

SELF-OPTIMIZING AND PARETO-OPTIMAL POLICIES IN GENERAL ENVIRONMENTS BASED ON BAYES-MIXTURES

Marcus Hutter

IDSIA, Galleria 2, CH-6928 Manno-Lugano, Switzerland*

marcus@idsia.ch – <http://www.idsia.ch/~marcus>

Keywords

Rational agents, sequential decision theory, reinforcement learning, value function, Bayes mixtures, self-optimizing policies, Pareto-optimality, unbounded effective horizon, (non) Markov decision processes.

Abstract

The problem of making sequential decisions in unknown probabilistic environments is studied. In cycle t action y_t results in perception x_t and reward r_t , where all quantities in general may depend on the complete history. The perception x_t and reward r_t are sampled from the (reactive) environmental probability distribution μ . This very general setting includes, but is not limited to, (partial observable, k-th order) Markov decision processes. Sequential decision theory tells us how to act in order to maximize the total expected reward, called value, if μ is known. Reinforcement learning is usually used if μ is unknown. In the Bayesian approach one defines a mixture distribution ξ as 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} . As an example application, we discuss ergodic Markov decision processes, which allow for self-optimizing policies. Furthermore, 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.

*This work was supported by SNF grant 2000-61847.00 to Jürgen Schmidhuber.

1 Introduction

Reinforcement learning: There exists a well developed theory for reinforcement learning agents in known probabilistic environments (like Blackjack) called sequential decision theory [Bel57, Ber95]. The optimal agent is the one which maximizes the future expected reward sum. This setup also includes deterministic environments (like static mazes). Even adversarial environments (like Chess or Backgammon) may be seen as special cases in some sense [Hut00, ch.6] (the reverse is also true [BT00]). Sequential decision theory deals with a wide range of problems, and provides a general formal solution in the sense that it is mathematically rigorous and (uniquely) specifies the optimal solution (leaving aside computational issues). The theory breaks down when the environment is unknown (like when driving a car in the real world). Reinforcement learning algorithms exist for unknown Markov decision processes (MDPs) with small state space, and for other restricted classes [KLM96, SB98, Ber95, KV86], but even in these cases their learning rate is usually far from optimum.

Performance measures: In this work we are interested in general (probabilistic) environmental classes \mathcal{M} . We assume \mathcal{M} is given, and that the true environment μ is in \mathcal{M} , but is otherwise unknown. The expected reward sum (value) V_μ^p when following policy p is of central interest. We are interested in policies \tilde{p} which perform well (have high value) independent of what the true environment $\mu \in \mathcal{M}$ is. A natural demand from an optimal policy is that there is no other policy yielding higher or equal value in *all* environments $\nu \in \mathcal{M}$ and a strictly higher value in one $\nu \in \mathcal{M}$. We call such a property *Pareto-optimality*. The other quantity of interest is how close $V_\mu^{\tilde{p}}$ is to the value V_μ^* of the optimal (but infeasible) policy p^μ which knows μ in advance. We call a policy whose average value converges asymptotically for all $\mu \in \mathcal{M}$ to the optimal value V_μ^* if μ is the true environment, *self-optimizing*.

Main new results for Bayes-mixtures: We define the Bayes-mixture ξ as a weighted average of the environments $\nu \in \mathcal{M}$ and analyze the properties of the Bayes-optimal policy p^ξ which maximizes the mixture value V_ξ . One can show that not all environmental classes \mathcal{M} admit self-optimizing policies. One way to proceed is to search for and prove weaker properties than self-optimizingness [Hut00]. Here we follow a different approach: 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 \tilde{p} at all. The main new result of this work is that this necessary condition is also sufficient for p^ξ to be self-optimizing. No other properties need to be imposed on \mathcal{M} . The other new result is that p^ξ is always Pareto-optimal, with no conditions at all imposed on \mathcal{M} .

Contents: Section 2 defines the model of agents acting in general probabilistic environments and defines the finite horizon value of a policy and the optimal value-maximizing policy. Furthermore, the mixture-distribution is introduced and the fundamental linearity and convexity properties of the mixture-values is stated. Section 3 defines and proves Pareto-optimality of p^ξ . The concept is refined to balanced Pareto-optimality, showing

that a small increase of the value for some environments only leaves room for a small decrease in others. Section 4 shows that p^ξ is self-optimizing if \mathcal{M} admits self-optimizing policies, and also gives the speed of convergence in the case of finite \mathcal{M} . The finite horizon model has several disadvantages. For this reason Section 5 defines the discounted (infinite horizon) future value function, and the corresponding optimal value-maximizing policy. Pareto-optimality and self-optimizingness of p^ξ are shown for this model. As an application we show in Section 6 that the class of ergodic MDPs admits self-optimizing policies w.r.t. the undiscounted model and w.r.t. the discounted model if the effective horizon tends to infinity. Together with the results from the previous sections this shows that p^ξ is self-optimizing for ergodic MDPs. Conclusions and outlook can be found in Section 7.

2 Rational Agents in Probabilistic Environments

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 $y_{1:k} \equiv y_1 \dots y_k$ of the agents policy p on input $x_{<k} \equiv x_1 \dots x_{k-1}$ 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*. 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 don't talk about states of the environment, only about observations. In the more general case of a *probabilistic environment*, given the history $\underline{y}_{<k} y_k \equiv \underline{y}_1 \dots \underline{y}_{k-1} y_k \equiv y_1 x_1 \dots y_{k-1} x_{k-1} y_k$, the probability that the environment leads to perception x_k in cycle k is (by definition) $\rho(\underline{y}_{<k} \underline{y}_k)$. The underlined argument \underline{y}_k in ρ is a random variable and the other non-underlined arguments $\underline{y}_{<k} y_k$ represent conditions.¹ We call probability distributions like ρ *chronological*. Since value optimizing policies can always be chosen deterministic, there is no real need to generalize the setting to probabilistic policies. Arbitrarily we formalize Sections 3 and 4 in terms of deterministic policies and Section 5 in terms of probabilistic policies.

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 $x_k \equiv x'_k r_k$ are divided into a regular part x'_k and some (possibly empty or delayed) reward $r_k \in [0, r_{max}]$.² We use the abbreviation

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

¹The standard notation $\rho(x_k | \underline{y}_{<k} y_k)$ for conditional probabilities destroys the chronological order and would become quite confusing in later expressions.

²In 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 ρ' depending on the complete history. It is very convenient to integrate r_k into x_k and ρ' into ρ .

which is essentially Bayes rules, and $\varepsilon = \mathbf{y}_{<1}$ for the empty string. The ρ -expected reward sum (value) of future cycles k to m with outputs $\mathbf{y}_{k:m}$ generated by the agent's policy p , the optimal policy p^ρ which maximizes the value, its action y_k and the corresponding value can formally be defined as follows.

Definition 1 (Value function and optimal policy) *We define the value of policy p in environment ρ given history $\mathbf{y}_{<k}$, or shorter, the ρ -value of p given $\mathbf{y}_{<k}$, as*

$$V_{km}^{p\rho}(\mathbf{y}_{<k}) := \sum_{x_{k:m}} (r_k + \dots + r_m) \rho(\mathbf{y}_{<k} \mathbf{y}_{k:m}) |_{y_{1:m}=p(x_{<m})}. \quad (2)$$

m is the lifespan or initial horizon of the agent. The ρ -optimal policy p^ρ which maximizes the (total) value $V_\rho^p := V_{1m}^{p\rho}(\varepsilon)$ is

$$p^\rho := \arg \max_p V_\rho^p, \quad V_{km}^{*\rho}(\mathbf{y}_{<k}) := V_{km}^{p^\rho\rho}(\mathbf{y}_{<k}). \quad (3)$$

Explicit expressions for the action y_k in cycle k of the ρ -optimal policy p^ρ and their value $V_{km}^{\rho}(\mathbf{y}_{<k})$ are*

$$y_k = \arg \max_{y_k} \sum_{x_k} \max_{y_{k+1}} \sum_{x_{k+1}} \dots \max_{y_m} \sum_{x_m} (r_k + \dots + r_m) \cdot \rho(\mathbf{y}_{<k} \mathbf{y}_{k:m}), \quad (4)$$

$$V_{km}^{*\rho}(\mathbf{y}_{<k}) = \max_{y_k} \sum_{x_k} \max_{y_{k+1}} \sum_{x_{k+1}} \dots \max_{y_m} \sum_{x_m} (r_k + \dots + r_m) \cdot \rho(\mathbf{y}_{<k} \mathbf{y}_{k:m}). \quad (5)$$

where $\mathbf{y}_{<k}$ is the actual history.

One can show [Hut00] that these definitions are consistent and correctly capture our intention. For instance, consider the expectimax expression (5): 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 $y_k x_k \dots y_m x_m$ to correctly incorporate the dependency of x_i and y_i on the history. Obviously

$$V_{km}^{*\rho}(\mathbf{y}_{<k}) \geq V_{km}^{p\rho}(\mathbf{y}_{<k}) \quad \forall p, \quad \text{especially} \quad V_\rho^* \geq V_\rho^p \quad \forall p. \quad (6)$$

Known environment μ : Let us now make a change in conventions and assume that μ is the true environment in which the agent operates and that we know μ (like in Blackjack).³ Then, policy p^μ is optimal in the sense that no other policy for an agent leads to higher μ -expected reward. This setting includes as special cases deterministic environments, Markov decision processes (MDPs), and even adversarial environments for special choices of μ [Hut00]. There is no principle problem in determining the optimal action y_k as long as μ is known and computable and \mathcal{X} , \mathcal{Y} and m are finite.

³If the existence of true objective probabilities violates the philosophical attitude of the reader he may assume a deterministic environment μ .

The mixture distribution ξ : Things drastically change if μ is unknown. For (parameterized) MDPs with small state (parameter) space, suboptimal reinforcement learning algorithms may be used to learn the unknown μ [KLM96, SB98, Ber95, KV86]. In the Bayesian approach the true probability distribution μ is not learned directly, but is replaced by a Bayes-mixture ξ . Let us assume that we know that the true environment μ is contained in some known set \mathcal{M} of environments. For convenience we assume that \mathcal{M} is finite or countable. The Bayes-mixture ξ is defined as

$$\xi(\underline{y}_{1:m}) = \sum_{\nu \in \mathcal{M}} w_\nu \nu(\underline{y}_{1:m}) \quad \text{with} \quad \sum_{\nu \in \mathcal{M}} w_\nu = 1, \quad w_\nu > 0 \quad \forall \nu \in \mathcal{M} \quad (7)$$

The weights w_ν may be interpreted as the prior degree of belief that the true environment is ν . Then $\xi(\underline{y}_{1:m})$ could be interpreted as the prior subjective belief probability in observing $x_{1:m}$, given actions $y_{1:m}$. It is, hence, natural to follow the policy p^ξ which maximizes V_ξ^p . If μ is the true environment the expected reward when following policy p^ξ will be $V_\mu^{p^\xi}$. The optimal (but infeasible) policy p^μ yields reward $V_\mu^{p^\mu} \equiv V_\mu^*$. 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^* . These are the main issues of the remainder of this work.

A universal choice of ξ and \mathcal{M} : One may also ask what the most general class \mathcal{M} and weights w_ν could be. Without any prior knowledge we should include *all* environments in \mathcal{M} . In this generality this approach leads at best to negative results. More useful is the assumption that the environment possesses some structure, we just don't know which. From a computational point of view we can only unravel effective structures which are describable by (semi)computable probability distributions. So we may include *all* (semi)computable (semi)distributions in \mathcal{M} . Occam's razor tells us to assign high prior belief to simple environments. Using Kolmogorov's universal complexity measure $K(\nu)$ for environments ν one should set $w_\nu \sim 2^{-K(\nu)}$, where $K(\nu)$ is the length of the shortest program on a universal Turing machine computing ν . The resulting policy p^ξ has been developed and intensively discussed in [Hut00]. It is a unification of sequential decision theory [Bel57, Ber95] and Solomonoff's celebrated universal induction scheme [Sol78, LV97]. In the following we consider generic \mathcal{M} and w_ν . The following property of V_ρ is crucial.

Theorem 1 (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 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 identity is definition (3), the equality uses linearity of $V_\rho^{p^\xi}$ just proven, and the last inequality follows from the dominance (6) and non-negativity of the weights w_ν . \square

One loose interpretation of the convexity is that a mixture can never increase performance. In the remainder of this work μ denotes the true environment, ρ any distribution, and ξ the Bayes-mixture of distributions $\nu \in \mathcal{M}$.

3 Pareto Optimality of policy p^ξ

The total μ -expected reward $V_\mu^{p^\xi}$ of policy p^ξ is of central interest in judging the performance of policy p^ξ . We know that there *are* policies (e.g. p^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 .

Theorem 2 (Pareto optimality) *Policy 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 .*

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 (Theorem 1). The strict inequality follows from the assumption and from $w_\nu > 0$. The identity is just Definition 1(3). The last inequality follows from the fact that p^ξ maximizes by definition the universal value (6). The contradiction $V_\xi^p > V_\xi^{p^\xi}$ proves Pareto-optimality of policy p^ξ . \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 3 (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 (Theorem 2) 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 (Theorem 1) and dominance $V_\xi^{p^\xi} \geq V_\xi^{\tilde{p}}$ (6). The remainder of Theorem 3 is obvious from $0 \leq \Delta = \Delta_{\mathcal{L}} - \Delta_{\mathcal{H}}$ and by bounding the weighted average Δ_η by its maximum. \square

4 Self-optimizing Policy p^ξ w.r.t. Average Value

In the following we study under which circumstances⁴

$$\frac{1}{m} V_{1m}^{p^\xi \nu} \rightarrow \frac{1}{m} V_{1m}^{*\nu} \quad \text{for } m \rightarrow \infty \quad \text{for all } \nu \in \mathcal{M}. \quad (8)$$

The least we must demand from \mathcal{M} to have a chance that (8) is true is that there exists some policy \tilde{p} at all with this property, i.e.

$$\exists \tilde{p} : \frac{1}{m} V_{1m}^{\tilde{p}\nu} \rightarrow \frac{1}{m} V_{1m}^{*\nu} \quad \text{for } m \rightarrow \infty \quad \text{for all } \nu \in \mathcal{M}. \quad (9)$$

Luckily, this necessary condition will also be sufficient. This is another (asymptotic) optimality property of policy p^ξ . If universal convergence in the sense of (9) is possible at all in a class of environments \mathcal{M} , then policy p^ξ converges in the sense of (8). We will call policies \tilde{p} with a property like (9) *self-optimizing* [KV86]. The following two Lemmas pave the way for proving the convergence Theorem.

Lemma 1 (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 $\delta_\nu(m) \geq 0$ converging to zero.

⁴Here 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.

Lemma 2 (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) $\delta_\nu(m) \leq f(m) \quad \forall \nu$ implies $\delta(m) \leq f(m)$.
- ii) $\delta_\nu(m) \xrightarrow{m \rightarrow \infty} 0 \quad \forall \nu$ implies $\delta(m) \xrightarrow{m \rightarrow \infty} 0$ if $0 \leq \delta_\nu(m) \leq c$.

Proof: (i) immediately follows from $\delta(m) = \sum_\nu w_\nu \delta_\nu(m) \leq \sum_\nu w_\nu f(m) \leq f(m)$. For (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 $\delta(m) \geq 0$. \square

Theorem 4 (Self-optimizing policy p^ξ w.r.t. average value) 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 (9) is possible at all in a class of environments \mathcal{M} , then policy p^ξ converges (in the sense of (8)). 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 p_m^ξ can be arbitrarily close to optimum for sufficiently large (pre-chosen) horizon m . This weakness will be resolved in the next subsection.

Proof: (i) $\Delta_\nu(m) = f(m)$ implies $\Delta(m) = f(m)$ by Lemma 2(i). Inserting this in Lemma 1 proves Theorem 4(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 assumed bounded rewards $0 \leq r \leq r_{max}$ 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 Theorem 4(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 1 and convergence to zero from Lemma 2(ii). This proves Theorem 4(ii). \square .

In Section 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 4).

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 main reason is that 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 fixed⁵. 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.

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 $\sum_{i=k}^{k+h_k^{eff}} \gamma_k \sim \sum_{i=k+h_k^{eff}}^{\infty} \gamma_k$. See [Hut00, Ch.4] 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 .

Definition 2 (Discounted value function and optimal policy) *We define the γ discounted weighted-average future value of (probabilistic) policy π in environment ρ given history $\underline{y}_{<k}$, or shorter, the ρ -value of π given $\underline{y}_{<k}$, as*

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

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

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

⁵A dynamic horizon like $m \rightsquigarrow m_k = k^2$ can lead to policies with very poor performance [Hut00, Ch.4].

Remarks:

- $\pi(\underline{y}_{<k}\underline{y}_{k:m})$ is actually independent of x_m , since π is chronological.
- Normalization of $V_{k\gamma}$ by Γ_k does not affect the policy p^ρ .
- The definition of p^ρ is independent of k .
- Without normalization by Γ_k the future values would converge to zero for $k \rightarrow \infty$ in every environment for every policy.
- For an MDP environment, a stationary policy, and geometric discounting $\gamma_k \sim \gamma^k$, 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 Definition 1 is possible.
- Setting $\gamma_k = 1$ for $k \leq m$ and $\gamma_k = 0$ for $k > m$ gives back the undiscounted model (1) with $V_{1\gamma}^{p\rho} = \frac{1}{m} V_{1m}^{p\rho}$.
- $V_{k\gamma}$ (and w_k^ν defined below) depend on the realized history $\underline{y}_{<k}$.

Similarly to the previous sections one can prove the following properties:

Theorem 5 (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

$$V_{k\gamma}^{\pi\xi} = \sum_{\nu \in \mathcal{M}} w_k^\nu V_{k\gamma}^{\pi\nu} \quad \text{and} \quad V_{k\gamma}^{*\xi} \leq \sum_{\nu \in \mathcal{M}} w_k^\nu V_{k\gamma}^{*\nu}$$

$$\text{where} \quad \xi(\underline{y}_{<k}\underline{y}_{k:m}) = \sum_{\nu \in \mathcal{M}} w_k^\nu \nu(\underline{y}_{<k}\underline{y}_{k:m}) \quad \text{with} \quad w_k^\nu := w_\nu \frac{\nu(\underline{y}_{<k})}{\xi(\underline{y}_{<k})}$$

The conditional representation of ξ can be proven by dividing the definition (7) of $\xi(\underline{y}_{1:m})$ by $\xi(\underline{y}_{<k})$ and by using Bayes rules (1). 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 6 (Pareto optimality) For every k and history $\underline{y}_{<k}$ 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 3 (Value difference relation)

$$0 \leq V_{k\gamma}^{*\nu} - V_{k\gamma}^{\tilde{\pi}_k\nu} =: \Delta_k^\nu \quad \Rightarrow \quad 0 \leq V_{k\gamma}^{*\nu} - V_{k\gamma}^{p^\xi\nu} \leq \frac{1}{w_k^\nu} \Delta_k^\nu \quad \text{with} \quad \Delta_k := \sum_{\nu \in \mathcal{M}} w_k^\nu \Delta_k^\nu$$

The proof of Theorem 6 and Lemma 3 follows the same steps as for Theorem 2 and Lemma 1 with appropriate replacements. The proof of the analogue of the convergence Theorem 4 involves one additional step. We abbreviate “with μ probability 1” by w. μ .p.1.

Theorem 7 (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 ν -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 true μ -environment, i.e.*

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

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 usually 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 7 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 Theorem 4(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 $\delta_\nu(k) := \Delta_k^\nu = V_{k\gamma}^{*\nu} - V_{k\gamma}^{\tilde{\pi}_k\nu}$. Since we assumed bounded rewards $0 \leq r \leq r_{max}$ and $V_{k\gamma}^{*\nu}$ is a weighted average of rewards we have

$$V_{k\gamma}^{*\mu} \leq r_{max} \quad \text{and} \quad V_{k\gamma}^{\tilde{\pi}_k\mu} \geq 0 \quad \implies \quad 0 \leq \delta_\nu(k) = \Delta_k^\nu \leq c := r_{max}.$$

The following inequalities follow from Lemma 3:

$$0 \leq V_{k\gamma}^{*\mu} - V_{k\gamma}^{p^\xi\mu} \leq \frac{1}{w_k^\mu} \Delta_k = \frac{1}{w_k^\mu} \delta(k) \stackrel{?}{\rightarrow} 0 \quad (10)$$

The premise in Theorem 7 is that $\delta_\nu(k) = V_{k\gamma}^{*\nu} - V_{k\gamma}^{\tilde{\pi}_k\nu} \rightarrow 0$ for $k \rightarrow \infty$ which implies $\delta(k) \rightarrow 0$ (w. μ .p.1) by Lemma 2(ii). What is new and what remains to be shown is that w_k^μ is bounded from below in order to have convergence of (10) to zero. 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 Bayes 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. \square

6 Markov Decision Processes

From all possible environments, Markov (decision) processes are probably the most intensively studied ones. To give an example, we apply Theorems 4 and 7 to ergodic Markov decision processes, but we will be very brief.

Definition 3 (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 $\mathbf{y}_{<k} y_k$ does only depend on the last action $y_k \in \mathcal{Y}$ and the last observation x_{k-1} , i.e. if $\mu(\mathbf{y}_{<k} y_k \underline{x}_k) = \mu(\mathbf{y}_{k-1} \underline{x}_k)$. In this case x_k is called a state, \mathcal{X} the state space, and $\mu(\mathbf{y}_{k-1} \underline{x}_k)$ 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 $\mu(\mathbf{y}_{k-1} \underline{x}_k)$ is independent of the action y_{k-1} 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 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 on 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 8 (Self-optimizing policies for ergodic MDPs) *There exist self-optimizing policies \tilde{p}_m for the class of ergodic MDPs in the sense that*

$$i) \quad \exists \tilde{p}_m \forall \nu \in \mathcal{M}_{MDP1} : \frac{1}{m} V_{1m}^{*\nu} - \frac{1}{m} V_{1m}^{\tilde{p}_m \nu} \leq c_\nu m^{-1/3} \xrightarrow{m \rightarrow \infty} 0,$$

where c_ν are some constants. In the discounted case, if the discount sequence γ_k has unbounded effective horizon $h_k^{\text{eff}} \xrightarrow{k \rightarrow \infty} \infty$, then there exist self-optimizing policies $\tilde{\pi}_k$ for the class of ergodic MDPs in the sense that

$$ii) \quad \exists \tilde{\pi}_k \forall \nu \in \mathcal{M}_{MDP1} : V_{k\gamma}^{\tilde{\pi}_k \nu} \xrightarrow{k \rightarrow \infty} V_{k\gamma}^{*\nu} \quad \text{if} \quad \frac{\gamma_{k+1}}{\gamma_k} \rightarrow 1.$$

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 3, demanding that the Markov process is ergodic under *every* policy. See, for instance, [KV86, Ber95]. We will only briefly outline one algorithm satisfying Theorem 8 without trying to optimize performance.

Proof idea: 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 $T_{ss'}^a \equiv \nu(a \underline{ss}')$ 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 $\geq \frac{k_0}{|\mathcal{X}|} (T_{min}/|\mathcal{Y}|)^{|\mathcal{X}|} \sim k_0$.⁶ 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 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 $\mathcal{Y}_{<k}$ is simply ignored and the analogue to $m \rightarrow \infty$ is $h_k^{eff} \rightarrow \infty$ for $k \rightarrow \infty$, which is ensured by $\frac{\gamma_{k+1}}{\gamma_k} \rightarrow \infty$. 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 policy⁷ based 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 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

$$V_{k_0 \gamma}^{\tilde{\pi}_k \nu} \rightarrow V_{k_0 \gamma}^{*\nu} \quad \text{for} \quad \Delta = k_0 - k \rightarrow \infty, \quad (11)$$

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 (11) this shows $V_{k\gamma}^{\tilde{\pi}_k \nu} \rightarrow V_{k\gamma}^{*\nu}$ for $k_0 := k + \Delta(k)$. \square

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. [Hut00] gives a detailed account on discount and horizon issues, and motivates $h_k^{eff} \rightarrow \infty$ philosophically.

Together with Theorems 4 and 7, Theorem 8 immediately implies that policy p^ξ is self-optimizing for the class of ergodic MDPs.

Corollary 1 (Policy p^ξ is self-optimizing for ergodic MDPs) *If \mathcal{M} is a finite or countable class of ergodic MDPs, and $\xi() := \sum_{\nu \in \mathcal{M}} w_\nu \nu()$, then policies 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*

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

If \mathcal{M} is finite, then the speed of the first convergence is at least $O(m^{-1/3})$.

⁶For $T_{ss'}^a = 0$ the estimate $\hat{T}_{ss'}^a = 0$ is exact.

⁷For non-geometric discounts as here, optimal policies are, in general, *not* stationary.

7 Conclusions

Summary: We studied agents acting in general probabilistic environments with reinforcement feedback. We only assumed that the true environment μ belongs to a known class of environments \mathcal{M} , but is otherwise unknown. We showed that the Bayes-optimal policy p^ξ based on the Bayes-mixture $\xi = \sum_{\nu \in \mathcal{M}} w_\nu \nu$ is Pareto-optimal and self-optimizing if \mathcal{M} admits self-optimizing policies. The class of ergodic MDPs admitted self-optimizing policies w.r.t. the average value and w.r.t. the discounted value if the effective horizon grew indefinitely.

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 $\xi = \int w_\theta \mu_\theta d\theta$. Under some mild (differentiability and existence) conditions, most results of this work remain valid in some form, especially Corollary 1 for *all* ergodic MDPs.

Bayesian self-optimizing policy: Policy p^ξ with unbounded effective horizon for ergodic MDPs is the first purely Bayesian self-optimizing consistent policy for ergodic MDPs. The policies of all previous approaches were either hand crafted, like the ones in the proof of Theorem 8, or were Bayesian with a pre-chosen horizon m , or with geometric discounting γ with finite effective horizon (which does not allow self-optimizing policies) [KV86, Ber95]. The combined conditions $\Gamma_k < \infty$ and $\frac{\gamma^{k+1}}{\gamma^k} \rightarrow 1$ allow a consistent self-optimizing Bayes-optimal policy based on mixtures.

Bandits: Bandits are a special subclass of ergodic MDPs. In a two-armed bandit problem you pull repeatedly one of two levers resulting in a gain of A\$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 1 shows that policy p^ξ yields asymptotically optimal payoff. The discounted unbounded horizon approach and result is, to the best of our knowledge, even new when restricted to Bandits.

Other environmental classes: Bandits, i.i.d. processes, classification tasks, and many more are all special (degenerate) cases of ergodic MDPs, for which Corollary 1 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 4 and 7.

Outlook: Future research could be the derivation of non-asymptotic bounds, possibly along the lines of [Hut01]. To get good bounds one may have to exploit extra properties of the environments, like the mixing rate of MDPs [KS98]. Another possibility is to search for other performance criteria along the lines of [Hut00, Ch.6], especially for the universal prior [Sol78] and for the Speed prior [Sch02]. Finally, instead of convergence of the expected reward sum, studying convergence with high probability of the actual reward sum would be interesting.

References

- [Bel57] R. Bellman. *Dynamic Programming*. Princeton University Press, New Jersey, 1957.
- [Ber95] D. P. Bertsekas. *Dynamic Programming and Optimal Control, Vol. (I) and (II)*. Athena Scientific, Belmont, Massachusetts, 1995. Volumes 1 and 2.
- [BT00] R. I. Brafman and M. Tennenholtz. A near-optimal polynomial time algorithm for learning in certain classes of stochastic games. *Artificial Intelligence*, 121(1–2):31–47, 2000.
- [Doo53] J. L. Doob. *Stochastic Processes*. John Wiley & Sons, New York, 1953.
- [Hut00] M. Hutter. A theory of universal artificial intelligence based on algorithmic complexity. Technical Report cs.AI/0004001, 62 pages, 2000. <http://arxiv.org/abs/cs.AI/0004001>.
- [Hut01] M. Hutter. General loss bounds for universal sequence prediction. *Proceedings of the 18th International Conference on Machine Learning (ICML-2001)*, pages 210–217, 2001.
- [KLM96] L. P. Kaelbling, M. L. Littman, and A. W. Moore. Reinforcement learning: a survey. *Journal of AI research*, 4:237–285, 1996.
- [KS98] M. 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.
- [KV86] P. R. Kumar and P. P. Varaiya. *Stochastic Systems: Estimation, Identification, and Adaptive Control*. Prentice Hall, Englewood Cliffs, NJ, 1986.
- [LV97] M. Li and P. M. B. Vitányi. *An introduction to Kolmogorov complexity and its applications*. Springer, 2nd edition, 1997.
- [RN95] S. J. Russell and P. Norvig. *Artificial Intelligence. A Modern Approach*. Prentice-Hall, Englewood Cliffs, 1995.
- [SB98] R. Sutton and A. Barto. *Reinforcement learning: An introduction*. Cambridge, MA, MIT Press, 1998.
- [Sch02] J. Schmidhuber. The Speed Prior: a new simplicity measure yielding near-optimal computable predictions. *Proceedings of the 15th Annual Conference on Computational Learning Theory (COLT 2002)*, 2002.
- [Sol78] R. J. Solomonoff. Complexity-based induction systems: comparisons and convergence theorems. *IEEE Trans. Inform. Theory*, IT-24:422–432, 1978.