to the enumerability of V_{km_k} , ensuring arbitrary close approximations of V_{km_k} we expect that the behavior of $\text{AI}\xi^{\tilde{t}\tilde{l}}$ converges to the behavior of $\text{AI}\xi$ in the limit $\tilde{t}, \tilde{l}, l_P \rightarrow \infty$ in a sense.
- Depending on what you know/assume that a program p of size \tilde{l} and computation time per cycle \tilde{t} is able to achieve, the computable $\text{AI}\xi^{\tilde{t}\tilde{l}}$ model will have the same capabilities. For the strongest assumption of the existence of a Turing machine, which outperforms human intelligence, the $\text{AI}\xi^{\tilde{t}\tilde{l}}$ will do too, within the same time frame up to a (unfortunately very large) constant factor.



Ray Solomonoff

“... in spite of its incomputability, Algorithmic Probability can serve as a kind of ‘Gold Standard’ for induction systems.” (Solomonoff, 1997)

“It’s hard to make predictions, especially about the future” (Niels Bohr)

“... we have the mathematical theory of decision making under uncertainty. What the mathematical theory is worth, it is hard to say. It does have the advantage, though, of providing definite rules.” (Richard E. Bellman)

Chapter 8

Discussion

8.1	What has been Achieved	802
8.1.1	Results	802
8.1.2	Comparison to other Approaches	803
8.2	General Remarks	804
8.2.1	Miscellaneous	805
8.2.2	Prior Knowledge	806
8.2.3	Universal Prior Knowledge	806
8.2.4	How AIXI(<i>tl</i>) Deals with Encrypted Information	807
8.2.5	Mortal Embodied Agents	807
8.3	Personal Remarks	808
8.3.1	On the Foundations of Machine Learning	809
8.3.2	In a World without Occam	810
8.4	Outlook & Open Questions	810
8.5	Assumptions, Problems, Limitations	812
8.5.1	Assumptions	812
8.5.2	Problems	813
8.5.3	Limitations	814
8.6	Philosophical Issues	814
8.6.1	Turing Test	814
8.6.2	On the Existence of Objective Probabilities	815
8.6.3	Free will versus Determinism	815
8.6.4	The Big Questions	817
8.7	Conclusions	818

This chapter critically reviews what has been achieved in the thesis and discusses

some otherwise unmentioned topics of general interest. We summarize our major results and compare performance and generality of $\text{AIXI}(tl)$ to those of other approaches to AI. We remark on various topics, including concurrent actions and observations, the choice of the I/O spaces, treatment of encrypted information, and peculiarities of mortal embodied agents. We also make some personal comments and speculations on the present status and the future of the research fields AI and machine learning themselves. We continue with an outlook on further research. Since many ideas have already been presented in the Problems and Conclusions sections of the various chapters, we concentrate on non-technical open questions of general importance, including optimality, down-scaling, implementation, approximation, elegance, extra knowledge, and training of/for $\text{AIXI}(tl)$. Furthermore, we collect and state all explicit or implicit assumptions, problems and limitations of $\text{AIXI}(tl)$. We briefly discuss some relevant philosophical issues: The free will versus determinism paradox, the existence of objective probabilities, and the Turing test. We also include some (personal) remarks on non-computable physics, the number of wisdom Ω , and consciousness. As it should be, the thesis concludes with conclusions.

8.1 What has been Achieved

8.1.1 Results

The major theme of the thesis was to develop a mathematical foundation of Artificial Intelligence. This is not an easy task since intelligence has many (often ill-defined) faces. More specifically, our goal was to develop a theory for rational agents acting optimally in any environment. Thereby we touched various scientific areas, including reinforcement learning, algorithmic information theory, Kolmogorov complexity, computational complexity theory, information theory and statistics, Solomonoff induction, Levin search, sequential decision theory, adaptive control theory, and many more. The major achievements have been the following:

- We derived various convergence results, (tight) loss bounds, and Pareto-optimality for predictors based on Bayes-mixture priors. We gave an Occam's razor argument that, using Solomonoff's prior leads to a universally optimal prediction scheme (Chapter 3).
- We presented sequential decision theory in a very general form in which actions and observations may depend on arbitrary past events. The development was more of a formal exercise and optimality of the $\text{AI}\mu$ model for known environment μ is obvious by construction (Chapter 4).
- We unified sequential decision theory and Solomonoff's theory of universal induction (both optimal in their own domain). The resulting parameter-free AIXI model constitutes an agent for which we gave strong arguments that it behaves optimally in any environment (it copes with exploration versus

exploitation, large state spaces, generalization and function approximation, non-stationary and partially observable environments, ...) (Chapter 5).

- Vice versa we defined a universal intelligence order relation \prec regarding which AIXI is the most intelligent agent and argued this order relation to be reasonable (Section 5.1.4).
- We discussed the difficulties in extending the optimality results from the prediction case to AIXI. In this course we suggested various potentially relevant environmental (separability) concepts (Sections 5.2 and 5.3).
- We discussed the choice of the horizon and came to the conclusion that a reward discounting (like near-harmonic) which leads to an effective horizon that increases in proportion to the current age of the agent is best (Section 5.7).
- For restricted environmental classes and Bayes-mixtures ξ we showed that AIXI_ξ is self-optimizing and Pareto-optimal (Section 5.4 and 5.6).
- We showed how AIXI is suitable for dealing with a number of important problem classes, including sequence prediction, strategic games, function minimization, and supervised learning (Chapter 6).
- Based on the mathematical (incomputable) AIXI model we developed a computable model, AIXI_{tl} , with optimal order of computation time, apart from a large multiplicative constant (Section 7.2).
- We developed a general purpose algorithm – the asymptotically fastest (and shortest) algorithm for all well-defined problems. We got rid of the large multiplicative constant as in Levin search and AIXI_{tl} , at the expense of an (unfortunately even larger) additive constant (Section 7.1).

All in all, the results show that Artificial Intelligence can be framed by an elegant mathematical theory. Some progress has also been made towards an elegant *computational* theory of intelligence.

8.1.2 Comparison to other Approaches

A different way to measure the achievements of this work is to compare $\text{AIXI}(tl)$ to other AI approaches. In Table 8.1 we compare various learning algorithms which are rather general in purpose, have an agent-like setup, are popular or otherwise interesting or promising. We subjectively rate the different approaches w.r.t. various performance and generality criteria. We use a grayscale between YES and NO, since the evaluation is often arguable, especially if there are many algorithm variants. For most table cells one can imagine a variant and an application for which rating YES would be justified, and one for which rating NO would be justified. Hence, the

presented ratings refer to typical algorithm variants and typical (intended) applications. It is beyond the scope of this work to describe all these approaches and to justify each rating in detail.

We consider the following *properties* in the different columns: An algorithm is *time efficient* if it runs on a today's computer in acceptable time for "interesting" applications. An algorithm is *data efficient* if it exploits (learns from) the information contained in the received data in a theoretically near optimal fashion. An algorithm gets a *yes* in the *exploration* column only if it addresses the exploration versus exploitation problem in a fundamental near optimal way (Many algorithms are greedy with some simple random exploration added). *Convergence* of an algorithm may be just to *some* policy, to a local optimum, or to a/the *global optimum*. An important issue is whether learning algorithms are capable of *generalizing* from previous experience to *similar* situations. We also indicate whether an algorithm is capable of or designed for dealing with non-Markovian environments, e.g. POMDPs. All selected algorithms are capable of *learning* by experience, except Value/Policy iteration which need an exact description of the environment in advance. The last column distinguishes between passive predictors and *active* agents.

The first group of algorithms in Table 8.1 contains "classical" reinforcement learning algorithms. See [SB98, BT96] for a description of Value and Policy iteration, Temporal Difference (TD) learning with finite \mathcal{S} state space versus linear/general function approximation. See [BB01, KHS01a] for an introduction to direct gradient-based reinforcement learning. The second group in Table 8.1 contains various other learning algorithms: Logic planners [RN95, Part IV], split trees [Rin94, McC95], adaptive Levin search [SZW97], optimal ordered problem solver (OOPS) [Sch02b], prediction with expert advice (PEA) [CB97], Market/economy based reinforcement learning [Bau99, KHS01b]. The third group in Table 8.1 lists the main models of this work: Sequence prediction based on Solomonoff's prior (SPXI) and the AIXI(tl) model(s). The last line lists the capabilities of human agents. Overall it can be said that the models in the first two groups are applicable in limited domains with feasible computation time, whereas the models of the last group are completely general, but computationally not feasible without further approximations.

8.2 General Remarks

This section remarks on some otherwise unmentioned topics of general interest. The logically disconnected subsections discuss concurrent actions and observations, the choice of the I/O spaces, (universal) prior knowledge, treatment of encrypted information, and peculiarities of mortal embodied agents.

Table 8.1 (Properties of learning algorithms) The table compares various important properties of learning algorithms limited to their typical domains. The evaluation is often subjective and arguable, especially because there are many variants, so we introduce the grayscale YES \rightarrow yes \rightarrow yes/no \rightarrow no/yes \rightarrow no \rightarrow NO (see Section 8.1.2 for further explanation).

Algorithm	time efficient	data efficient	explo-ration	conver-gence	global optimum	genera-lization	POMDP	learning	active
Value/Policy iteration	yes/no	yes	–	YES	YES	NO	NO	NO	yes
TD w/ finite \mathcal{S}	yes/no	NO	NO	YES	YES	NO	NO	YES	YES
TD linear func.approx.	yes/no	NO	NO	yes	yes/no	YES	NO	YES	YES
TD general func.approx.	no/yes	NO	NO	no/yes	NO	YES	NO	YES	YES
Direct Policy Search	no/yes	YES	NO	no/yes	NO	YES	no	YES	YES
Logic Planners	yes/no	YES	yes	YES	YES	no	no	YES	yes
RL with Split Trees	yes	YES	no	YES	NO	yes	YES	YES	YES
Pred.w. Expert Advice	yes/no	YES	–	YES	yes/no	yes	NO	YES	NO
Adaptive LS	no/yes	no	no	yes	yes/no	yes	YES	YES	YES
OOPS	yes/no	no	–	yes	yes/no	YES	YES	YES	YES
Market/Economy RL	yes/no	no	NO	no	no/yes	yes	yes/no	YES	YES
SPXI	no	YES	–	YES	YES	YES	NO	YES	NO
AIXI	NO	YES	YES	yes	YES	YES	YES	YES	YES
AIXI $_{tl}$	no/yes	YES	YES	YES	yes	YES	YES	YES	YES
Human	yes	yes	yes	no/yes	NO	YES	YES	YES	YES

8.2.1 Miscellaneous

Game theory. In game theory [OR94] one often wants to model the situation of simultaneous actions, whereas the AI ξ models have serial I/O. Simultaneity can be simulated by withholding the environment from the current agent’s output y_k , until x_k has been received by the agent. Formally, this means that $\mu(y_{<k}y_k)$ is independent of the last output y_k . The AI ξ agent is already of simultaneous type in an abstract view if the behavior p is interpreted as the action. In this sense, AIXI is the action p^* which maximizes the utility function (reward), under the assumption that the environment acts according to ξ . The situation is different from game theory as the environment ξ is not a second ‘player’ that tries to optimize his own utility (see Section 6.3).

Input/output spaces. In various examples we have chosen differently specialized input and output spaces \mathcal{X} and \mathcal{Y} . It should be clear that, in principle, this is unnecessary, as large enough spaces \mathcal{X} and \mathcal{Y} (e.g. the set of strings of length 2^{32}) serve every need and can always be Turing-reduced to the specific presentation

needed internally by the AIXI agent itself. But it is clear that, using a generic interface, such as camera and monitor for learning tic-tac-toe, for example, adds the task of learning vision and drawing.

8.2.2 Prior Knowledge

In many problems in practice we have extra information about the problem at hand, which could and should be used to guide the forecasting. If the prior knowledge is of the form that it includes only environments of certain structures, e.g. MDPs one can use the appropriate Bayes-mixture over these environments. If there is reason to belief that certain environments are less or more likely than Occam's razor tells us, then this could be coded in the weights w_ν . Unfortunately, this procedure is often intractable in practice, since one has only a (possibly vague) description of prior facts, which are hard to translate into classes \mathcal{M} and/or weights w_ν . Fortunately there is a simple way of incorporating all prior knowledge D in an easy and optimal way. The trick is to get rid of all prior knowledge by prefixing the observation sequence $x_1x_2\dots$ by *some* binary coding $d_{1:l}$ of D . Using then Solomonoff's prior ξ_U on $d_{1:l}x_{1:n}$ for prediction on cycles $l+1$ to $n+l$ one gets loss bounds (to logarithmic accuracy) in terms of $K(\mu|D)$. If D contains information about μ it will reduce the Kolmogorov complexity of μ , if not we cannot expect D to improve prediction accuracy. This also solves the often mentioned concern of how to make good predictions for short sequences (sparse data) of length $n=O(1)$. As long as $n+l$ is larger than the typical compiler constants, predictions based on ξ_U are good. It *seems* that in science one often faces problems with data of information content, say 200 bits only, and none or very little prior knowledge, say only 100 bits is available. For instance having thrown a biased coin for 200 times, and describing our prior knowledge as "i.i.d. with uniform second order prior over bias θ " *seems* not to contribute more than 100 bits on prior information. Laplace's law of succession [Lap1814] leads to reasonable estimates of θ and predictions of further tosses, whereas ξ_U is far from μ in cycle 300 for typical U . But this is an illusion that we only have 100 bits of prior knowledge. We spent at least 15 years in school, before having heard about Dirichlet priors and Laplace rule. Our whole scientific knowledge serves as prior knowledge. If we take for D a representative collection of scientific books (+ some language books), l is much larger than the typical compiler constants and $\xi_U(x_n|d_{1:l}x_{<n})$ will be very close to the true bias θ ! This holds true also for more complex examples. We can make non-arbitrary predictions given a sequence of $l+n-1$ bits only if ξ_U leads to the same prediction for all "reasonably complex" universal Turing-machines U .

8.2.3 Universal Prior Knowledge

There are people who believe universal AI is not possible, that one *has* to incorporate some/sufficient prior knowledge. I disagree in a sense described in Section 8.4. A

different approach is to exclude only those environments which we are sure not to be realized. This approach is worth considering, but has the following problems:

1) Physical knowledge is never 100% sure. For instance, 100 years ago everybody would have assumed a flat 3D universe. I'm not too concerned about this, since today's physical theories (at least the parts which seem to be *relevant* for (in a very broad sense) human-sized and equipped agents) are very accurate, and reliable. Instead of eliminating universes which seem to be excluded by our observations and theories, one may only reduce the prior belief in odd universes, but this does not help in substantially increasing the prior belief of the likely universes.

2) More seriously, μ does not describe the total universe, but only a small fraction, from the subjective perspective of the agent. It is (somewhat/much?) harder to characterize the set of possible universes \mathcal{M} from the subjective agent perspective.

3) One may take into account only general properties of the universe like locality, continuity, or the existence of (manipulable) objects with properties and relations in a manifold. The major problem is that, although the universe seems to be a local continuous MDP (ignoring quantum effects), μ is neither an MDP, nor local. What the agent directly observes (with his sensors, like a camera) is not the complete MDP state and often appears non-local. So probably very little *really* exploitable can be said about μ .

Of course, the scientific approach is to simply *assume* some properties (whether true in real life or not) and analyze the performance of the resulting models.

8.2.4 How AIXI(tl) Deals with Encrypted Information

If you encrypt a message with a key of size k but then you also give away the key, then the overall algorithm A for encryption has size $O(1)$. A = "Take key and message and encrypt message with key". The system only has to "learn" A^{-1} which can itself be described in length $O(1)$. Only very little information is needed to learn $O(1)$ bits. In this sense decryption is easy. The problem is that A^{-1} may be an extremely slow algorithm (e.g. finding the prime factors from the public key in RSA). But note, in AIXI we are not talking about computation time, we are only talking about information efficiency (learning in the least number of interaction cycles). This is maybe one of the key ideas to separate data efficiency from computation time efficiency. Of course in the real world computation time matters, so we invented AIXI tl . AIXI tl can do every job as good as the best length l and time t bounded agent, apart from time factor 2^l and a huge offset time. This offset time would be used, for instance, in the RSA example, to (once-and-for-all) find the factorization, and then, decryption is easy, of course.

8.2.5 Mortal Embodied Agents

The examples we gave in this thesis, particularly those in Chapter 6, were mainly bodiless agents: predictors, gamblers, optimizers, learners. There are some pecu-

liarities with reinforcement learning, autonomous, embodied robots in real environments.

We can still reward the robot according to how well it solves the task we want it to do. A minimal requirement is that the robot's hardware functions properly. If the robot starts to malfunction its capabilities degrade, resulting in lower reward. So in an attempt to maximize reward the robot will also maintain itself. The problem is that some parts will malfunction rather quickly when no appropriate actions are performed, e.g. flat batteries, if not recharged in time. Even worse, the robot may work perfectly until the battery is nearly empty, and then suddenly stop its operation (death), resulting in (minimal) zero reward from then on. There is too little time to learn how to maintain itself before it's too late. An autonomous embodied robot cannot start from scratch but must have some rudimentary built-in capabilities (which may not be that rudimentary at all) which allows it to at least survive. This is similar to the problem discussed in Section 6.4.5 of using AIXI in the FMF setting with too late reward. Using FMF ξ corresponds to incorporating some rudimentary capability. Animals survive due to reflexes, innate behavior, an internal reward attached to the condition of their organs, and a guarding environment during childhood. Different species emphasize different aspects. Reflexes and innate behaviors are stressed in lower animals versus years of safe childhood for humans. The same variety of solutions is available for constructing autonomous robots (which we will not detail here).

Another problem connected, but possibly not limited to embodied agents, especially if they are rewarded by humans, is the following: Sufficiently intelligent agents may increase their rewards by psychologically manipulating their human "teachers", or by threatening them. This is a general sociological problem which successful AI will cause, which has nothing specifically to do with AIXI. Every intelligence superior to humans is capable of manipulating the latter. In the absence of manipulable humans, e.g. where the reward structure serves a survival function, AIXI may directly hack into its reward feedback. Since this is unlikely to increase its long-term survival, AIXI will probably resist this kind of manipulation (like most humans don't take hard drugs, due to their long-term catastrophic consequences).

8.3 Personal Remarks

It is hard to predict the future, as it is to predict the development of research areas like AI. Nevertheless I would like to risk a try. To be more specific, in the following I suggest a framework for machine learning research. It is a mixture of how I expect and would like the field to look in the near future.

8.3.1 On the Foundations of Machine Learning

Instead of addressing machine learning directly, let us first consider a different research area such as algorithm and complexity theory. The goal of algorithm theory is to find and analyze fast algorithms, the goal of complexity theory is to show lower bounds on the time needed to solve certain problem classes. All concepts are rigorously defined: algorithm, Turing machine, problem class, computation time, ... Most disciplines generally start with an informal way of attacking a subject. With time the discipline becomes more and more formalized, often up to a point where it is completely rigorous. Examples are number theory, set theory, proof theory, probability theory, infinitesimal calculus, quantum field theory, ... Each theory experienced a time in which it was dealt with in an informal way, but after a while it was made rigorous, is now completely axiomatized, and rarely questioned¹. Of course not all disciplines are axiomatized yet or are axiomatizable at all (e.g. biology), and new research areas emerge, starting in an informal condition, but the point is that once a field has emerged, the path is towards increasing rigor.

What can be said about machine learning? In machine learning one tries to build and understand systems which learn from past data and make good prediction, which are able to generalize, act intelligently, ... Many terms are only vaguely defined or have many alternative definitions. As discussed in Chapter 2 and elsewhere, from a formal point of view, all learning tasks can be unified in the framework of sequence prediction. We propose Occam's razor, quantified in terms of Solomonoff-Kolmogorov complexity, combined with the chain rule for conditional probabilities, and possibly sequential decision theory as a rigorous mathematical/axiomatic definition of machine learning. More precisely, Solomonoff's induction scheme should be "used" for sequence prediction tasks, and, when combined with sequential decision theory for making sequential decisions. The results of this work and also the results of the more applied MML, MDL, and SRM principles support the power of Occam's razor. As long as there is no convincing evidence against Occam's razor, and even more importantly, as long as there is no alternate suggestion of how to define machine learning rigorously, it is worth assuming Occam's razor and studying its consequences. Indeed, we have shown in Chapter 3 that the performance of Solomonoff's universal induction scheme, as compared to any other prediction scheme in any environment, is so good that one may be tempted to take the results as proof of Occam's razor. Whereas the theorems were proven with mathematical rigor, one has to be careful about their interpretation and the underlying assumptions, especially in Theorem 3.70 which was more a self-consistency or bootstrap result.

We expect that in the future, machine learning will, by default, be based on Occam's razor. Real-world machine learning tasks will with overwhelming majority be solved by developing algorithms which approximate Kolmogorov complexity /

¹Quantum field theory may be argued not to be in a completely mathematically satisfactory condition, yet.

Solomonoff’s prior (e.g. MML, MDL, SRM, and more specific ones, like SVM, ZIP, Neural/Bayes nets with complexity penalty, ...). Machine learning theory will derive results on convergence speed and approximation quality of the various approximation schemes. Only a minority will investigate non-standard ML by modifying or replacing Occam’s razor “axiom” in the hope of finding something better.

8.3.2 In a World without Occam

Finally I would like to remark on an analogy to Peano’s axioms for the natural numbers and especially the induction axiom. Remove this axiom and replace it with a vague concept which resembles this axiom². It is then not unreasonable to call this concept *induction principle* since it infers properties valid for *all* natural numbers from a local $n \rightarrow n+1$ property. Imagine arithmetic were still in this situation. Most modern mathematical theorems would evaporate and, with them nearly all modern technology. Fortunately we have this induction axiom, but as a formal rule it is now purely deductive!

I believe the same is/will be true for Occam’s razor. Without the vague concept of Occam’s razor, science and, hence, machine learning would probably be not existent at all. Informal Occam’s razor is directly or indirectly the basis of all scientific induction. The establishment of a formal version of Occam’s razor would give machine learning in particular, and maybe even science in general, a significant boost. I anticipate this by looking at the practical success of MML, MDL, SRM, and SVM and the theoretical impact of Kolmogorov complexity and Solomonoff induction, which are all formalizations of Occam’s razor.

8.4 Outlook & Open Questions

Many ideas for further studies have already been stated in the various chapters of the thesis, especially in the Problems and Conclusions sections. This outlook only contains non-technical open questions regarding $\text{AIXI}(tl)$ of general importance.

Value bounds. Rigorous proofs for non-asymptotic value bounds for $\text{AI}\xi$ are the major theoretical challenge – general ones, as well as tighter bounds for special environments μ , e.g. for rapidly mixing MDPs. For AIXI other performance criteria have to be found and proved. Although not necessary from a practical point of view, the study of continuous classes \mathcal{M} , restricted policy classes, and/or infinite \mathcal{Y} , \mathcal{X} and m may lead to useful insights.

Scaling AIXI down. A direct implementation of the $\text{AIXI}(tl)$ model is, at best, possible for toy environments due to the large factor 2^l in computation time. But there are other applications of the AIXI theory. We have seen in several examples how to integrate problem classes into the AIXI model. Conversely, one can downscale

²I guess that there was a time in history when arithmetic was exactly in this condition.

the $AI\xi$ model by using more restricted forms of ξ . This could be done in the same way as the theory of universal induction has been downscaled with many insights to the Minimum Description Length principle [LV92a, Ris89] or to the domain of finite automata [FMG92]. The AIXI model might similarly serve as a super model or as the very definition of (universal unbiased) intelligence, from which specialized models could be derived.

Implementation and approximation. With a reasonable computation time, the AIXI model would be a solution of AI (see next point if you disagree). The $AIXItl$ model was the first step, but the elimination of the factor 2^l without introducing a large additive constant like in $M_{p^*}^\varepsilon$ and without giving up universality will (almost certainly) be a very difficult task. One could try to select programs p and prove $VA(p)$ in a more clever way than by mere enumeration, to improve performance without destroying universality. All kinds of ideas like genetic algorithms, advanced theorem provers and many more could be incorporated. It remains to be seen whether these hand waving suggestions can be substantiated.

Computability. We seem to have transferred the AI problem just to a different level, to proving $VA(p)$. This shift has some advantages (and also some disadvantages) but presents, in no way, a solution. Nevertheless, we want to stress that we have reduced the AI problem to (mere) computational questions. Even the most general other systems the author is aware of, depend on some (more than complexity) assumptions about the environment or it is far from clear whether they are, indeed, universally optimal. Although computational questions are themselves highly complicated, this reduction is a non-trivial result. A formal theory of something, even if not computable, is often a great step toward solving a problem and also has merits of its own, and AI should not be different in this respect (see previous item).

Elegance. Many researchers in AI believe that intelligence is something complicated and cannot be condensed into a few formulas. It is more a combining of enough *methods* and much explicit *knowledge* in the right way. From a theoretical point of view we disagree, as the AIXI model is simple and seems to serve all needs. From a practical point of view we agree to the following extent: To reduce the computational burden one should provide special purpose algorithms (*methods*) from the very beginning, probably many of them related to reduce the complexity of the input and output spaces \mathcal{X} and \mathcal{Y} by appropriate pre/post-processing *methods*.

Extra knowledge. There is no need to incorporate extra *knowledge* from the very beginning. It can be presented in the first few cycles in *any* format. As long as the algorithm to interpret the data is of size $O(1)$, the AIXI agent will ‘understand’ the data after a few cycles (see Section 8.2.2 and 6.5). If the environment μ is complicated but extra knowledge z makes $K(\mu|z)$ small, one can show that the bound (5.9,5.10) reduces roughly to $\ln 2 \cdot K(\mu|z)$ when $x_1 \equiv z$, i.e. when z is presented in the first cycle. The special purpose algorithms could be presented in x_1 , too, but it would be cheating to say that no special purpose algorithms had been implemented

in AIXI. The boundary between implementation and training is unsharp in the AIXI model.

Training. We have not said much about the training process itself, as it is not specific to the AIXI model and has been discussed in literature in various forms and disciplines [Sol86, Sch02b]. By a training process we mean a sequence simple-to-complex tasks to solve, with the simpler ones hopefully helping in learning the more complex ones. A serious discussion would be out of place. To repeat a truism, it is, of course, important to present enough knowledge x'_k and evaluate the agent output y_k with r_k in a reasonable way. To maximize the information content in the reward, one should start with simple tasks and give positive reward to approximately the better half of the outputs y_k .

8.5 Assumptions, Problems, Limitations

Just as every approach to AI (or any other field) makes assumptions and has its problems and limitations, so does AIXI(tl). It is time to take a critical look at all explicit or implicit assumptions, problems and limitations.

8.5.1 Assumptions

- The central assumption of this work is Occam's razor. Since Occam's razor seems to be at the heart of science and intelligent behavior in any case, it is not a restrictive assumption, but nevertheless a profound one. Occam's razor actually only serves as a motivation in this work; the actual assumption we use is different/weaker, see next item.
- The environment is sampled from a computable probability distribution with a reasonable program size on a natural Turing machine. Assumption 2.5 ensures that AIXI(tl) is essentially independent of whatever universal Turing machine is chosen.
- We assumed the existence of objective randomness/probabilities respecting Kolmogorov's probability Axioms 2.14. As remarked in Section 2.3 this assumption is not essential, since we can restrict the setting to deterministic environments μ . Using Bayes-mixtures as subjective probabilities also did not involve any assumptions, since they were justified decision-theoretically.
- All reinforcement learning approaches we are aware of define the total reward as a *sum* of rewards $r_1 + \dots + r_m$ over cycles, and so do we. In finance, where money can be reinvested, a product is common, but this can be converted to a sum by taking the logarithm. In Section 6.4.3 we encountered a case where the minimum/maximum may be appropriate. One may replace the reward

sum in AIXI (5.3) by other functions, but we take a pragmatic view and stick to the reward sum as long as there is no evidence of any serious deficiency.

- For probabilistic environments we defined the value of a policy which shall be maximized in the standard way as the *expected* reward sum. More generally one may define for each policy a probability distribution for the total reward. The question is how to compare these distributions. Besides the popular mean, one may want to compare medians or quantiles. The most frequent argument for departing from the mean is to achieve more robust policies, e.g. policies which, with high probability, have a high lower bound on their reward sum. We believe that robustness is never a primary goal in itself. The reason for wanting robustness is that one dislikes low rewards more than the rewards themselves express. A natural solution is to take the expectation of f -transformed rewards, where f is a monotone increasing concave function (like log). f penalizes small rewards and leads to more robustness.
- We assumed finite action/observation spaces \mathcal{X}/\mathcal{Y} which are sufficient for all practical purposes. From a theoretical point of view infinite spaces may be attractive in certain situations. Countable \mathcal{X} should cause no problems, in case of countable \mathcal{Y} only ε -optimal policies may exist. For continuous ξ one has to somehow generalize the notion of Kolmogorov complexity and Solomonoff prior.
- We assumed bounded non-negative rewards $r_k \in [0, r_{max}]$. Non-negativity is not essential, but boundedness is essential for ensuring existence of values. Again, from a practical point of view this should not be restrictive.
- We assumed finite horizon or near-harmonic discounting to ensure the existence of values. We provided motivations for the choice of the latter, but we are not sure whether it represents a final answer.

After all this one should not forget that most other approaches to AI implicitly or explicitly make the same, and usually even more, assumptions.

8.5.2 Problems

- Assume AIXI is used in a multi-agent setup interacting with other agents. For simplicity we only discuss the case of a single other agent in a competitive setup, i.e. a two-person zero-sum game situation. We can entangle agents A and B by $x'_k(A) = y_k(B)$, $x'_{k+1}(B) = y_k(A)$. The rewards $r_k(A)$ and $r_k(B)$ are provided externally by the rules of the game. The situation where A is AIXI and B is a perfect minimax player has been analyzed in Section 6.3. In multi-agent systems one is mostly interested in a symmetric setup, i.e. B is also an AIXI. Whereas both AIXIs *may* be able to learn the game and improve their strategies (to optimal minimax), this setup violates one of our

basic assumptions. Since AIXI is incomputable $\text{AIXI}(B)$ does not constitute a computable environment for $\text{AIXI}(A)$. More generally, starting with any class of environments \mathcal{M} , then $\mu \triangleq \text{AIXI}_{\mathcal{M}}$ seems not to belong to class \mathcal{M} for most (all?) choices of \mathcal{M} . Various results of the thesis can no longer be applied, since $\mu \notin \mathcal{M}$ when coupling two AIXIs. Many questions arise: Are there interesting environmental classes for which $\text{AIXI}_{\mathcal{M}} \in \mathcal{M}$ or $\text{AIXI}_{\mathcal{M}} \in \mathcal{M}$? Do $\text{AIXI}(A/B)$ converge to optimal minimax players? Do AIXIs perform well in general multi-agent setups?

8.5.3 Limitations

- Although AIXI may be regarded as a formal definition or a mathematical solution of AI, it is not a practical solution due to its incomputability. $\text{AIXI}_{\mathcal{M}}$ is a step in the direction of a computable theory of AI, but is also not practically feasible. Whether AIXI can be scaled down in a systematic way to yield practical AI algorithms or whether it will only serve as a guiding principle in attacking difficult AI problems remains to be seen.

8.6 Philosophical Issues

Many arguments against (strong and weak) AI have been proposed: Lucas' and Penrose's arguments based on Gödel's incompleteness theorem [Luc61, Pen89, Pen94], Searle's Chinese room argument [Sea80], the lookup-table argument [Chu86], Moravec's brain prosthesis experiment [Mor88], the free will argument, and, of course, various religious reasons. All of them have loopholes and can be refuted, but here is not the place to repeat this discussion. We only discuss the free will paradox. There are also objections to the existence of objective probabilities. We present a possibly new one below. We also briefly comment on the Turing test, which also fits under the heading of this section. Finally, we speculate on the *big* questions of AI in general and the AIXI model in particular, related to non-computable physics, the number of wisdom Ω , and consciousness.

8.6.1 Turing Test

The Turing test [Tur50] was designed to decide whether an AI system is intelligent. We should concede a machine true intelligence if it passes the Turing test, but to deny intelligence in case of failure may be too hard. The true problem in using the Turing test (e.g. instead of AIXI) as a *definition* of intelligent systems is another. The test involves a human interrogator and, hence, cannot be formalized mathematically, therefore it does also not allow the development of a computational theory of intelligence.

8.6.2 On the Existence of Objective Probabilities

Throughout the thesis we have assumed the existence of objective probabilities respecting Kolmogorov's probability Axioms 2.14. As remarked in Section 2.3 we could have restricted the development to classes \mathcal{M} of deterministic environments, thus avoiding objective probabilities³. In the following we give an argument which makes the belief in objective probabilities look somewhat “unscientific”. The assumption that an event occurs with some objective probability expresses the opinion that the occurrence of an individual stochastic event has no explanation, i.e. is inherently impossible to predict for sure. One central goal of science is to *explain* things. Often we do not have an explanation (yet) that's acceptable, but to say that “something can principally not be explained” is to stop even *trying* to find an explanation. From a distance, tossing a coin looks objectively random, but looking at it closer the outcome is just subjectively unknown due to most observers' lack of knowledge of initial conditions and external influences on the coin during its throw. When knowing the exact initial conditions and the exact equations of motion, classical physics is predictable (this includes chaotic systems). Physicists claim that quantum mechanics is truly random, and there is indeed quite some evidence to suggest this, but experiments cannot exclude the possibility that quantum events are only pseudo-random [Sch02c]. It seems safer and more honest to say that with our current technology and understanding we can only determine (subjective) outcome probabilities. If a sufficiently large community of people arrive at the same subjective probabilities from their prior knowledge, one may want to call these probabilities objective. For instance for most people (those with no special equipment and education) a fair coin comes up head in 50% of the cases. And for *all* people so far, if they measure the spin of one photon in a para-positronium decay, it is up in 50% of the cases. On one hand we have to abandon objective probabilities because their assumption seems unscientific, but on the other hand their assumption is very convenient. Without objective probabilities there would be no (objective) unbiased coins, dice, MDPs, radio-active decays, etc. Maybe one should admit a grayscale of more or less subjective probabilities.

8.6.3 Free will versus Determinism

For illustrational purpose we replace determinism with computability.

The paradox. If the brain of a human were computable we could predict the action of the human with a computer. If we tell the human his action in advance he is forced to perform this action and hence loses his free will. Assuming humans have free will refutes the computability assumption of the brain.

This paradox between computability and free will is sometimes used as an argument against the possibility of AI. However, it vanishes by a more careful reasoning:

³Using Bayes-mixtures as subjective probabilities did not need any (e.g. Cox's) axioms for justification, but received a decision-theoretic justification.

That a part of the universe is computable is defined as follows:

Assumption 1. Given a box (part of the universe) in state s at time t we can compute the next (or some farther future) state s' at time $t' > t$ if there is no interaction between the box and the rest of the universe during time $t...t'$.

Without this independence assumption in time-interval $[t, t']$ the possibility of correct prediction cannot be guaranteed.

Assumption 2. Assume that the brain is computable. It receives input x at time t and computes action y at time t' . During the thinking period $[t, t']$ it is completely separated from the environment.

After input x , the brain B is in a state s and Assumption 1 applies, i.e. we can compute, say with algorithm $p: \mathcal{X} \rightarrow \mathcal{Y}$, the brain's decision y . We can't inform the brain in period $[t, t']$ of this decision without violating Assumption 2. We are free to hand out y in a closed envelope to B . After B has made its decision, he is allowed to open the letter and realizes that his decision was predictable. Such an experiment will have enormous psychological, social, and legal consequences, and looks paradoxical, but does not lead to any contradictions!

Assume we allow intermittent interaction, then the brain B' maps input (x, y) to an action y' , which is possibly different from y . There is no contradiction, since p maps \mathcal{X} to \mathcal{Y} whereas B' maps $\mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{Y}$, so these functions have nothing to do with each other.

Consider a variant of the paradox, where B itself reliably predicts/precomputes its own action, and then “decides” to deviate from its own prediction. More formally, the assumption is that a part B_2 of brain $B \equiv B_1$ can simulate B_1 . By assumption B_2 is functionally identical to B_1 . Only for illustrational purposes we make the further assumption that B_2 also operates identically to B_1 in the sense that B_2 itself contains a part, say B_3 which simulates B_2 , etc. We have an infinite recursion. The first question is *not what* the output of B is and whether it is finitely computable but whether this infinite recursion has a value *at all*. What we need is a fixed point. Insert a function f into brain B (as a possible candidate for B_2) and test whether B computes the same function. If it does, then f (and hence B) is a fixed point of the recursion. If such a fixed point exists (and is unique) we may define the value of the infinite recursion as this fixed point value. Finally we would have to check whether this fixed point can be found by a finite algorithm. It is well known that not every recursion $y = f(y)$ has a fixed point. The paradox in our case is just that we implicitly assumed the existence of a (unique) fixed point. The paradox is resolved by noting that this fixed point simply does not exist. Sometimes fixed points can be found by iteration. One starts with some value y_1 for y and iterates $y_2 = f(y_1)$, ..., $y_n = f(y_{n-1})$. If the limit y_∞ exists, then it is a fixed point. Assume our function B acts with $y = 1$ if B_2 predicts $y = 0$ and vice versa. In this case $y_n = 1 - y_{n-1}$ oscillates and y_∞ does not exist. (In the case of a binary decision this proves that a fixed point does not exist). So this paradox is about non-existent fixed-points. A self-contradictory brain simply does not exist (when starting from Assumptions 1

and 2). It is not possible to set up a brain with a part predicting its own behavior reliably in every situation. The same analysis holds for an infinite regression of *external* predictors, telling the human his action in advance.

Note that neither the paradox, nor the solution has anything specific to do with *computable functions*. We could have formulated the paradox in terms of general mathematical functions (mappings).

8.6.4 The Big Questions

On non-computable physics & brains. There are two possible objections to AI in general and, therefore, to AIXI in particular. Non-computable physics (which is not too odd) could make Turing computable AI impossible. As at least the world that is relevant for humans seems mainly to be computable we do not believe that it is necessary to integrate non-computable devices into an AI system. The (clever and nearly convincing) ‘Gödel’ argument by Penrose [Pen89, Pen94] (refining Lucas [Luc61]) that non-computational physics *must* exist and *is* relevant to the brain, has (in our opinion convincing) loopholes.

Evolution & the number of wisdom. A more serious problem is the evolutionary information gathering process. It has been shown that the ‘number of wisdom’ Ω contains a very compact tabulation of 2^n undecidable problems in its first n binary digits [Cha91]. Ω is only enumerable with computation time increasing more rapidly with n than any recursive function. The enormous computational power of evolution could have developed and coded something like Ω into our genes, which significantly guides human reasoning. In short: Intelligence could be something complicated and evolution toward it from an even cleverly designed algorithm of size $O(1)$ could be too slow. As evolution has already taken place we could add the information from our genes or brain structure to any/our AI system, but this means that the important part is still missing and that it is principally impossible to derive an efficient algorithm from a simple formal definition of AI.

Consciousness. For what is probably the *biggest question*, that of *consciousness*, we want to give a physical analogy. Quantum (field) theory is the most accurate and universal physical theory ever invented. Although already developed in the 1930s the *big* question, regarding the interpretation of the wave function collapse, is still open. Although extremely interesting from a philosophical point of view, it is completely irrelevant from a practical point of view⁴. We believe the same to be valid for *consciousness* in the field of Artificial Intelligence: Philosophically highly interesting but practically unimportant. Whether consciousness *will* be explained some day is another question.

⁴In the Theory of Everything, the collapse might become of “practical” importance and must or will be solved.

8.7 Conclusions

All tasks which require intelligence to be solved can naturally be formulated as a maximization of some expected utility in the framework of agents. We presented a functional (4.7) and an iterative (4.17) formulation of such a decision theoretic agent in Chapter 4, which is general enough to cover all AI problem classes, as has been demonstrated by several examples. The main remaining problem is the unknown prior probability distribution μ of the environment(s). Conventional learning algorithms are unsuitable, because they can neither handle large (unstructured) state spaces, nor do they converge in the theoretically minimal number of cycles, nor can they handle non-stationary environments appropriately. On the other hand, Solomonoff's universal semimeasure $M \triangleq \xi_U$ (2.24), based on ideas from algorithmic information theory, solves the problem of the unknown prior distribution for induction problems as has been demonstrated in Chapters 2 and 3. No explicit learning procedure is necessary, as ξ_U automatically converges to μ . We unified the theory of universal sequence prediction with the decision theoretic agent by replacing the unknown true prior μ by an appropriately generalized universal semimeasure ξ in Chapter 5. We gave various arguments that the resulting AIXI model is the most intelligent, parameter-free and environmental/application independent model possible. We defined an intelligence order relation (5.14) to give a rigorous meaning to this claim. Furthermore, possible solutions to the horizon problem have been discussed. In Chapter 6 we outlined how the AIXI model solves various problem classes. These included sequence prediction, strategic games, function minimization and, especially, learning to learn supervised. The list could easily be extended to other problem classes like classification, function inversion and many others. The major drawback of the AIXI model is that it is uncomputable, or more precisely, only asymptotically computable, which makes an implementation impossible. To overcome this problem, we constructed a modified model $\text{AIXI}t_l$, which is still effectively more intelligent than any other time t and length l bounded algorithm (Section 7.2). The computation time of $\text{AIXI}t_l$ is of the order $t \cdot 2^l$. A way of overcoming the large multiplicative 2^l constant has been presented at the expense of an (unfortunately even larger) additive constant (Section 7.1). Possible further research has been discussed. The main directions could be to prove general and special reward bounds, use AIXI as a super model and explore its relation to other specialized models and finally improve performance with or without giving up universality.

Bibliography

- [ACBFS95] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. Gambling in a rigged casino: The adversarial multi-armed bandit problem. In *Proceedings of the 36th Annual Symposium on Foundations of Computer Science (FOCS 1995)*, pages 322–331, Los Alamitos, CA, 1995. IEEE Computer Society Press. 535
- [Açz66] J. Aćzel. *Lectures on Functional Equations and Their Applications*. Academic Press, New York, 1966. 214, 228
- [AG00] P. Auer and C. Gentile. Adaptive and self-confident on-line learning algorithms. In *Proceedings of the 13th Conference on Computational Learning Theory*, pages 107–117. Morgan Kaufmann, San Francisco, 2000. 342
- [AS83] D. Angluin and C. H. Smith. Inductive inference: Theory and methods. *ACM Computing Surveys*, 15(3):237–269, 1983. 126, 228, 351
- [Bab91] L. Babai et al. Checking computations in polylogarithmic time. *STOC: 23rd ACM Symp. on Theory of Computation*, 23:21–31, 1991. 712
- [Bar00] A. R. Barron. Limits of information, markov chains, and projection. In *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, pages 25–25, Sorrento, Italy, 2000. 352
- [Bau99] Eric B. Baum. Toward a model of intelligence as an economy of agents. *Machine Learning*, 35(2):155–185, 1999. 804
- [Bay63] T. Bayes. An essay towards solving a problem in the doctrine of chances. *Philos. Trans. Royal Soc.*, 53:376–398, 1763. 203
- [BB01] J. Baxter and P. L. Bartlett. Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research*, 15:319–350, 2001. 804
- [BD62] D. Blackwell and L. Dubins. Merging of opinions with increasing information. *Annals of Mathematical Statistics*, 33:882–887, 1962. 351
- [BEHW87] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. K. Warmuth. Occam’s razor. *Information Processing Letters*, 24(6):377–380, April 1987. 228
- [BEHW89] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. K. Warmuth. Learnability and the Vapnik-Chervonenkis dimension. *Journal of the ACM*, 36(4):929–965, October 1989. 228
- [Bel57] R. E. Bellman. *Dynamic Programming*. Princeton University Press, New Jersey, 1957. 112, 402, 415
- [Ber13] J. Bernoulli. *Ars Conjectandi*. Thurnisiorum, Basel, 1713. [Reprinted in: “Die Werke von Jakob Bernoulli”, pages 106–286, volume 3, Birkhäuser Verlag, Basel, 1975 – and – “A Source Book in Mathematics”, pages 85–90, Dover, New York, 1959. English translation of part IV (with limit theorem) by Bing Sung, Harvard University Dept. of Statistics, Technical Report #2, 1966]. 212

- [Ber95a] D. P. Bertsekas. *Dynamic Programming and Optimal Control, Vol. (I)*. Athena Scientific, Belmont, Massachusetts, 1995. 126, 416
- [Ber95b] D. P. Bertsekas. *Dynamic Programming and Optimal Control, Vol. (II)*. Athena Scientific, Belmont, Massachusetts, 1995. 126, 529, 538
- [BF85] D. A. Berry and B. Fristedt. *Bandit Problems: Sequential Allocation of Experiments*. Chapman and Hall, London, 1985. 535
- [BGHK92] F. Bacchus, A. Grove, J. Y. Halpern, and D. Koller. From statistics to beliefs. In *Proceedings of the Tenth National Conference on Artificial Intelligence (AAAI-92)*, pages 602–608, San Jose, California, July 1992. AAAI Press. 226
- [Blu67] M. Blum. A machine-independent theory of the complexity of recursive functions. *Journal of the ACM*, 14(2):322–336, 1967. 120, 122, 702, 711
- [Blu71] M. Blum. On effective procedures for speeding up algorithms. *Journal of the ACM*, 18(2):290–305, 1971. 120, 122, 702, 711
- [BM98] A. A. Borovkov and A. Moullagaliev. *Mathematical Statistics*. Gordon & Breach, 1998. 336, 352
- [BS84] B. G. Buchanan and E. H. Shortliffe. *Rule-Based Expert Systems: The MYCIN Experiments of the Stanford Heuristic Programming Project*. Addison Wesley, Reading, Massachusetts, 1984. 227
- [BSA83] A. G. Barto, R. S. Sutton, and C. W. Anderson. Neuronlike adaptive elements that can solve difficult learning control problems. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-13:834–846, 1983. 126
- [BT96] D. P. Bertsekas and J. N. Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, Belmont, MA, 1996. 102, 112, 126, 127, 402, 415, 524, 526, 804
- [Can74] G. Cantor. Über eine Eigenschaft des Inbegriffs aller reellen algebraischen Zahlen. *Journal für reine und angew. Math.*, 77:258–262, 1874. [English translation: “On a property of the set of real algebraic numbers”. In W. B. Ewald, editor, *A Source Book in the Foundations of Mathematics*, volume 2, pages 839–843, Oxford, Clarendon Press].
- [Car63] G. Cardano. *Liber de ludo aleae*, 1565/1663. Published in 1663 but completed already around 1565.
- [Car48] R. Carnap. On the application of inductive logic. *Philosophy and Phenomenological Research*, 8:133–148, 1948. 226
- [Car50] R. Carnap. *Logical Foundations of Probability*. University of Chicago Press, Chicago, 1950. 226
- [CB90] B. S. Clarke and A. R. Barron. Information-theoretic asymptotics of Bayes methods. *IEEE Transactions on Information Theory*, 36:453–471, 1990. 110, 303, 340, 341
- [CB97] N. Cesa-Bianchi et al. How to use expert advice. *Journal of the ACM*, 44(3):427–485, 1997. 303, 342, 804
- [CBL01] N. Cesa-Bianchi and G. Lugosi. Worst-case bounds for the logarithmic loss of predictors. *Machine Learning*, 43(3):247–264, 2001. 344
- [Cha66] G. J. Chaitin. On the length of programs for computing finite binary sequences. *Journal of the ACM*, 13(4):547–569, 1966. 126, 225, 710
- [Cha69] G. J. Chaitin. On the length of programs for computing finite binary sequences: Statistical considerations. *Journal of the ACM*, 16(1):145–159, 1969. 208, 225

- [Cha75] G. J. Chaitin. A theory of program size formally identical to information theory. *Journal of the ACM*, 22(3):329–340, 1975. 104, 126, 225, 226, 317
- [Cha91] G. J. Chaitin. Algorithmic information and evolution. in *O.T. Solbrig and G. Nicolis, Perspectives on Biological Complexity*, IUBS Press, pages 51–60, 1991. 126, 229, 817
- [Che85] P. Cheeseman. In defense of probability. In A. Joshi, editor, *Proceedings of the Ninth International Joint Conference on Artificial Intelligence*, pages 1002–1009, Los Altos, California, 1985. Morgan Kaufmann. 227, 622
- [Che88] P. Cheeseman. An inquiry into computer understanding. *Computational Intelligence*, 4(1):58–66, 1988. 227, 622
- [Chu40] A. Church. On the concept of a random sequence. *Bulletin of the American Mathematical Society*, 46:130–135, 1940. 226
- [Chu86] P. S. Churchland. *Neurophilosophy: Toward a Unified Science of the Mind-Brain*. MIT Press, Cambridge, Massachusetts, 1986. 814
- [Con97] M. Conte et al. Genetic programming estimates of Kolmogorov complexity. In *Genetic Algorithms: Proceedings of the 17th International Conference*, pages 743–750, 1997. 126, 229
- [Cov74] T. M. Cover. Universal gambling schemes and the complexity measures of Kolmogorov and Chaitin. Technical Report 12, Statistics Department, Stanford University, Stanford, CA, 1974. 226
- [Cox46] R. T. Cox. Probability, frequency, and reasonable expectation. *American Journal of Physics*, 14(1):1–13, 1946. 105, 214, 226, 227
- [Csi67] I. Csiszr. Information-type measures of difference of probability distributions and indirect observations. *Studia Sci. Math. Hungar*, 2:299–318, 1967. 352
- [CT91] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. Wiley Series in Telecommunications. John Wiley & Sons, New York, NY, USA, 1991. 226, 352
- [CW82] D. Coppersmith and S. Winograd. On the asymptotic complexity of matrix multiplication. *SIAM Journal on Computing*, 11(3):472–492, August 1982. 705
- [CW90] D. Coppersmith and S. Winograd. Matrix multiplication via arithmetic progressions. *Journal of Symbolic Computation*, 9(3):251–280, March 1990. 705
- [Dal73] R. P. Daley. Minimal-program complexity of sequences with restricted resources. *Information and Control*, 23(4):301–312, 1973. 126, 226, 714
- [Dal77] R. P. Daley. On the inference of optimal descriptions. *Theoretical Computer Science*, 4(3):301–319, 1977. 126, 226, 714
- [Dau90] J. W. Dauben. *Georg Cantor: His mathematics and philosophy of the infinite*. Princeton University Press, Princeton, 1990.
- [Daw84] A. P. Dawid. Statistical theory. The prequential approach. *J.R. Statist. Soc. A*, 147:278–292, 1984. 103, 203, 228, 342
- [Dem68] A. P. Dempster. A generalization of Bayesian inference. *Journal of the Royal Statistical Society*, 30 (Series B):205–247, 1968. 227
- [Doo53] J. L. Doob. *Stochastic Processes*. John Wiley & Sons, New York, 1953. 305, 311, 521
- [EF98] G. W. Erickson and J. A. Fossa. *Dictionary of Paradox*. University Press of America, Lanham, MD, 1998. 535

- [Fel68] W. Feller. *An Introduction to Probability Theory and its Applications*. John Wiley & Sons, New York, 3 edition, 1968. 226
- [Fer67] T. S. Ferguson. *Mathematical Statistics: A Decision Theoretic Approach*. Academic Press, New York, 3rd edition, 1967. 336
- [Fin37] B. de Finetti. Le prévision: ses lois logiques, ses sources subjectives. *Ann. Inst. Poincaré*, 7:1–68, 1937. [English translation: “Foresight: Its logical laws, its subjective sources” in *Studies in Subjective Probability*, H. E. Kyburg and H. E. Smokler, editors. Krieger, New York, pages 55–118, 1980]. 226
- [Fin73] T. L. Fine. *Theories of Probability*. Academic Press, 1973. 227, 228
- [Fis22] R. A. Fisher. On the mathematical foundations of theoretical statistics. *Philosophical Transactions of the Royal Society of London*, Series A 222:309–368, 1922. 226
- [Fit96] Melvin C. Fitting. *First-Order Logic and Automated Theorem Proving*. Graduate Texts in Computer Science. Springer-Verlag, Berlin, 2nd edition, 1996. 707
- [FMG92] M. Feder, N. Merhav, and M. Gutman. Universal prediction of individual sequences. *IEEE Transactions on Information Theory*, 38:1258–1270, 1992. 124, 126, 226, 317, 318, 811
- [FS97] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997. 343
- [FT91] D. Fudenberg and J. Tirole. *Game Theory*. The MIT Press, Cambridge, Massachusetts, 1991. 608
- [Gác74] P. Gács. On the symmetry of algorithmic information. *Russian Academy of Sciences Doklady. Mathematics (formerly Soviet Mathematics–Doklady)*, 15:1477–1480, 1974. 104, 225, 317
- [Gal68] Robert G. Gallager. *Information Theory and Reliable Communication*. John Wiley & Sons, New York, NY, USA, 1968. 225
- [GCSR95] A. Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin. *Bayesian Data Analysis*. Chapman, 1995. 227
- [GH02] P. D. Grünwald and J. Y. Halpern. Updating probabilities. In *Proceedings of the 18th Conference on Uncertainty in Artificial Intelligence (UAI-2002)*, pages 187–196, San Francisco, CA, 2002. Morgan Kaufmann. 535
- [Gin87] M. L. Ginsberg, editor. *Readings in Nonmonotonic Reasoning*. Morgan Kaufmann, Los Altos, California, 1987. 227
- [Git89] J. C. Gittins. *Multi-Armed Bandit Allocation Indices*. John Wiley & Sons, 1989. 535
- [GJ74] J. C. Gittins and D. M. Jones. A dynamic allocation index for the sequential design of experiments. *Progress in Statistics*, pages 241–266, 1974. 535
- [Göd31] K. Gödel. Über formal unentscheidbare Sätze der principia mathematica und verwandter systeme I. *Monatshefte für Mathematik und Physik*, 38:173–198, 1931. [English translation by E. Mendelsohn: “On undecidable propositions of formal mathematical systems”. In M. Davis, editor, *The undecidable*, pages 39–71, New York, 1965. Raven Press, Hewlett]. 225
- [Grü98] P. D. Grünwald. *The Minimum Description Length Principle and Reasoning under Uncertainty*. PhD thesis, Universiteit van Amsterdam, 1998. 228, 351

- [GTV01] P. Gács, J. Tromp, and P. M. B. Vitányi. Algorithmic statistics. *IEEE Transactions on Information Theory*, 47(6):2443–2463, 2001. 203, 344
- [Hac75] I. Hacking. *The Emergence of Probability*. Cambridge University Press, Cambridge, 1975. 226
- [Hal90] A. Hald. *A History of Probability and Statistics and Their Applications Before 1750*. Wiley, New York, 1990. 226
- [Hal99] Joseph Y. Halpern. A counterexample to theorems of Cox and Fine. *Journal of AI research*, 10:67–85, 1999. 228
- [Har79] J. Hartmanis. Relations between diagonalization, proof systems, and complexity gaps. *Theoretical Computer Science*, 8(2):239–253, April 1979. 711
- [Hec88] D. E. Heckerman. An axiomatic framework for belief updates. In John F. Lemmer and Laveen N. Kanal, editors, *Uncertainty in artificial intelligence 2*, volume 5 of *Machine intelligence and pattern recognition*, pages 11–22, Amsterdam, 1988. North-Holland. 228
- [HHL86] E. J. Horvitz, D. E. Heckerman, and C. P. Langlotz. A framework for comparing alternative formalisms for plausible reasoning. In *Proceedings of the Fifth National Conference on Artificial Intelligence (AAAI-86)*, volume 1, pages 210–214, Philadelphia, Pennsylvania, August 1986. Morgan Kaufmann. 228
- [HKW98] D. Haussler, J. Kivinen, and M. K. Warmuth. Sequential prediction of individual sequences under general loss functions. *IEEE Transactions on Information Theory*, 44(5):1906–1925, 1998. 342
- [HMU01] J. E. Hopcroft, R. Motwani, and J. D. Ullman. *“Introduction to Automata Theory, Language, and Computation”*. Addison–Wesley, 2nd edition edition, 2001. 126, 205, 230
- [Hug89] R. I. G. Hughes. *Structure and Interpretation of Quantum Mechanics*. Harvard Univ Press, 1989. 227
- [Hum39] D. Hume. *A Treatise of Human Nature, Book I*. Edited version by L. A. Selby-Bigge and P. H. Nidditch, Oxford University Press, 1978., 1739. 203
- [Hut00] M. Hutter. A theory of universal artificial intelligence based on algorithmic complexity. Technical Report cs.AI/0004001, München, 62 pages, 2000. <http://arxiv.org/abs/cs.AI/0004001>. 127
- [Hut01a] M. Hutter. Convergence and error bounds of universal prediction for general alphabet. *Proceedings of the 12th European Conference on Machine Learning (ECML-2001)*, pages 239–250, 2001. 127, 352
- [Hut01b] M. Hutter. General loss bounds for universal sequence prediction. In C. E. Brodley and A. P. Danyluk, editors, *Proceedings of the 18th International Conference on Machine Learning (ICML-2001)*, pages 210–217, Manno(Lugano), CH, 2001. Morgan Kaufmann. 127, 527
- [Hut01c] M. Hutter. New error bounds for Solomonoff prediction. *Journal of Computer and System Sciences*, 62(4):653–667, 2001. 4, 127, 309, 319, 355
- [Hut01d] M. Hutter. Towards a universal theory of artificial intelligence based on algorithmic probability and sequential decisions. *Proceedings of the 12th European Conference on Machine Learning (ECML-2001)*, pages 226–238, 2001. 127, 344, 704

- [Hut01e] M. Hutter. Universal sequential decisions in unknown environments. *Proceedings of the 5th European Workshop on Reinforcement Learning (EWRL-5)*, 27:25–26, 2001. 127
- [Hut02a] M. Hutter. The fastest and shortest algorithm for all well-defined problems. *International Journal of Foundations of Computer Science*, 13(3):431–443, 2002. 127
- [Hut02b] M. Hutter. Optimality of universal Bayesian prediction for general loss and alphabet. Technical Report IDSIA-02-02, Istituto Dalle Molle di Studi sull'Intelligenza Artificiale (IDSIA), Manno(Lugano), Switzerland, 2002. 127
- [Hut02c] M. Hutter. Self-optimizing and Pareto-optimal policies in general environments based on Bayes-mixtures. In *Proceedings of the 15th Annual Conference on Computational Learning Theory (COLT 2002)*, Lecture Notes in Artificial Intelligence, pages 364–379, Sydney, Australia, 2002. Springer. 127
- [Jay78] E. T. Jaynes. Where do we stand on maximum entropy? In R. D. Levine and M. Tribus, editors, *The Maximum Entropy Formalism*, pages 15–118. MIT Press, Cambridge, MA, 1978. 228
- [Jay96] E. T. Jaynes. *Probability theory: the logic of science*. online, 1996. 226, 228
- [Jef83] R. C. Jeffrey. *The Logic of Decision*. University of Chicago Press, Chicago, Illinois, second edition, 1983. 226
- [Key21] J. M. Keynes. *A Treatise on Probability*. Macmillan, London, 1921. 226
- [KHS01a] I. Kwee, M. Hutter, and J. Schmidhuber. Gradient-based reinforcement planning in policy-search methods. *Proceedings of the 5th European Workshop on Reinforcement Learning (EWRL-5)*, 27:27–29, 2001. 127, 804
- [KHS01b] I. Kwee, M. Hutter, and J. Schmidhuber. Market-based reinforcement learning in partially observable worlds. *Proceedings of the International Conference on Artificial Neural Networks (ICANN-2001)*, pages 865–873, 2001. 127, 804
- [KLC98] L. P. Kaelbling, M. L. Littman, and A. R. Cassandra. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 101:99–134, 1998. 126
- [Kle36] S. Kleene. General recursive functions of natural numbers. *Mathematische Annalen*, 112:727–742, 1936. 225
- [KLM96] L. P. Kaelbling, M. L. Littman, and A. W. Moore. Reinforcement learning: a survey. *Journal of AI research*, 4:237–285, 1996. 112, 126, 415
- [Knu73] D. E. Knuth. *The Art of Computer Programming, Vol. I: Fundamental Algorithms*. Addison-Wesley, Reading, MA, 1973. 4
- [Ko86] K.-I. Ko. On the notion of infinite pseudorandom sequences. *Theoretical Computer Science*, 48(1):9–33, 1986. 126, 226, 714
- [Kol33] A. N. Kolmogorov. *Grundlagen der Wahrscheinlichkeitsrechnung*. Springer-Verlag, Berlin, 1933. [English translation: *Foundations of the theory of probability*. Chelsea, New York, second edition, 1956]. 226
- [Kol63] A. N. Kolmogorov. On tables of random numbers. *Sankhya, the Indian Journal of Statistics*, Series A 25, 1963. 226
- [Kol65] A. N. Kolmogorov. Three approaches to the quantitative definition of information. *Problems of Information and Transmission*, 1(1):1–7, 1965. 104, 126, 208, 225, 226, 317, 710

- [Kol83] A. N. Kolmogorov. Combinatorial foundations of information theory and the calculus of probabilities. *Russian Mathematical Surveys*, 38(4):27–36, 1983. 225, 226
- [Kra49] L. G. Kraft. A device for quantizing, grouping and coding amplitude modified pulses. Master’s thesis, Electrical Engineering Department, Massachusetts Institute of Technology, Cambridge, MA, 1949. 225, 708
- [KS98] M. J. Kearns and S. Singh. Near-optimal reinforcement learning in polynomial time. In *Proc. 15th International Conf. on Machine Learning*, pages 260–268. Morgan Kaufmann, San Francisco, CA, 1998. 126, 527
- [KU63] A. N. Kolmogorov and V. A. Uspenskii. On the definition of an algorithm. *American Mathematical Society Translations*, 29:216–245, 1963. Translated from Russian Original Uspekhi Mat. Nauk. 13(4):3–28, 1958. 710
- [KU87] A. N. Kolmogorov and V. A. Uspenskii. Algorithms and randomness. *Theory of Probability and its Applications*, 3(32):389–412, 1987. 229
- [KV86] P. R. Kumar and P. P. Varaiya. *Stochastic Systems: Estimation, Identification, and Adaptive Control*. Prentice Hall, Englewood Cliffs, NJ, 1986. 115, 127, 204, 507, 509, 516, 524, 526, 535
- [KW99] J. Kivinen and M. K. Warmuth. Averaging expert predictions. In P. Fischer and H. U. Simon, editors, *Proceedings of the 4th European Conference on Computational Learning Theory (Eurocolt-99)*, volume 1572 of *LNAI*, pages 153–167, Berlin, 1999. Springer. 342, 343
- [Kyb77] H. E. Kyburg. Randomness and the right reference class. *The Journal of Philosophy*, 74(9):501–521, 1977. 226
- [Kyb83] H. E. Kyburg. The reference class. *Philosophy of Science*, 50:374–397, 1983. 226
- [Lam87] M. van Lambalgen. *Random Sequences*. PhD thesis, University of Amsterdam, 1987. 225
- [Lap14] P. Laplace. Théorie analytique des probabilités, 1814. English translation: “A Philosophical Essay on Probabilities”, F. W. Truscott & F. L. Emory, Dover, pages 16–17, 1952.
- [Lev73a] L. A. Levin. On the notion of a random sequence. *Soviet Math. Dokl.*, 14(5):1413–1416, 1973. 224, 225
- [Lev73b] L. A. Levin. Universal sequential search problems. *Problems of Information Transmission*, 9:265–266, 1973. 120, 124, 126, 225, 228, 317, 704
- [Lev74] L. A. Levin. Laws of information conservation (non-growth) and aspects of the foundation of probability theory. *Problems of Information Transmission*, 10(3):206–210, 1974. 104, 225, 226, 317
- [Lev84] L. A. Levin. Randomness conservation inequalities: Information and independence in mathematical theories. *Information and Control*, 61:15–37, 1984. 120, 124, 317, 704
- [Lov69a] D. W. Loveland. On minimal-program complexity measures. In ACM, editor, *First ACM Symposium on Theory of Computing*, pages 61–78, New York, 1969. ACM Press. 226
- [Lov69b] D. W. Loveland. A variant of the Kolmogorov concept of complexity. *Information and Control*, 15(6):510–526, 1969. 226
- [Luc61] J. R. Lucas. Minds, machines, and Gödel. *Philosophy*, 36:112–127, 1961. 713, 814, 817

- [LV77] L. A. Levin and V. V. V'yugin. Invariant properties of informational bulks. In *Proceedings of the 6th Symposium on Mathematical Foundations of Computer Science*, volume 53 of *LNCS*, pages 359–364. Springer, 1977. 228
- [LV91] M. Li and P. M. B. Vitányi. Learning simple concepts under simple distributions. *SIAM Journal on Computing*, 20(5):911–935, 1991. 122, 126, 228, 714
- [LV92a] M. Li and P. M. B. Vitányi. Inductive reasoning and Kolmogorov complexity. *Journal of Computer and System Sciences*, 44:343–384, 1992. 124, 126, 228, 351, 504, 811
- [LV92b] M. Li and P. M. B. Vitányi. Philosophical issues in Kolmogorov complexity (invited lecture). In W. Kuich, editor, *Proceedings on Automata, Languages and Programming (ICALP '92)*, volume 623 of *LNCS*, pages 1–15, Berlin, Germany, 1992. Springer. 126, 228
- [LV97] M. Li and P. M. B. Vitányi. *An introduction to Kolmogorov complexity and its applications*. Springer, 2nd edition, 1997. 7, 102, 122, 123, 126, 202, 204, 217, 221, 222, 223, 225, 226, 228, 229, 230, 303, 304, 305, 309, 311, 312, 317, 328, 336, 351, 352, 532, 536, 703, 704, 708, 710, 714
- [LW89] N. Littlestone and M. K. Warmuth. The weighted majority algorithm. In *30th Annual Symposium on Foundations of Computer Science*, pages 256–261, Research Triangle Park, North Carolina, 1989. IEEE. 342
- [LW94] N. Littlestone and M. K. Warmuth. The weighted majority algorithm. *Information and Computation*, 108(2):212–261, 1994. 126, 228, 342, 713
- [MA93] A. W. Moore and C. G. Atkeson. Prioritized sweeping: Reinforcement learning with less data and less time. *Machine Learning*, 13:103–130, 1993. 126
- [McC80] J. McCarthy. Circumscription—A form of non-monotonic reasoning. *Artificial Intelligence*, 13(1–2):27–39, 1980. 227
- [McC95] A. K. McCallum. Instance-based utile distinctions for reinforcement learning with hidden state. In *Proceedings of the 12th International Conference on Machine Learning*, pages 387–395, 1995. 804
- [MD80] D. McDermott and J. Doyle. Nonmonotonic logic 1. *Artificial Intelligence*, 13:41–72, 1980. 227
- [MF98] N. Merhav and M. Feder. Universal prediction. *IEEE Transactions on Information Theory*, 44(6):2124–2147, 1998. 303, 304, 323, 341, 342, 351
- [Mic66] D. Michie. Game-playing and game-learning automata. In L. Fox, editor, *Advances in Programming and Non-Numerical Computation*, pages 183–200. Pergamon, New York, 1966. 126, 410
- [Mis19] R. von Mises. Grundlagen der wahrscheinlichkeitsrechnung. *Mathematische Zeitschrift*, 5:52–99, 1919. Correction, *Ibid.*, volume 6, 1920, [English translation in: *Probability, Statistics, and Truth*, Macmillan, 1939]. 226
- [Mis28] R. von Mises. *Wahrscheinlichkeit, Statistik und Wahrheit*. J. Springer, Berlin, 1928. [English translation: *Probability, Statistics, and Truth*, Allen and Unwin, London, 1957]. 226
- [ML66] P. Martin-Löf. The definition of random sequences. *Information and Control*, 9(6):602–619, 1966. 212, 226, 228
- [ML69] P. Martin-Löf. The definition of random sequences. *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 19:225–230, 1969. 228

- [Mor88] H. Moravec. *Mind Children: The Future of Robot and Human Intelligence*. Harvard University Press, Cambridge, 1988. 814
- [Mos65] F. Mosteller. *Fifty Challenging Problems in Probability with Solutions*. Addison-Wesley, Reading, MA, 1965. 535
- [NM44] J. Von Neumann and O. Morgenstern. *Theory of games and economic behavior*. Princeton University Press, New Jersey, 1944. 126, 402, 608
- [Odi89] Piergiorgio Odifreddi. *Classical Recursion Theory (Volume 1)*. North-Holland, Amsterdam, 1989. 225
- [Odi99] Piergiorgio Odifreddi. *Classical Recursion Theory (Volume 2)*. Elsevier, Amsterdam, 1999. 225
- [OR94] M. J. Osborne and A. Rubenstein. *A Course in Game Theory*. The MIT Press, Cambridge, Massachusetts, 1994. 416, 608, 805
- [Par95] J. B. Paris. *The Uncertain Reasoner's Companion: A Mathematical Perspective*. Cambridge University Press, Cambridge, England, 1995. 214, 228
- [Pas54] B. Pascal. Letters to Fermat, 1654.
- [Pen89] R. Penrose. *The Emperor's New Mind*. Oxford U. P., 1989. 713, 814, 817
- [Pen94] R. Penrose. *Shadows of the mind, A search for the missing science of consciousness*. Oxford Univ. Press, 1994. 126, 713, 814, 817
- [PF97] X. Pintado and E. Fuentes. A forecasting algorithm based on information theory. Technical report, Centre Universitaire d'Informatique, University of Geneva, 1997. 126, 226
- [Pin64] M. S. Pinsker. *Information and Information Stability of Random Variables and Processes*. Holden-Day, San-Francisco, CA, 1964. Russian original, Izd. Akad. Nauk, 1960. 352
- [Pop34] K. R. Popper. *Logik der Forschung*. Springer, Berlin, 1934. [English translation: *The Logic of Scientific Discovery* Basic Books, New York, 1959 – and – Hutchinson, London, revised edition, 1968]. 226
- [Pos44] E. L. Post. Recursively enumerable sets of positive integers and their decision problems. *Bulletin of the American Mathematical Society*, 50:284–316, 1944. 225
- [Put63] H. Putnam. 'Degree of confirmation' and inductive logic. In P. A. Schilpp, editor, *The Philosophy of Rudolf Carnap*. Open Court, La Salle, Illinois, 1963. 226
- [Ram31] F. P. Ramsey. Truth and probability. In R. B. Braithwaite, editor, *The Foundations of Mathematics: Collected Papers of Frank P. Ramsey*, pages 156–198. Routledge and Kegan Paul, London, 1931. 226
- [Rei49] H. Reichenbach. *The Theory of Probability: An Inquiry into the Logical and Mathematical Foundations of the Calculus of Probability*. University of California Press, Berkeley and Los Angeles, second edition, 1949. 226, 228
- [Rei80] R. Reiter. A logic for default reasoning. In *Artificial Intelligence*, pages 81–132, 1980. 227
- [Res01] N. Rescher. *Paradoxes: Their Roots, Range, and Resolution*. Open Court Pub Co (Sd), Lanham, MD, 2001. 535
- [Rin94] M. Ring. *Continual Learning in Reinforcement Environments*. PhD thesis, University of Texas at Austin, Austin, Texas., 1994. 804
- [Ris78] J. J. Rissanen. Modeling by shortest data description. *Automatica*, 14:465–471, 1978. 228

- [Ris89] J. J. Rissanen. *Stochastic Complexity in Statistical Inquiry*. World Scientific Publ. Co., 1989. 124, 126, 209, 228, 342, 811
- [Ris96] J. J. Rissanen. Fisher Information and Stochastic Complexity. *IEEE Trans on Information Theory*, 42(1):40–47, January 1996. 341
- [RN95] S. J. Russell and P. Norvig. *Artificial Intelligence. A Modern Approach*. Prentice-Hall, Englewood Cliffs, 1995. 110, 112, 126, 227, 401, 410, 415, 608, 804
- [Rob52] H. Robbins. Some aspects of the sequential design of experiments. *Bulletin of the American Mathematical Society*, 58:527–535, September 1952. 535
- [Rog67] H. Rogers. *Theory of Recursive Functions and Effective Computability*. McGraw-Hill, New York, 1967. 225
- [RV01] A. Robinson and A. Voronkov, editors. *Handbook of Automated Reasoning*. Elsevier Science B.V., 2001. 712
- [Sam59] A. L. Samuel. Some studies in machine learning using the game of checkers. *IBM Journal on Research and Development*, 3:210–229, 1959. 126
- [Sav54] L. J. Savage. *The Foundations of Statistics*. Wiley, New York, 1954. 226
- [SB98] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. Cambridge, MA, MIT Press, 1998. 102, 112, 126, 127, 402, 804
- [Sch71] C. P. Schnorr. *Zufälligkeit und Wahrscheinlichkeit*. Springer, Berlin, 1971. 225
- [Sch73] C. P. Schnorr. Process complexity and effective random tests. *Journal of Computer and System Sciences*, 7(4):376–388, 1973. 225
- [Sch80] A. Schönhage. Storage modification machines. *SIAM Journal on Computing*, 9(3):490–508, 1980. 710
- [Sch95] J. Schmidhuber. Discovering solutions with low Kolmogorov complexity and high generalization capability. In *Proc. 12th International Conference on Machine Learning*, pages 488–496. Morgan Kaufmann, 1995. 124, 126
- [Sch97] J. Schmidhuber. Discovering neural nets with low Kolmogorov complexity and high generalization capability. *Neural Networks*, 10(5):857–873, 1997. 124, 126, 229, 706, 712
- [Sch99] M. Schmidt. Time-bounded Kolmogorov complexity may help in search for extra terrestrial intelligence (SETI). *Bulletin of the European Association for Theoretical Computer Science*, 67:176–180, 1999. 126, 229
- [Sch00] J. Schmidhuber. Algorithmic theories of everything. Report IDSIA-20-00, quant-ph/0011122, IDSIA, Manno (Lugano), Switzerland, 2000. 221, 228, 229, 710
- [Sch02a] J. Schmidhuber. Hierarchies of generalized Kolmogorov complexities and nonenumerable universal measures computable in the limit. *International Journal of Foundations of Computer Science*, 13(4):587–612, 2002. 107, 221, 225, 228, 317, 336, 337, 710
- [Sch02b] J. Schmidhuber. Optimal ordered problem solver. Technical Report IDSIA-12-02, IDSIA, 2002. 124, 126, 228, 229, 706, 712, 804, 812
- [Sch02c] J. Schmidhuber. The Speed Prior: a new simplicity measure yielding near-optimal computable predictions. In J. Kivinen and R. H. Sloan, editors, *Proceedings of the 15th Annual Conference on Computational Learning Theory (COLT 2002)*, Lecture Notes in Artificial Intelligence, pages 216–228, Sydney, Australia, July 2002. Springer. 126, 226, 228, 229, 317, 352, 714, 815
- [Sea80] J. Searle. Minds, brains, and programs. *Behavioral & Brain Sciences*, 3:417–458, 1980. 814

- [SH02] J. Schmidhuber and M. Hutter. Universal learning algorithms and optimal search. *NIPS 2001 Workshop*, 2002. <http://www.idsia.ch/~marcus/idsia/nipsws.htm>. 337
- [Sha48] C. E. Shannon. A mathematical theory of communication. *Bell System Technical Journal*, 27:379–423, 623–656, 1948. 226, 708
- [Sha76] G. Shafer. *A Mathematical Theory of Evidence*. Princeton University Press, Princeton, 1976. 227
- [Sha85] G. Shafer. Conditional probability. *International Statistical Review*, 53(3):261–277, 1985. 535
- [Sho67] J. R. Shoenfield. *Mathematical logic*. Addison-Wesley, 1967. 225, 707
- [Sho76] E. H. Shortliffe. *Computer-Based Medical Consultations: MYCIN*. Elsevier/North-Holland, Amsterdam, London, New York, 1976. 227
- [Sol64] R. J. Solomonoff. A formal theory of inductive inference: Part 1 and 2. *Inform. Control*, 7:1–22, 224–254, 1964. 105, 126, 203, 208, 216, 225, 226, 228, 304, 351, 710
- [Sol78] R. J. Solomonoff. Complexity-based induction systems: comparisons and convergence theorems. *IEEE Trans. Inform. Theory*, IT-24:422–432, 1978. 105, 106, 126, 216, 217, 218, 222, 225, 226, 228, 303, 304, 309, 311, 352, 504
- [Sol86] R. J. Solomonoff. Applications of algorithmic probability to artificial intelligence. In *Uncertainty in Artificial Intelligence*, pages 473–491. Elsevier Science Publishers, 1986. 120, 126, 704, 812
- [Sol97] R. J. Solomonoff. The discovery of algorithmic probability. *Journal of Computer and System Sciences*, 55(1):73–88, 1997. 126, 228, 351
- [Sol99] R. J. Solomonoff. Two kinds of probabilistic induction. *Computer Jnl.*, 42(4):256–259, 1999. 615
- [Sto01] D. Stork. Foundations of Occam’s razor and parsimony in learning. *NIPS 2001 Workshop*, 2001. <http://www.rii.ricoh.com/~stork/OccamWorkshop.html>. 337
- [Str69] V. Strassen. Gaussian elimination is not optimal. *Numerische Mathematik*, 13:354–356, 1969. 705
- [Sut88] R. S. Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*, 3:9–44, 1988. 126
- [SV88] J. I. Seiferas and P. M. B. Vitányi. Counting is easy. *Journal of the ACM*, 35(4):985–1000, 1988. 709
- [Szé86] G. J. Székely. *Paradoxes in Probability Theory and Mathematical Statistics*. D. Reidel, Dordrecht, 1 edition, 1986. 535
- [SZW97] J. Schmidhuber, J. Zhao, and M. A. Wiering. Shifting inductive bias with success-story algorithm, adaptive Levin search, and incremental self-improvement. *Machine Learning*, 28:105–130, 1997. 124, 126, 229, 706, 712, 804
- [Tes94] G. Tesauro. “TD”-Gammon, a self-teaching backgammon program, achieves master-level play. *Neural Computation*, 6(2):215–219, 1994. 126
- [Tri69] M. Tribus. *Rational Descriptions, Decisions and Designs*. Pergamon Press, New York, 1969. 214, 228
- [Tur36] A. M. Turing. On computable numbers, with an application to the Entscheidungsproblem. *Proc. London Math. Soc.*, 2(42):230–265, 1936. 207, 225
- [Tur50] A. M. Turing. Computing machinery and intelligence. *Mind*, October 1950. 814

- [USS90] V. A. Uspenskii, A. L. Semenov, and A. K. Shen. Can an individual sequence of zeros and ones be random? *Russian Mathematical Surveys*, 45, 1990. 229
- [Val84] L. G. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984. 126, 228
- [Vap99] V. N. Vapnik. *The Nature of Statistical Learning Theory*. Springer, New York, second edition, 1999. 203
- [VL00] P. M. B. Vitányi and M. Li. Minimum description length induction, Bayesianism, and Kolmogorov complexity. *IEEE Transactions on Information Theory*, 46(2):446–464, 2000. 228, 230, 312, 344, 710
- [Vog60] W. Vogel. An asymptotic minimax theorem for the two-armed bandit problem. *Ann. Math. Statist.*, 31:444–451, 1960. 535
- [Vov87] V. G. Vovk. On a randomness criterion. *Soviet Mathematics Doklady*, 35(3):656–660, 1987. 229, 312, 354
- [Vov92] V. G. Vovk. Universal forecasting algorithms. *Information and Computation*, 96(2):245–277, 1992. 126, 228, 342, 713
- [Vov99] V. G. Vovk. Competitive on-line statistics. Technical report, CLRC and DoCS, University of London, 1999. 342, 351
- [VV02] N. Vereshchagin and P. M. B. Vitányi. Kolmogorov’s structure functions with an application to the foundations of model selection. Technical report, CWI, Amsterdam, 2002. <http://arxiv.org/abs/cs.CC/0204037>. 203
- [VW98] V. G. Vovk and C. Watkins. Universal portfolio selection. In *Proceedings of the 11th Annual Conference on Computational Learning Theory (COLT-98)*, pages 12–23, New York, 1998. ACM Press. 126, 229
- [Wal37] A. Wald. Die Widerspruchsfreiheit des Kollektivbegriffs in der Wahrscheinlichkeitsrechnung. In *Ergebnisse eines Mathematischen Kolloquiums*, volume 8, pages 38–72, 1937. 226
- [Wal91] P. Walley. *Statistical Reasoning with Imprecise Probabilities*. Chapman and Hall, London, 1991. 227
- [Wan96] Y. Wang. *Randomness and Complexity*. PhD thesis, 1996. 225
- [Wat89] C. Watkins. *Learning from Delayed Rewards*. PhD thesis, King’s College, Oxford, 1989. 126
- [WB68] C. S. Wallace and D. M. Boulton. An information measure for classification. *Computer Jnl.*, 11(2):185–194, August 1968. 228
- [WD92] C. Watkins and P. Dayan. Q-learning. *Machine Learning*, 8:279–292, 1992. 126
- [WM97] D. H. Wolpert and W. G. Macready. No free lunch theorems for optimization. *IEEE Transactions on Evolutionary Computation*, 1(1):67–82, 1997. 109, 332, 337
- [WS96] M. A. Wiering and J. Schmidhuber. Solving POMDPs with Levin search and EIRA. In *Machine Learning: Proceedings of 13th International Conference*, pages 534–542, Bari, Italy, 1996. 124, 126
- [WS98] M. A. Wiering and J. Schmidhuber. Fast online “Q”(λ). *Machine Learning*, 33(1):105–116, 1998. 126
- [Yam98] K. Yamanishi. A decision-theoretic extension of stochastic complexity and its applications to learning. *IEEE Transactions on Information Theory*, 44:1424–1439, 1998. 342, 344

- [YFY01] R. Yaroshinsky and R. El-Yaniv. Smooth online learning of expert advice. *Submitted for publication*, 2001. 342
- [Zad65] L. A. Zadeh. Fuzzy sets. *Information and Control*, 8:338–353, 1965. 227
- [Zad78] L. A. Zadeh. Fuzzy sets as a basis for a theory of possibility. *Fuzzy sets and systems*, 1978, 1:3–28, 1978. 227
- [Zim91] H.-J. Zimmermann. *Fuzzy Set Theory—And Its Applications*. Kluwer, Dordrecht, The Netherlands, second revised edition, 1991. 227
- [ZL70] A. K. Zvonkin and L. A. Levin. The complexity of finite objects and the development of the concepts of information and randomness by means of the theory of algorithms. *Russian Mathematical Surveys*, 25(6):83–124, 1970. 217, 218, 219, 225, 226, 228, 304, 317, 710

Index

Symbols

α -norm

loss function, 324

σ -algebra, 216, 305

A

absolute

distance, 308

loss function, 324

absorbing

environment, 536

accessibility, 415

action, 402, 404

random, 531, 536

actions

concurrent, 805

active

agent, 804

system, 344

adaptive

control, 404, 507

Levin search, 229, 804

agent, 404

active, 804

algorithmic, 501

most intelligent, 507

prewired, 404

reactive, 404

universal, 501

agents, 402

bodiless, 807

embodied, 807

immortal, 530

lazy, 530

mortal, 807

aggregating strategy, 342

$AI\mu$ model

equivalence, 410

functional form, 402

recursive & iterative form, 408

special aspects, 412

$AI\rho$ model, 513

discounted, 519

functional, 513

iterative, 513

$AI\xi$ model

axiomatic approach, 531

loss bound, 622

Pareto-optimality, 514, 520

prediction, 622

structure, 531

AIXI model, 502

approximation, 811

computability, 811

general Bayes-mixture, 508

generality, 804

implementation, 811

optimality, 507

performance, 804

AIXI tl

optimality, 718

algorithm

best vote, 716

convergence, 804

fastest, 702, 707

hegde, 342

I/O stream, 712

incremental, 716

inversion, 704

learning, 804

non-incremental, 715

optimal, 804

optimization, 704

repeated evaluation, 712

search, 704

short, 711

simple, 704

speedup, 704

Strassen, V., 705

weighted majority, 342, 804

algorithmic

agent, 501

information theory, 225

probability, 228

specification, 702

algorithms

- parallel, 704

alphabet, 403

- continuous, 344
- countable, 344
- infinite, 344

amplitude

- signal, 330

Anderson, C. W., 127, 902

Angluin, D., 126, 228, 352, 901

animals, 808

application

- classification, 342
- games of chance, 328
- i.i.d. experiments, 342
- Kolmogorov complexity, 229
- Levin search, 229
- partial sequence prediction, 342

approximable, 209

- (semi)measure, 218
- probability distribution, 317

approximation

- AIXI model, 811
- value, valid, 717

arbitrary horizon

- convergence, 339

arithmetic, 810

artificial intelligence

- elegant \leftrightarrow complex, 811

Asimov, I., 501

asymmetry, 405

asymptotic

- convergence, 507
- learnability, 511
- optimality, 333
- runtime, 702

Atkeson, C. G., 127, 908

Auer, P., 343, 536, 901

autonomous

- robots, 807

average

- profit, 329
- reward, 529

axiom

- induction, 810

axiomatic approach

- AI ξ model, 531

axiomatic treatment, 809

axiomatize, 809

axioms, 707, 708

Aczel, J., 214, 228, 901

B

Babai et al., L., 713, 901

Bacchus, F., 226, 902

background

- knowledge, 806

balanced

- Pareto-optimality, 335, 515

Bandit problem, 508, 527

Barron, A. R., 110, 303, 341, 352, 901, 902

Bartlett, P. L., 804, 901

Barto, A. G., 102, 112, 126, 127, 402, 804, 902, 910

Baum, E. B., 804, 901

Baxter, J., 804, 901

Bayes rule, 203, 213

Bayes, T., 203, 901

Bayes-mixture

- general, 508

Bayes-optimal

- prediction, 317, 321

Bayesian

- self-optimizing policy, 527

behaviour

- innate, 808

belief

- contamination, 538
- probability, 306
- state, 404
- update, 228

Bellman equations, 404, 415

Bellman, R. E., 112, 402, 415, 801, 901

Bernardo's prior, 341

Bernoulli, 212

- process, 340

- sequence, 332

Bernoulli, J., 212, 901

Berry, D. A., 536, 902

Bertsekas, D. P., 102, 112, 126, 127, 402, 415, 416, 525, 527, 530, 539, 804, 902

bet, 328

bias, 504

Blackwell, D., 352, 901

Blum, M., 120, 122, 703, 711, 902

Blumer, A., 229, 901

Bohr, N., 801

boosting, 342

- bound, 355, 511, 531

bootstrap, 337

Borovkov, A. A., 336, 352, 902

Boulton, D. M., 229, 912

bound

- boost, 355, 511, 531
- entropy, 341
- error, 317, 318
- loss, 321–323, 328
- lower, 331
- probabilistic, 355
- sharp, 331
- tight, 331
- time, 703
- value, 509
- bounded
 - horizon, 338
- bounds
 - Kolmogorov complexity, 208
 - relative entropy, 308
 - value, 810
- brain
 - non-computable, 816, 817
- brain prosthesis
 - paradox, 814
- Buchanan, B. G., 227, 902
- butterfly effect, 321

C

- Cantor, G., 902
- Cardano, G., 902
- Carlin, J. B., 227, 904
- Carnap, R., 226, 902
- Cassandra, A. R., 127, 906
- Certainty factors, 227
- Cesa-Bianchi, N., 345, 536, 901, 902
- Cesa-Bianchi et al., N., 303, 342, 343, 804, 902
- chain rule, 105
- Chaitin, G. J., 104, 126, 208, 225, 226, 229, 317, 710, 817, 902, 903
- chaos, 414
- Cheeseman, P., 227, 228, 622, 903
- chess, 608, 611, 619
- Chinese room
 - paradox, 814
- chronological, 403
 - function, 405
 - order, 408
 - semimeasure, 504, 532, 537
 - Turing machine, 405
- Church thesis, 206
- Church, A., 226, 903
- Church-Turing thesis
 - extended, 206
- Churchland, P. S., 814, 903
- circumscription, 227
- Clarke, B. S., 110, 303, 341, 902
- class
 - problem, 702
- classical physics, 815
- classification, 342
- closed loop
 - control, 404
- code
 - universal, 207
- combining experts
 - prediction, 343
- complete
 - history, 415
- complexity
 - incomputable, 710
 - increase, 231
 - input sequence, 414
 - Kolmogorov, 104, 204
 - of functions, 710
 - of game, 611
 - parametric, 340
- compression
 - Lempel-Ziv, 317
- computability
 - AIXI model, 811
- computable
 - \leftrightarrow free will, 815
 - (semi)measure, 218
 - finite, 209
 - probability distribution, 216
 - recursive, 209
- computation
 - time, 707
- concept class
 - restricted, 508
- concepts
 - separability, 509
- concurrent
 - actions and observations, 805
- consciousness, 817
- consistency, 507
- consistent
 - control, 404
 - policy, 506
- constants, 414
- contamination
 - belief, 538
- Conte et al., M., 126, 229, 903
- continuity, 807
- continuous
 - alphabet, 344
 - entropy bound, 341
 - forecast, 343

- hypothesis class, 340
- probability class, 340, 526
- semimeasure, 216
- value, 521, 538
- weights, 340
- control, 404
 - adaptive, 404, 507
 - closed loop, 404
 - consistent, 404
 - open loop, 404
 - self-optimizing, 404
 - self-tuning, 404
 - stochastic, 404
- controlled
 - Markov chain, 404
- controller, 404
- convergence
 - \mathcal{M} , 307
 - ξ to μ , 310
 - ξ^{AI} to μ^{AI} , 505
 - algorithm, 804
 - arbitrary horizon, 339
 - asymptotic, 507
 - bounded horizon, 338
 - finite, 507
 - generalized, 312
 - in mean sum, 307
 - in probability, 307
 - in the mean, 307
 - individual, 354
 - Martin-Löf, 230, 307, 312, 354
 - of averages, 516, 537
 - of instantaneous loss, 326
 - random sequence, 307
 - rate, 517, 520
 - relations, 308
 - semi-martingale, 311
 - speed, 311, 355
 - unbounded horizon, 356
 - uniform, 512
 - value, 515, 517, 518, 520, 523
 - with high probability, 353
 - with probability 1, 307
- convexity
 - value, 513, 519
- Coppersmith, D., 705, 903
- cost
 - expected, 404
 - immediate, 404
 - total, 404
- countable
 - alphabet, 344

- probability class, 316
- counting, 710
- Cover, T. M., 226, 352, 903
- Cox's axioms, 214, 228
 - variants, 228
- Cox's theorem, 214, 228
 - loopholes, 228
- Cox, R. T., 105, 214, 226, 228, 903
- creator, 401
- cryptography, 807
 - RSA, 807
- Csiszr, I., 352, 903
- cumulative
 - reward, 404
- cumulatively enumerable
 - semi-measure, 317
- curvature, 534
 - Gauss, 535
- curvature matrix, 340
- cybernetic systems, 402
- cycle, 403
- cylinder sets, 216, 305

D

- Daley, R. P., 126, 226, 714, 903
- data
 - efficiency, 804
- Dauben, J. W., 903
- Dawid, A. P., 103, 203, 228, 342, 903
- Dayan, P., 127, 912
- dead code, 710
- decision
 - suboptimal, 511
 - wrong, 511
- decomposition, 412
- decryption, 807
- Default reasoning, 227
- degree of belief, 105, 211, 214
- delayed
 - prediction, 339
- Dempster, A. P., 227, 903
- Dempster-Shafer theory, 227
- density
 - error, 319
- deterministic, 403
 - \leftrightarrow free will, 815
 - environment, 405
 - optimal policy, 519
- dice
 - example, 330
- differential
 - gain, 529

- discounted
 - $AI\rho$ model, 519
 - value, 519
- discounting, 519, 525
 - finite, 528
 - general, 528
 - geometric, 528
 - harmonic, 528
 - power, 528
 - universal, 528
- discrete
 - (semi)measure, 221
 - probability class, 316
- distance
 - absolute, 308
 - Euclidian, 308
 - Hellinger, 308
 - Kullback-Leibler, 308
 - quadratic, 308
 - relative entropy, 308
 - square, 308
- distance measures
 - probability distribution, 308
- dominance
 - value, 416
- Doob, J. L., 305, 311, 521, 903
- Doyle, J., 227, 908
- Dubins, L., 352, 901
- dynamic
 - horizon, 528
- dynamic programming, 404

E

- economy based RL, 804
- effective
 - horizon, 526
- efficiency, 507
 - data, 804
 - time, 804
- Ehrenfeucht, A., 229, 901
- Einstein, A., 101, 701
- El-Yaniv, R., 343, 913
- embodied
 - agents, 807
- encrypted
 - information, 807
- entropy
 - bound, 341
 - inequalities, 308
 - inequality, 505, 533
 - relative, 308, 341
- enumerable, 209

- (semi)measure, 218
- chronological semimeasure, 537
- semi-measure, 317
- semimeasure, 533
- weights, 337
- enumeration
 - proof, 707
- environment, 404
 - absorbing, 536
 - deterministic, 405
 - ergodic, 524, 539
 - factorizable, 512
 - farsighted, 512
 - forgetful, 512, 539
 - general, 513
 - incomputable, 538
 - inductive, 510
 - influence, 344
 - known, 401
 - Markovian, 512
 - passive, 510
 - probabilistic, 406
 - pseudo-passive, 509, 510
 - random, 337
 - real, 807
 - relevant, 538
 - self-optimizing, 538, 622
 - stationary, 512
 - uniform, 512
 - weakly forgetful, 230
- environmental class, 622
 - limited, 508
 - others, 527
- Epicurus' principle, 203
- episode, 413
- equivalence
 - non-provable, 711
- equivalent
 - provably, 707
- ergodic
 - environment, 524, 539
 - MDP, 524, 538
- Erickson, G. W., 535, 903
- error
 - bound, 317, 318
 - density, 319
 - expected, 318
 - finite, 319
 - instantaneous, 318
 - loss function, 324
 - minimize, 317
 - posterization, 622

- probabilistic bound, 355
 - regret, 319
 - total, 318
 - error bound
 - exponential in Km , 607
 - lower, 331, 354
 - sharp, 331
 - tight, 331
 - estimable, 209
 - (semi)measure, 218
 - estimate
 - parameter, 341
 - Euclidian
 - distance, 308
 - loss function, 324
 - evaluation
 - of function, 613
 - event, 212
 - evolution, 817
 - example
 - dice, 330
 - expected
 - cost, 404
 - error, 318
 - loss, 322
 - utility, 415
 - expectimax
 - algorithm, 410
 - tree, 410
 - experiment, 105, 212
 - experiments
 - i.i.d., 342
 - expert advice
 - prediction, 530, 804
 - expert systems, 227
 - exploitation, 404, 416, 804
 - exploration, 416, 804
- F**
- factorizable
 - environment, 412, 512
 - fair coin flips, 105, 216
 - farsighted
 - environment, 512
 - farsightedness
 - dynamic, 407, 409
 - fast
 - matrix multiplication, 705
 - fastest
 - algorithm, 702
 - Feder, M., 124, 126, 226, 303, 304, 317, 318, 323, 342, 352, 811, 904, 908
 - feedback
 - more, 611
 - negative, 405
 - positive, 405
 - Feller, W., 227, 904
 - Ferguson, T. S., 336, 904
 - finance, 812
 - Fine, T. L., 227, 228, 904
 - Finetti, B., 226, 904
 - finite
 - computable, 209
 - convergence, 507
 - discounting, 528
 - error, 319
 - state space, 804
 - finite-state automata, 317
 - Fisher information, 340
 - Fisher, R. A., 226, 904
 - Fitting, M. C., 707, 904
 - fixed
 - horizon, 528
 - fixed point
 - prediction, 816
 - forecast
 - continuous, 343
 - probabilistic, 343
 - forgetful
 - environment, 230, 512, 539
 - formal
 - specification, 702
 - formula, 707
 - Fossa, J. A., 535, 903
 - free lunch, 337
 - free will
 - paradox, 815
 - frequentist, 211, 212
 - Freund, Y., 343, 536, 901, 904
 - Fristedt, B., 536, 902
 - Fudenberg, D., 608, 904
 - Fuentes, E., 126, 226, 909
 - function
 - complexity of, 710
 - concave, 323
 - evaluations, 613
 - loss, 321
 - minimize, 612
 - recursive, 532
 - function approximation
 - general, 804
 - linear, 804
 - function minimization, 612
 - greedy, 613

- inventive, 616
 - with AIXI, 617
- functional
 - AI ρ model, 513
- functional form, 408
- Furst, M., 701
- future
 - reward, 518
 - value, 518
- Fuzzy
 - logic, 227
 - sets, 227
 - systems, 227
- G**
- Gödel incompleteness, 814
- Gödel incompleteness, 817
- gain
 - differential, 529
- Gallager, R. G., 225, 904
- game playing
 - with AIXI, 611
- game theory, 608, 805
- games
 - chess, 619
 - complexity, 611
 - of chance, 328
 - repeated, 610
 - strategic, zero-sum, 607
 - variable lengths, 610
- Gauss
 - curvature, 535
- Gelman, A., 227, 904
- general
 - Bayes-mixture, 508
 - discounting, 528
 - environment, 513
 - loss bound, 328
 - loss function, 327
 - property, 807
 - weights, 513
- general Bayes-mixture
 - AIXI model, 508
- generalization, 804
- generalization techniques, 416
- generalized
 - convergence, 312
 - random sequence, 224
- generalized universal prior, 228, 503
- genetic algorithms, 811
- Gentile, C., 343, 901
- geometric

- discounting, 528
- geometric discounting
 - self-optimizing, 537
- Ginsberg, M. L., 227, 904
- Gittins, J. C., 536, 904
- Gold Standard, 801
- greedy, 407
 - function minimization, 613
 - strategy, 344
- Grove, A., 226, 902
- Grünwald, P. D., 229, 352, 535, 904
- Gutman, M., 124, 126, 226, 317, 318, 811, 904
- Gödel, K., 225, 904
- Gács, P., 104, 203, 225, 317, 344, 904, 905

H

- Hacking, I., 226, 905
- Hald, A., 226, 905
- Halpern, J. Y., 226, 228, 535, 902, 904, 905
- halting
 - probability, 230
 - sequence, 230
- harmonic
 - discounting, 528
- Hartmanis, J., 712, 905
- Haussler, D., 229, 342, 343, 901, 905
- HeavenHell example, 509, 536
- Heckerman, D. E., 228, 905
- hedge algorithm, 342
- Hellinger
 - distance, 308
 - loss function, 324
- history, 406
 - complete, 415
- holographic proofs, 713
- Hopcroft, J. E., 126, 205, 230, 905
- horizon, 409
 - arbitrary, 339
 - bounded, 338
 - choice, 414
 - dynamic, 528
 - effective, 526
 - fixed, 528
 - infinite, 529
 - problem, 518, 528
 - unbounded, 356
- Horvitz, E. J., 228, 905
- Hughes, R. I. G., 227, 905
- human, 414
- humans, 808
- Hume, D., 203, 905

Hutter, M., 4, 127, 309, 319, 338, 344, 352, 356,
527, 705, 804, 905, 906, 911

hypothesis class
continuous, 340

I

i.i.d., 524

process, 340

I/O sequence, 405

identification
system, 404

ignorance, 227

image, 414

immediate
cost, 404

immortal
agents, 530

implementation
AIXI model, 811

imprecise probabilities, 227

impreciseness, 227

incompleteness theorem, 711

incomputable
complexity, 710
environment, 538

inconsistent
policy, 506, 520

increase
complexity, 231

incremental
algorithm, 716

independent
episodes, 413

indifference principle, 215

individual
convergence, 354

induction, 303
axiom, 810
principle, 810
universal, 228

inductive
environment, 510

inequality
distance measures, 308
entropy, 505, 533
Jensen, 311, 323

inference
rules, 707

infinite
action space, 813
alphabet, 344
horizon, 529

observation space, 813
prediction space, 344

influence
environment, 344

information, 330
encrypted, 807
Fisher, 340
state, 404
symmetry, 208
transmission, 330

information theory
algorithmic, 225

informed
prediction, 318, 322

input, 402, 404
device, 414
regular, 405
reward, 405
word, 403

input space
choice, 414, 805

instantaneous
error, 318
loss, 322, 326
loss bound, 326
reward, 404

intelligence, 402
aspects, 621
effective order, 717
intermediate, 507
order relation, 506

interference, 227

intermediate
intelligence, 507

internal
reward, 808

inventive
function minimization, 616

inversion
problem, 704

investment, 328

iterative
AI ρ model, 513
iterative formulation, 408

J

Jaynes, E. T., 226–228, 906

Jeffrey, R. C., 226, 906

Jeffreys' prior, 341

Jensen's inequality, 311, 323

Jones, D. M., 536, 904

K

Kaelbling, L. P., 112, 127, 415, 906
 Kearns, M. J., 127, 527, 907
 Keynes, J. M., 226, 906
 Kivinen, J., 342, 343, 905, 907
 Kleene, S., 225, 906
 knowledge
 background, 806
 incorporate, 811
 physical, 807
 prior, 806
 universal prior, 806
 Knuth, D. E., 4, 906
 Ko, K.-I., 126, 226, 714, 906
 Koller, D., 226, 902
 Kolmogorov axioms, 212
 Kolmogorov complexity, 104, 204, 207, 208
 algorithmic properties, 211
 application, 229
 bounds, 208
 definition, 208
 information properties, 208
 oracle properties, 230
 time-bounded, 226
 time-limited, 714
 variants, 225
 Kolmogorov, A. N., 104, 126, 208, 225, 226, 229, 317, 710, 906, 907
 Kolmogorov-Uspenskii
 machine, 710
 Kraft, L. G., 225, 709, 907
 Kullback-Leibler
 divergence, 308
 Kumar, P. R., 115, 127, 204, 507, 509, 516, 525, 527, 536, 907
 Kwee, I., 127, 804, 906
 Kyburg, H. E., 226, 907

L

Lagrange
 multiplier, 534
 Lagrange multiplier, 351
 Lagrange, J., 401
 Lambalgen, M., 225, 907
 Langlotz, C. P., 228, 905
 Laplace, P., 301, 401, 907
 lazy
 agents, 530
 learnable
 asymptotically, 511
 task, 506

 Turing machine, 355
 learning, 307, 344
 a relation, 619
 algorithm, 804
 by reinforcement, 404, 415, 804
 by temporal difference, 804
 model, 404
 rate, 416
 supervised, 619
 with expert advice, 342
 Lempel-Ziv compression, 317
 Levin search, 704
 adaptive, 229, 804
 application, 229
 Levin, L. A., 104, 120, 124, 126, 217–219, 224–226, 228, 229, 304, 317, 704, 710, 907, 908, 913
 Li, M., 13, 102, 122–124, 126, 202, 204, 217, 221–223, 225, 226, 228–231, 303–305, 309, 311, 312, 317, 328, 337, 344, 352, 504, 532, 537, 703–705, 709, 710, 714, 811, 908, 912
 lifetime, 405, 414
 limited
 environmental class, 508
 limits, 414
 linear
 function approximation, 804
 linearity
 value, 513, 519
 Littlestone, N., 126, 229, 342, 343, 714, 908
 Littman, M. L., 112, 127, 415, 906
 locality, 807
 logarithmic
 loss function, 324
 logic
 planners, 804
 system, 707
 lookahead
 multi-step, 338
 lookup-table
 paradox, 814
 loss
 arbitrary, 322
 bound, 321–323, 328
 expected, 322
 function, 321
 instantaneous, 322, 326
 minimize, 322
 total, 322
 loss bound
 AI ξ model, 622

- general, 328
 - instantaneous, 326
 - structure, 343
 - with high probability, 353
- loss function
 - α -norm, 324
 - absolute, 324
 - arbitrary, 322
 - error, 324
 - Euclidian, 324
 - examples, 324
 - general, 327
 - Hellinger, 324
 - logarithmic, 324
 - quadratic, 324
 - square, 324
 - static, 327
 - time-dependent, 327
- Loveland, D. W., 226, 907
- lower
 - error bound, 354
- Lucas, J. R., 713, 814, 817, 907
- Lugosi, G., 345, 902
- M**
- machine
 - Kolmogorov-Uspenskii, 710
 - pointer, 710
- machine learning, 809
 - application, 809
 - categorization, 809
 - framework, 808
 - non-standard, 809
 - theory, 809
- machine model, 710
- Macready, W. G., 109, 333, 337, 338, 912
- majorization
 - multiplicative, 216
- manipulation, 808
- market based RL, 804
- Markov, 415
 - decision process, 524
- Markov chain
 - controlled, 404
- Markov decision process, 404
- Markovian
 - k -th order, 512
 - environment, 512
- Martin-Löf, 212, 229
 - convergence, 230, 312, 354
 - random sequence, 224, 230
- Martin-Löf, P., 212, 226, 229, 908
- martingales, 311, 352
- mathematical
 - specification, 702
- matrix multiplication
 - fast, 705
- maximal
 - profit, 329
- maximize
 - reward, 405
- maximum entropy principle, 215
- Maximum-Likelihood
 - prediction, 344
- McCallum, A. K., 804, 908
- McCarthy, J., 227, 908
- McDermott, D., 227, 908
- MDP, 404, 524
 - ergodic, 524, 538
 - stationary, 524
 - uniform mixture, 539
- measure
 - approximable, 218
 - computable, 218
 - discrete, 221
 - enumerable, 218
 - estimable, 218
 - universal, 219
- Merhav, N., 124, 126, 226, 303, 304, 317, 318, 323, 342, 352, 811, 904, 908
- Michie, D., 126, 410, 908
- minimize
 - error, 317
 - function, 612
 - loss, 322
- minimum description length, 810
- minimum message length, 810
- Mises, R., 226, 908
- mixing rate, 527
- mixture
 - uniform MDP, 539
- mixtures, 353
- model
 - AIXI, 502
 - causal, 203
 - learning, 203, 404
 - predictive, 203
 - true, 203
 - universal, 502
- modus ponens, 708
- money, 328
- monitor, 414
- Monte-Carlo, 527
- Moore, A. W., 112, 127, 415, 906, 908

Moravec, H., 814, 909
 Morgenstern, O., 127, 402, 608, 909
 mortal
 agents, 807
 most intelligent
 agent, 507
 Mosteller, F., 535, 909
 Motwani, R., 126, 205, 230, 905
 Moullagaliev, A., 336, 352, 902
 multi-agent system, 813
 multi-step
 lookahead, 338
 prediction, 338
 multiplicative
 domination, 306
 majorization, 216
 multiplier
 Lagrange, 534

N
 nanobots, 227
 Napoleon, B., 401
 natural
 Turing machine, 206
 natural numbers, 810
 Neumann, J. V., 127, 402, 608, 701, 909
 neural net, 229
 no free lunch, 337
 noise, 414
 noisy world, 414
 non-computable
 brain, 816, 817
 physics, 817
 non-deterministic world, 414
 non-perfect, 414
 non-provable
 equivalence, 711
 nonmonotonic logic, 227
 nonprobabilistic approaches, 228
 Norvig, P., 110, 112, 126, 227, 401, 410, 415, 608, 804, 910
 number of wisdom, 817

O

object, 807
 objectivist, 105, 211, 212
 observation, 404
 observations
 concurrent, 805
 Occam's razor, 201, 203, 215, 337, 809
 Ockham, W., 201

Odifreddi, P., 225, 909
 OnlyOne example, 510
 open loop
 control, 404
 optimal
 algorithm, 804
 deterministic policy, 519
 policy, 415
 problem solver, 229
 value, 513
 weights, 337
 optimality
 AIXI model, 507
 AIXItl, 718
 asymptotic, 333
 by construction, 508
 Pareto, 333
 universal, 506, 507
 optimality properties, 331
 optimization
 problem, 612, 704
 oracle properties
 Kolmogorov complexity, 230
 order relation
 effective intelligence, 717
 intelligence, 506
 universal, 507
 Osborne, M. J., 416, 608, 805, 909
 outcome, 212
 output, 402, 404
 device, 414
 word, 403
 output space
 choice, 414, 805

P

paradox
 brain prosthesis, 814
 Chinese room, 814
 free will, 815
 lookup-table, 814
 parallel
 algorithms, 704
 parameter
 estimate, 341
 parametric
 complexity, 340
 Pareto-optimality, 333, 353, 507, 537
 AI ξ model, 514, 520
 balanced, 335, 515
 Paris, J. B., 214, 228, 909
 Pascal, B., 909

- passive
 - environment, 510
- Peano axioms, 810
- Penrose, R., 126, 713, 814, 817, 909
- perception, 402, 404
- perfect, 414
- performance
 - AIXI model, 804
 - sequence prediction, 804
- perspective
 - subjective, 807
- physical
 - knowledge, 807
- physical random processes, 105, 212
- physics
 - classical, 815
 - non-computable, 817
 - quantum, 414, 815
 - wave function collapse, 817
- Pinsker, M. S., 352, 909
- Pintado, X., 126, 226, 909
- planners
 - logic, 804
- plausibility, 214
- player, rational, 607
- pointer
 - machine, 710
- policy, 403, 404, 415
 - consistent, 506
 - extended chronological, 716
 - inconsistent, 506, 520
 - optimal, 415
 - optimal deterministic, 519
 - probabilistic, 416, 518
 - restricted class, 527
 - self-optimizing, 508, 517, 520, 524
 - stationary, 524
- policy iteration, 404, 415, 804
- Popper, K. R., 226, 909
- portfolio, 328
- Possibility theory, 227
- Post, E. L., 225, 909
- posterior probability, 306
- posterization, 355, 511, 531, 806
 - prediction error, 622
- power
 - discounting, 528
- prediction
 - AI ξ model, 622
 - arbitrary, 322
 - Bayes-optimal, 317, 321
 - combining experts, 343
 - delayed, 339
 - expert advice, 530, 804
 - fixed point, 816
 - informed, 318, 322
 - Maximum-Likelihood, 344
 - multi-step, 338
 - self-contradictory, 816
 - universal, 318
 - with expert advice, 342
- prediction error
 - posterization, 622
- prediction space
 - infinite, 344
- prefix property, 403
- prequential approach, 103, 203, 228
- prewired
 - agent, 404
- principle
 - of indifference, 215
 - of maximum entropy, 215
 - of parsimony, 215
 - of simplicity, 215
 - of symmetry, 215
- prior
 - Bernando, 341
 - determination, 215
 - Jeffreys, 341
 - knowledge, 806
 - probability, 305
 - Solomonoff, 337
 - Speed, 317, 353
 - universal knowledge, 806
- probabilistic
 - environment, 406
 - forecast, 343
 - policy, 416, 518
- probability, 211
 - algorithmic, 228
 - axioms, 212
 - belief, 306
 - classes, 316
 - complex valued, 227
 - conditional, 213
 - distribution, 213, 408
 - frequency, 212
 - halting, 230
 - mass function, 213
 - measure, 213, 216, 305
 - nearby distribution, 316
 - objective, 212, 815
 - prior, 215, 305
 - second order, 227

- subjective, 214, 815
- probability class
 - continuous, 340, 526
 - countable, 316
 - discrete, 316
- probability distribution, 213, 408
 - approximable, 317
 - computable, 216, 316
 - conditional, 306, 408
 - distance measures, 308
 - generating, 305
 - generic class, 317
 - nearby, 316
 - over values, 813
 - posterior, 306
 - simple, 317
 - Solomonoff, 317
 - true, 305
 - universal, 305, 317
 - unknown, 305
- probability theory
 - alternatives, 227
 - history, 226
- problem
 - class, 702
 - horizon, 518, 528
 - inversion, 704
 - optimal solver, 229
 - optimization, 612, 704
 - poorly specified, 707
 - relevant, 511
 - satisfiability, 706
 - solvable, 506
 - traveling salesman, 612, 618
- process
 - Bernoulli, 340
 - i.i.d., 340
- profit, 328
 - average, 329
 - maximal, 329
- program
 - extended chronological, 716
 - fastest, 702
 - search, 705
- proof, 703, 811
 - enumeration, 707
 - search, 705
 - theory, 707
- property
 - general, 807
- provably
 - equivalent, 707

- pseudo-passive
 - environment, 509, 510
 - pseudo-random, 815
 - Putnam, H., 226, 909

Q

- quadratic
 - distance, 308
 - loss function, 324
- quantum logic, 227
- quantum physics, 414, 815

R

- Ramsey, F. P., 226, 909
- random
 - action, 531, 536
 - environment, 337
 - pseudo, 815
 - sequence, 304
- random sequence
 - convergence, 307
 - convergence relations, 308
 - generalized, 224
 - individual, 230
 - Martin-Löf, 224, 230
- rate
 - convergence, 517, 520
- rational player, 607
- reactive
 - agent, 404
- real
 - environment, 807
- recursive
 - computable, 209
 - function, 532
- recursive formulation, 408
- reduction
 - state space, 415
- redundance, 604
- redundant information, 604
- reflex, 808
- regression, 203
- regret
 - error, 319
- regularization, 318
- Reichenbach, H., 226, 228, 909
- reinforcement learning, 404, 415
 - classical algorithms, 804
 - economy based, 804
 - marked based, 804
 - policy gradient, 804

Reiter, R., 227, 909
 relation
 value difference, 516, 520
 relative
 entropy, 308, 341
 relative entropy
 distance, 308
 relevant
 environment, 538
 problem, 511
 Rescher, N., 535, 909
 restricted
 concept class, 508
 restricted domains, 416
 reward, 403, 415
 average, 529
 cumulative, 404
 future, 405, 518
 instantaneous, 404
 internal, 808
 maximize, 405
 mean, median, quantile, 813
 total, 404, 405
 Ring, M., 804, 909
 Rissanen, J. J., 124, 126, 209, 229, 341, 342, 811, 909, 910
 Robbins, H., 535, 910
 Robinson, A., 713, 910
 robots
 autonomous, 807
 the 3 laws, 501
 robustness, 813
 Rogers, H., 225, 910
 RSA
 cryptography, 807
 Rubenstein, A., 416, 608, 805, 909
 Rubin, D. B., 227, 904
 rules
 inference, 707
 runtime
 asymptotic, 702
 Russell, S. J., 110, 112, 126, 227, 401, 410, 415, 608, 804, 910

S

sample space, 212, 216, 305
 Samuel, A. L., 127, 910
 Savage, L. J., 226, 910
 scaling
 AIXI down, 810
 Schapire, R. E., 343, 536, 901, 904

Schmidhuber, J., 107, 124, 126, 127, 221, 225, 226, 228, 229, 317, 337, 338, 353, 706, 711, 713, 714, 804, 812, 815, 906, 910–912
 Schmidt, M., 126, 229, 910
 Schnorr, C. P., 225, 226, 910
 Schönhage, A., 710, 910
 science, 810
 search
 adaptive Levin, 804
 algorithm, 704
 Levin, 704
 program, 705
 proof, 705
 time bound, 705
 Searle, J., 814, 910
 Seiferas, J. I., 710, 911
 self-optimizing, 204
 Bayesian policy, 527
 control, 404
 environment, 538, 622
 geometric discounting, 537
 policy, 508, 517, 520, 524
 self-optimizingness, 507
 self-tuning, 204
 control, 404
 self-tuningness, 507
 Semenov, A. L., 229, 912
 semi-computable, 209
 semi-martingale
 convergence, 311
 semi-measure
 cumulatively enumerable, 317
 enumerable, 317
 Solomonoff, 317
 semimeasure
 approximable, 218
 chronological, 504, 532, 537
 computable, 218
 continuous, 216
 discrete, 221
 enumerable, 218, 505, 533, 537
 estimable, 218
 universal, 216, 219
 universal, time-limited, 714
 semimeasures, 352
 separability
 concepts, 509
 sequence
 Bernoulli, 332
 halting, 230
 random, 304

- training, 812
- sequence prediction
 - generality, 804
 - performance, 804
 - Solomonoff, 217
 - universal, 304
 - with $AI\mu$, 604
 - with AIXI, 603, 606
- sequences, 403
- sequential decision theory, 404, 415
- set
 - fuzzy, 227
 - prefix free, 403
- Shafer, G., 227, 535, 911
- Shannon, C. E., 226, 709, 911
- Shen, A. K., 229, 912
- Shoenfield, J. R., 225, 707, 911
- short
 - algorithm, 711
- Shortliffe, E. H., 227, 902, 911
- side information, 208
- signal
 - amplitude, 330
- simple
 - algorithm, 704
- simplex, 351
- Singh, S., 127, 527, 907
- Smith, C. H., 126, 228, 352, 901
- Solomonoff, 203
 - prior, 337
 - semi-measure, 317
- Solomonoff, R. J., 105, 106, 120, 126, 203, 208, 216–218, 222, 225, 226, 228, 303, 304, 309, 311, 352, 504, 615, 705, 710, 801, 812, 911
- solution
 - verification, 704
- solvable
 - problem, 506
- specification
 - algorithmic, 702
 - formal, 702
 - mathematical, 702
- speed
 - convergence, 355
- Speed prior, 229, 317, 353
- speedup
 - algorithm, 704
 - theorem, 702, 710
- split trees, 804
- square
 - distance, 308
 - loss function, 324
- state, 415
 - belief, 404
 - environmental, 415
 - information, 404
 - internal, 402
- state space, 524
 - finite, 804
- static
 - loss function, 327
- stationarity, 415
- stationary
 - environment, 512
 - MDP, 524
 - policy, 524
- statistics, 305
- Stern, H. S., 227, 904
- stochastic
 - control, 404
- stock market, 328
- Stork, D., 338, 911
- stragegy
 - aggregating, 342
- Strassen, V., 705, 911
 - algorithm, 705
- strategy
 - die selection, 330
 - games, 607
 - greedy, 344
 - winning, 328
- string
 - empty, 304, 405
 - length, 403
- strings, 403
- structural risk minimization, 810
- structure
 - $AI\xi$ model, 531
 - loss bound, 343
- subjective
 - perspective, 807
- subjectivist, 105, 211, 214
- suboptimal
 - decision, 511
- supervised learning
 - from examples, 619
 - with AIXI, 619
- support vector machines, 810
- Sutton, R., 601
- Sutton, R. S., 102, 112, 126, 127, 402, 804, 902, 910, 911
- symmetry
 - information, 208

symmetry principle, 215

system, 404

active, 344

identification, 404

logic, 707

Székely, G. J., 535, 911

T

tape

bidirectional, 206, 216, 403

unidirectional, 206, 216, 403

task

learnable, 506

temporal difference learning, 404, 804

term, 707

Tesauro, G., 127, 911

theorem prover

elaborate, 713

theorem provers, 811

theory, *see* particular theory 805

proof, 707

Thomas, J. A., 226, 352, 903

threshold

winning, 330

time

bound, 703

computation, 707

efficiency, 804

to win, 328, 329

time bound

AIXI $_{tl}$, 718

search, 705

time-bounded

Kolmogorov complexity, 226

time-dependent

loss function, 327

Tirole, J., 608, 904

total

cost, 404

error, 318

loss, 322

reward, 404

trader, 328

training

sequence, 812

transductive inference, 203

transition matrix, 404, 524

transmission

information, 330

transparent proofs, 713

traveling salesman

problem, 612, 618

tree

split, 804

Tribus, M., 214, 228, 911

Tromp, J., 203, 344, 905

Tsitsiklis, J. N., 102, 112, 127, 402, 415, 525, 527, 804, 902

Turing

test, 814

thesis, 206

Turing machine, 403

chronological, 405

head, 403

learnable, 355

natural, 206

prefix, 206, 216

universal, 207, 216

unnatural, 229

Turing, A. M., 207, 225, 401, 814, 911

typing monkeys, 714

U

Ullman, J. D., 126, 205, 230, 905

unbiasedness, 507

unbounded horizon

convergence, 356

uncertainty, 211

sources, 211

underline, 408

uniform

convergence, 512

environment, 512

MDP mixture, 539

universal

(semi)measure, 219

agent, 501

AIXI, 504

AIXI model, 502

code, 207

discounting, 528

element, 505

generalized prior, 503

induction, 228

optimality, 506, 507

order relation, 507

prediction, 318

prior knowledge, 806

probability distribution, 305, 317

time-limited semimeasure, 714

universality property, 504

universe, 401, 415

update weights, 343

Uspenskii, V. A., 229, 710, 907, 912

utility, 405
 expected, 415

V

V'yugin, V. V., 228, 908
 vagueness, 227
 Valiant, L. G., 126, 229, 912
 valid approximation
 value, 717
 value
 bound, 509
 bounds, 810
 continuous, 521, 538
 convergence, 515, 517, 518, 520, 523
 convexity, 513, 519
 difference relation, 516, 520
 discounted, 519
 dominance, 416
 function, 513
 future, 518
 justification, 719
 linearity, 513, 519
 median, 813
 optimality results, 513
 quantile, 813
 valid approximation, 717
 value iteration, 404, 415, 804
 Vapnik, V. N., 203, 912
 Varaiya, P. P., 115, 127, 204, 507, 509, 516, 525,
 527, 536, 907
 Vereshchagin, N., 203, 912
 verification
 of solution, 704
 video camera, 414
 Vitányi, P. M. B., 102, 122–124, 126, 202–204,
 217, 221–223, 225, 226, 228–231, 303–
 305, 309, 311, 312, 317, 328, 337, 344,
 352, 504, 532, 537, 703–705, 709, 710,
 714, 811, 905, 908, 911, 912
 Vogel, W., 536, 912
 Voronkov, A., 713, 910
 vote
 best, 716
 democratic, 714
 Vovk, V. G., 126, 229, 312, 342, 343, 352, 355,
 714, 912

W

Wald, A., 226, 912
 Wallace, C. S., 229, 912
 Walley, P., 227, 912

Wang, Y., 225, 912
 Warmuth, M. K., 126, 229, 342, 343, 714, 901,
 905, 907, 908
 Watkins, C., 126, 127, 229, 912
 wave function collapse, 817
 weather forecasts, 321
 weighted majority algorithm, 342, 804
 weights, 306
 continuous, 340
 enumerable, 337
 general, 513
 optimal, 336, 337
 posterior, 306
 time dependent, 306
 universal choice, 354
 update rule, 343
 Wheeler, J., 701
 Wiering, M. A., 124, 126, 127, 229, 706, 713,
 804, 911, 912
 winning
 strategy, 328
 threshold, 330
 Winograd, S., 705, 903
 witness, 704
 Wolpert, D. H., 109, 333, 337, 338, 912
 worst-case
 reasoning, 227

Y

Yamanishi, K., 343, 345, 912
 Yaroshinsky, R., 343, 913

Z

Zadeh, L. A., 227, 913
 Zhao, J., 124, 126, 229, 706, 713, 804, 911
 Zimmermann, H.-J., 227, 913
 Zvonkin, A. K., 217–219, 225, 226, 228, 304, 317,
 710, 913

Decision theory formally solves the problem of rational agents in uncertain worlds if the true environmental prior probability distribution is known. Solomonoff's theory of universal induction formally solves the problem of sequence prediction for unknown prior distribution. In this thesis both ideas are unified to one parameter-free theory for universal Artificial Intelligence. We give strong arguments that the resulting AIXI model is the most intelligent unbiased agent possible. We outline for a number of problem classes, including sequence prediction, strategic games, function minimization, reinforcement and supervised learning, how the AIXI model can formally solve them. The major drawback of the AIXI model is that it is uncomputable. To overcome this problem, we construct a modified algorithm $AIXI_{tl}$, which is still effectively more intelligent than any other time t and length l bounded agent. The computation time of $AIXI_{tl}$ is of the order $t \cdot 2^l$. The discussion includes formal definitions of intelligence order relations, the horizon problem and relations of the AIXI theory to other AI approaches.

$$\begin{array}{rcl}
 \text{Decision Theory} & = & \text{Probability} + \text{Utility Theory} \\
 + & & + \\
 \text{Universal Induction} & = & \text{Occam} + \text{Epicurus} + \text{Bayes} \\
 || & & || \\
 \text{Universal Artificial Intelligence without Parameters} & &
 \end{array}$$

Marcus Hutter received his masters in computer sciences in 1992 at TU-Munich, Germany. After his PhD in theoretical particle physics he developed algorithms in a medical software company for 5 years. For three years he has been working as a researcher at the AI institute IDSIA in Lugano, Switzerland. His current interests are centered around reinforcement learning, algorithmic information theory and statistics, universal induction schemes, adaptive control theory, and related areas.

