

Contents lists available at [ScienceDirect](https://www.sciencedirect.com)

Information Sciences

journal homepage: www.elsevier.com/locate/insDiscounted fully probabilistic design of decision rules [☆]Miroslav Kárný^{*}, Soňa Molnárová

The Czech Academy of Sciences, Institute of Information Theory and Automation, Department of Adaptive Systems, vodárenskou věží 4, Prague 8, 182 00, Czech Republic

ARTICLE INFO

Keywords:

Design principles
Kullback-Leibler's divergence
Probabilistic techniques
Discounting
Closed loop

ABSTRACT

Axiomatic fully probabilistic design (FPD) of optimal decision rules strictly extends the decision making (DM) theory represented by Markov decision processes (MDP). This means that any MDP task can be approximated by an explicitly found FPD task whereas many FPD tasks have no MDP equivalent. MDP and FPD model the closed loop — the coupling of an agent and its environment — via a joint probability density (pd) relating the involved random variables, referred to as behaviour. Unlike MDP, FPD quantifies agent's aims and constraints by an *ideal pd*. The ideal pd is high on the desired behaviours, small on undesired behaviours and zero on forbidden ones. FPD selects the optimal decision rules as the minimiser of Kullback-Leibler's divergence of the closed-loop-modelling pd to its ideal twin. The proximity measure choice follows from the FPD axiomatics.

MDP minimises the expected total loss, which is usually the sum of discounted partial losses. The discounting reflects the decreasing importance of future losses. It also diminishes the influence of errors caused by:

- ▶ the imperfection of the employed environment model;
- ▶ roughly-expressed aims;
- ▶ the approximate learning and decision-rules design.

The established FPD cannot currently account for these important features. The paper elaborates the missing discounted version of FPD. This non-trivial filling of the gap in FPD also employs an extension of dynamic programming, which is of an independent interest.

1. Introduction

An agent — a human, a device or a mixed group of both, referred to as “it” — chooses its actions in order to meet its aims. This is the core of any decision making that always runs under uncertainty. The inspected prescriptive Bayesian DM theory has its roots in [1,2]. It underlies the theory of Markov decision processes [3], which is the standard way of designing optimal, action-generating, decision rules. Stochastic control does the same [4,5] but its stress and vocabulary differ. They are often used interchangeably or jointly¹ [6].

[☆] This research was supported by EU-COST Action CA21169.

^{*} Corresponding author.

E-mail addresses: school@utia.cas.cz (M. Kárný), molnarova@utia.cas.cz (S. Molnárová).

URLs: <https://www.utia.cas.cz/people/karny> (M. Kárný), <https://www.utia.cas.cz/people/moln-rov> (S. Molnárová).

¹ References are just samples from a much more extensive set.

<https://doi.org/10.1016/j.ins.2024.121578>

Received 2 January 2024; Received in revised form 27 August 2024; Accepted 19 October 2024

Available online 22 October 2024

0020-0255/© 2024 Elsevier Inc. All rights are reserved, including those for text and data mining, AI training, and similar technologies.

The MDP design often employs discounting, which perceives DM results at a distant future as less important [7–9]. This weakens the adverse impact of imprecision (of any origin) on the design of decision rules [10,11]. It also simplifies the error analysis of various approximate designs of decision rules [12,13].

This text and [14] show that FPD, an abstract axiomatic version [14,15] of the model reference control [16], strictly extends MDP. Discounted FPD has not been yet developed regardless of its rich history [17–21], its tight connection with the independently proposed KL control² [23–26] and with the KL regularised control [27,28]. The paper fills this gap by elaborating the discounted FPD.

MDP tasks mapped on FPD tasks, see Proposition 3 below, naturally come with discounting. However, FPD tasks with no MDP equivalents lack it. The discounting extension to *any* FPD task is highly desirable to reach its positive effects.

The choice of the discounting rate is hard even if it has a monetary interpretation [29]. A sound choice is vital especially when DM is connected with sensitive aspects, say health [9]. The choice is even harder when discounting copes with doubts about persistency of the used environment model or of the quantified aims.

The importance of a sound choice of the discounting rate lies in the impact of the effective shortening of the design horizon. It may decrease the DM quality up to the closed-loop instability [30]. On the other hand, the optimisation without discounting may excite modelling errors and lead to poor DM quality up to the instability [31]. Sec. 4 illustrates the positive effect of discounting on diminishing the consequences of modelling errors.

The choice of the discounting rate is expected to be easier for FPD. The unified probabilistic language of FPD concerns both the closed-loop modelling and the aim expression. Thus FPD inputs, including the discounting rate, are quantified by probabilities. The meaning and role of the discounting rate in the design part is fully analogical to the data and time varying forgetting rate used in learning [32,33]. At the same time, the on-line learning of the forgetting rate is feasible [34]. These facts open a way to a data-based choice of the discounting rate within FPD. Due to the inclusion of MDP tasks in the set of FPD tasks this will help to choose the discounting rate for MDP as well. Initial attempts are in [35].

Layout

Sec. 2 recalls FPD and refines its relation to MDP. A generalisation of the usual dynamic programming [4] arises during this refinement. This overture serves to core Sec. 3, which presents the discounted FPD. Sec. 4 illustrates the impact of the theory numerically. Sec. 5 provides the concluding remarks.

Notation

The text uses decorated mnemonic labels: *a* action, D and d divergences, c closed-loop model, E expectation, *h* horizon, m model, N and n normalisations, L loss, *p* pointer, *r* decision rule(s), *s* state, *t* time, *v* value function, *w* weight. Sanserif fonts mark mappings. *c*, *m* and *r* are probability densities (pds³). Bold fonts mark the set: \mathbf{x} is the set of *x*s and is defined only if needed. \mathbf{x}^h is Cartesian product of *h* sets *x*, *h* is a natural number. := defines the left-hand side by assigning the right-hand side. \propto is proportionality. $f(x_t) := f_t(x_t)$, the double time subscript $t \in \mathbf{t}$ is dropped. $r_t := (r(a_t|s_{t-1}))_{a_t \in \mathbf{a}, s_{t-1} \in \mathbf{s}}$ are decision rules.⁴ The initial state $s_0 \in \mathbf{s}$ implicitly conditions all pds. Superscript ⁱ concerns the ideal pd, ^o indicates optimality.

2. Preliminaries

The FPD recall makes the paper self-reliant. The proved version of the dynamic programming generalises the usual one. It is applicable to a total loss with partial losses dependent on the optimised decision rules.

Note that works [24,37,38] provide rich commented lists of references to a wide range of FPD-related results.

2.1. Fully probabilistic design

DM concerns the closed loop consisting of an agent and its uncertain environment. The agent gradually applies actions $a_t \in \mathbf{a} \neq \emptyset$ at time epochs $t \in \mathbf{t} := \{1, \dots, h\}$, $h \leq \infty$, influencing the closed loop. They stimulate a transition of the closed-loop state $s_{t-1} \in \mathbf{s} \neq \emptyset$ to the state $s_t \in \mathbf{s}$. The states and actions, thought up to the decision horizon *h*, form the closed-loop behaviours

$$\mathbf{b} := (s_h, a_h, \dots, s_1, a_1) \in \mathbf{b} := (\mathbf{s}, \mathbf{a})^h.$$

The agent selects chooses $(a_t)_{t \in \mathbf{t}}$ via randomised decision rules

$$\mathbf{r} \in \mathbf{r} := \left\{ r(\mathbf{b}) := \prod_{t \in \mathbf{t}} r(a_t | s_{t-1}) \right\}.$$

The conditional pds $r(a_t | s_{t-1})$ model the causal decision rules *r*. The decision-rules-dependent joint pd $c^r(\mathbf{b})$ completely describes random behaviours $\mathbf{b} \in \mathbf{b}$. The state definition and the chain rule for pds [39] imply the factorised closed-loop model

$$c^r(\mathbf{b}) = \prod_{t \in \mathbf{t}} m(s_t | a_t, s_{t-1}) r(a_t | s_{t-1}) := m(\mathbf{b}) r(\mathbf{b}). \quad (1)$$

² KL stands for Kullback-Leibler's divergence [22]. The use of KL control is adopted.

³ Pd is Radon-Nikodým's derivative with respect to either Lebesgue's or counting measure [36]. Lebesgue's notation is taken as the generic one.

⁴ As usual, | separates conditions in pds as well as in expectations.

The known environment model $m(b) := \prod_{t \in \mathcal{T}} m(s_t | a_t, s_{t-1})$ consists of the conditional pds $m(s_t | a_t, s_{t-1})$, $s_t, s_{t-1} \in \mathcal{S}$, $a_t \in \mathcal{A}$, $t \in \mathcal{T}$, modelling state transitions.

Remark 1 (On the Closed-Loop State Relevance). Many DM setups neglect the need to use the *closed-loop* state whenever the environment model results from on-line estimation [39]. Indeed, if learning and acting run in parallel, the state includes values of the used statistic as actions influence them. The optimal decision rules thus balance exploration and exploitation efforts [40]. Without a reflection of this fact, the exploration becomes the extra challenging task requiring sophisticated ad hoc techniques [41]. \square

FPD quantifies agent's aims and constraints by an ideal (desired) joint pd $c^i(b)$, $b \in \mathcal{B}$. The pd $c^i(b)$ assigns high values to desired behaviours b , small values to undesired bs and zero to forbidden bs . The ideal closed-loop model factorises in the same way as the pd $c^f(b)$ in (1)

$$c^i(b) = \prod_{t \in \mathcal{T}} m^i(s_t | a_t, s_{t-1}) r^i(a_t | s_{t-1}) := m^i(b) r^i(b). \quad (2)$$

The agent-selected ideal environment model

$$m^i(b) := \prod_{t \in \mathcal{T}} m^i(s_t | a_t, s_{t-1}), \quad s_t, s_{t-1} \in \mathcal{S}, \quad a_t \in \mathcal{A},$$

combines the conditional pds $m^i(s_t | a_t, s_{t-1})$, $t \in \mathcal{T}$, modelling the desired state transitions. The agent-selected ideal decision rules

$$r^i(b) := \prod_{t \in \mathcal{T}} r^i(a_t | s_{t-1}), \quad a_t \in \mathcal{A}, \quad s_{t-1} \in \mathcal{S}$$

consist of the desired decision rules given by the pds $r^i(a_t | s_{t-1})$. Note that the definition of $m^i(b)$ in (2) implies the normalisation used in Sec. 2.2

$$\int_{s^h} m^i(b) d(s_h, \dots, s_1) = 1. \quad (3)$$

The FPD-optimal decision rules $r^\circ \in \mathcal{R}$ [14] minimise Kullback-Leibler's divergence (KL, [22]) $D(c^f || c^i)$ of c^f to c^i

$$r^\circ \in \text{Arg min}_{r \in \mathcal{R}} D(c^f || c^i) := \text{Arg min}_{r \in \mathcal{R}} \int_{\mathcal{B}} c^f(b) \ln \left(\frac{c^f(b)}{c^i(b)} \right) db. \quad (4)$$

Dynamic programming [4] provides the optimal decision rules r° (4). They are designed in a few steps using the following lemmas.

Lemma 1 (Additive Form of KL). For decision rules $r \in \mathcal{R}$ and the function

$$L_D^r(s_t, a_t, s_{t-1}) := \ln \left(\frac{m(s_t | a_t, s_{t-1}) r(a_t | s_{t-1})}{m^i(s_t | a_t, s_{t-1}) r^i(a_t | s_{t-1})} \right),$$

the functional D in (4) is the decision-rules-dependent expectation (E^r) of the total loss L_D^{hr} with the decision-rules-dependent addends L_D^r . KL in (4) reads

$$\begin{aligned} D(c^f || c^i) &= \sum_{t \in \mathcal{T}} \int_{(s, a, s)} c^f(s_t, a_t, s_{t-1}) L_D^r(s_t, a_t, s_{t-1}) d(s_t, a_t, s_{t-1}) \\ &= \sum_{t \in \mathcal{T}} \int_{\mathcal{S}} c^f(s_{t-1}) \left[\int_{(s, a)} m(s_t | a_t, s_{t-1}) r(a_t | s_{t-1}) \right. \\ &\quad \left. \times \ln \left(\frac{m(s_t | a_t, s_{t-1}) r(a_t | s_{t-1})}{m^i(s_t | a_t, s_{t-1}) r^i(a_t | s_{t-1})} \right) d(s_t, a_t) \right] ds_{t-1} := E^r[L_D^{hr}]. \end{aligned} \quad (5)$$

The employed marginal pd $c^f(s_{t-1})$ of $c^f(b)$ is independent of the decision rules $r(a_t | s_{t-1})$ with time indices $\tau \geq t$.

Proof. Formula (5) follows from:

- ▶ the product in the ratio $c^f(b)/c^i(b)$, see (1), (2), yielding the sum of logarithms;
- ▶ linearity of the integration in its argument;
- ▶ the dependence of $L_D^r(s_t, a_t, s_{t-1})$ on (s_t, a_t, s_{t-1}) allows other entries of $b \in \mathcal{B}$ to integrate out of the joint pd $c^f(b)$ reducing it to the marginal pd $c^f(s_t, a_t, s_{t-1})$;
- ▶ the chain rule for pds applied to $c^f(s_t, a_t, s_{t-1})$ using definitions of $m(s_t | a_t, s_{t-1})$, $r(a_t | s_{t-1})$, see (1), and of $c^f(s_{t-1})$;
- ▶ Fubini's theorem for multiple integrals [36].

The marginal pd $c^r(s_{t-1})$ results from the next multiple integration that runs over variables listed in differentials

$$c^r(s_{t-1}) = \int_{(s^{h-1}, a^h)} \prod_{\tau \in t} m(s_\tau | a_\tau, s_{\tau-1}) \times r(a_\tau | s_{\tau-1}) d(s_h, a_h, \dots, s_t, a_t, a_{t-1}, s_{t-2}, a_{t-2}, \dots, s_1, a_1).$$

The pds for $\tau \geq t$ integrate to unity. This shows the independence of $c^r(s_{t-1})$ of $r(a_\tau | s_{\tau-1})$, $\tau \geq t$. \square

The expression (5) and the proved independence allow us to find the optimal decision rules (4) using the backward induction known as dynamic programming [4]. It provides the desired, causal, optimal decision rules. We need its slight generalisation, which considers the minimisation of the decision-rules-dependent expectation of the total loss L^{hr} with its addends dependent on the optimised decision rules. The optimal randomised decision rules r^o are

$$r^o \in \text{Arg min}_{r \in \mathcal{R}} \int_b \sum_{t \in t} L^r(s_t, a_t, s_{t-1}) c^r(b) db := \text{Arg min}_{r \in \mathcal{R}} E^r[L^{hr}], \quad (6)$$

with the partial loss $L^r(s_t, a_t, s_{t-1})$ dependent on the rules $r(a_\tau | s_{\tau-1})$ with $\tau \leq t$.

Lemma 1 shows that the optimisation (4) is a subcase of (6). The existence of other cases might be important, see Remark 2. The next lemma solves the task (6).

Lemma 2 (Dynamic programming). *Let there exist stabilising decision rules \bar{r} , making the expected total loss $E^{\bar{r}}[L^{hr}]$ in (6) finite. Let us define the value functions*

$$v(s_{t-1}) := \min_{(r_\tau \in \mathcal{R})_{\tau \geq t}} \sum_{\tau \geq t} E^r[L^r(s_\tau, a_\tau, s_{\tau-1}) | s_{t-1}], \quad s_{t-1} \in \mathcal{S}, \quad t \in t. \quad (7)$$

Then, $v(s_{t-1})$ is bounded, $v(s_h) = 0$, and the next functional recursion holds

$$v(s_{t-1}) = \min_{r_t \in \mathcal{R}} \int_{(s, a)} m(s_t | a_t, s_{t-1}) r_t(a_t | s_{t-1}) [L^r(s_t, a_t, s_{t-1}) + v(s_t)] d(s_t, a_t). \quad (8)$$

The optimal decision rules r^o (6) are minimisers r_t^o in (8) and

$$v(s_0) = E^{r^o}[L^{hr^o}]. \quad (9)$$

Proof. The existence of stabilising decision rules validates the next manipulations. For any $t \in t$, the minimisation in (6) unfolds as $\min_{(r_\tau \in \mathcal{R})_{\tau < t}} \left[\min_{(r_\tau \in \mathcal{R})_{\tau \geq t}} E^r[L^{hr}] \right]$. The inner minimisation reduces to

$$\zeta_{t-1} := \min_{(r_\tau \in \mathcal{R})_{\tau \geq t}} E^r \left[\sum_{\tau \geq t} L^r(s_\tau, a_\tau, s_{\tau-1}) \right] = E^r[v(s_{t-1})]. \quad (10)$$

The second equality in (10) follows from:

- ▶ the chain rule for expectation $E^r[\bullet] = E^r[E^r[\bullet | s_{t-1}]]$ [36];
- ▶ the independence of $c^r(s_{t-1})$ from $(r_\tau)_{\tau \geq t}$, which implies the commutativity of the outer unconditional expectation and the minimisation over $(r_\tau)_{\tau \geq t}$;
- ▶ the value-function definition (7).

The smallest total expected loss in (6) coincides with ζ_0 . Thus, the minimising rules in the ζ_0 definition (10) form the optimal decision rules. Moreover, the implicit conditioning by the known s_0 makes $v(s_0) = \zeta_0$. This proves (9). It remains to prove the validity of the recursion (8).

The independence of $L^r(s_t, a_t, s_{t-1})$ from rules $(r_\tau)_{\tau > t}$ implies

$$\begin{aligned} v(s_{t-1}) &= \min_{r_t \in \mathcal{R}} \left[E^r[L^r(s_t, a_t, s_{t-1}) | s_{t-1}] + \min_{(r_\tau \in \mathcal{R})_{\tau > t}} \sum_{\tau > t} E^r[L^r(s_\tau, a_\tau, s_{\tau-1}) | s_{t-1}] \right] \\ &= \min_{r_t \in \mathcal{R}} \left[E^r \left[L^r(s_t, a_t, s_{t-1}) + \min_{(r_\tau \in \mathcal{R})_{\tau > t}} \sum_{\tau > t} E^r[L^r(s_\tau, a_\tau, s_{\tau-1}) | s_t] \right] | s_{t-1} \right] \\ &= \min_{r_t \in \mathcal{R}} E^r[L^r(s_t, a_t, s_{t-1}) + v(s_t) | s_{t-1}]. \end{aligned}$$

The justification of the respective equalities above is as follows.

- ▶ The first equality uses (7) and the independence of $E^r[L^r(s_t, a_t, s_{t-1}) | s_{t-1}]$ from the rules $(r_\tau)_{\tau > t}$.
- ▶ The second equality utilises the chain rule for expectations and the state definition. They imply $E^r[\bullet | s_{t-1}] = E^r[E^r[\bullet | s_t, s_{t-1}] | s_{t-1}] = E^r[E^r[\bullet | s_t] | s_{t-1}]$.

► The third equality results from the fact that the minimisation over $(r_t)_{t>1}$ and the outer expectation in $E^r[E^i[\bullet|s_t]|s_{t-1}]$ commute, and repetitive application of the definition (7). ◻

Remark 2 (On Possible FPD Generalisations). The FPD axiomatics [14] singled out KL among all f-divergences [42]. It avoids an unjustified coupling of optimal decision rules of two completely independent DMs solved as a single DM task. If this degenerate case is avoided by the problem formulation, then the axiomatics and Lemma 2 allow us to use any divergence in the role of the expected partial loss. This opens a way to a rich generalisation of FPD [43]. ◻

Lemmas 1, 2 serve to the novel proof of the solution of FPD. For an alternative proof, see [17].

Proposition 1 (FPD). The backward induction reduces to the next backward functional recursion for $n(s_{t-1}) \in [0, 1]$, $s_{t-1} \in s$, $t \in t$, with $n(s_h) := 1$, and $a_t \in a$,

$$\begin{aligned}
 d(a_t, s_{t-1}) &:= \int_s m(s_t|a_t, s_{t-1}) \ln \left(\frac{m(s_t|a_t, s_{t-1})}{n(s_t)m^i(s_t|a_t, s_{t-1})} \right) ds_t \\
 n(s_{t-1}) &:= \int_a r^i(a_t|s_{t-1}) \exp[-d(a_t, s_{t-1})] da_t.
 \end{aligned}
 \tag{11}$$

The reached minimum is $v(s_0) := -\ln(n(s_0)) = D(c^r^\circ || c^i)$ in (4). The optimal decision rules are

$$r^\circ(a_t|s_{t-1}) = \frac{r^i(a_t|s_{t-1}) \exp[-d(a_t, s_{t-1})]}{n(s_{t-1})} \propto r^i(a_t|s_{t-1}) \exp[-d(a_t, s_{t-1})].
 \tag{12}$$

Proof. For non-negative KL, the value function $v(s_{t-1}) \geq 0$. Thus, it can be expressed in the form $-\ln(n(s_{t-1})) := v(s_{t-1})$. The non-negativity of $v(s_{t-1})$ and the identity $v(s_h) = 0$ give $n(s_t) \in [0, 1]$ and $n(s_h) = 1$. For the addends $L^r(s_t, a_t, s_{t-1}) := L_D^r(s_t, a_t, s_{t-1})$ in the expression (5), and for $r_t := r(a_t|s_{t-1})$, $a_t \in a$, $s_{t-1} \in s$, the recursion (8) gets the form

$$\begin{aligned}
 v(s_{t-1}) &= \min_{r_t \in \mathcal{E}^r} \int_{(s,a)} m(s_t|a_t, s_{t-1}) \\
 &\quad \times r(a_t|s_{t-1}) \left[\ln \left(\frac{m(s_t|a_t, s_{t-1})r(a_t|s_{t-1})}{m^i(s_t|a_t, s_{t-1})r^i(a_t|s_{t-1})} \right) + v(s_t) \right] d(s_t, a_t).
 \end{aligned}$$

The definition $v(s_t) := -\ln(n(s_t))$ and the fact that the pd $m(s_t|a_t, s_{t-1})$ integrates to one over the states s_t in the set s imply

$$\begin{aligned}
 v(s_{t-1}) &= \min_{r_t \in \mathcal{E}^r} \int_a r(a_t|s_{t-1}) \left[\ln \left(\frac{r(a_t|s_{t-1})}{r^i(a_t|s_{t-1})} \right) \right. \\
 &\quad \left. + \underbrace{\int_s m(s_t|a_t, s_{t-1}) \ln \left(\frac{m(s_t|a_t, s_{t-1})}{n(s_t)m^i(s_t|a_t, s_{t-1})} \right) ds_t}_{d(a_t, s_{t-1}) :=} \right] da_t.
 \end{aligned}$$

The second integral above is the function $d(a_t, s_{t-1})$ used in the proposition formulation. The final recursion for the value function uses this function and a few elementary manipulations

$$\begin{aligned}
 v(s_{t-1}) &= \min_{r_t \in \mathcal{E}^r} \left[\int_a r(a_t|s_{t-1}) \ln \left(\frac{r(a_t|s_{t-1})}{r^\circ(a_t|s_{t-1})} \right) da_t \right. \\
 &\quad \left. - \ln \left(\underbrace{\int_a r^i(a_t|s_{t-1}) \exp(-d(a_t, s_{t-1})) da_t}_{n(s_{t-1}) :=} \right) \right] \\
 &= \min_{r_t \in \mathcal{E}^r} \int_a r(a_t|s_{t-1}) \ln \left(\frac{r(a_t|s_{t-1})}{r^\circ(a_t|s_{t-1})} \right) da_t - v(s_{t-1}).
 \end{aligned}$$

In the first equality above, just the logarithm $\ln(n(s_{t-1}))$ of the $r(a_t|s_{t-1})$ -independent normalisation factor $n(s_{t-1})$ is added and subtracted. Then, KL of the optimised pd $r(a_t|s_{t-1})$ to $r^\circ(a_t|s_{t-1})$ (12) of $a_t \in a$, conditioned on $s_{t-1} \in s$, appears. It gets its smallest zero value for $r(a_t|s_{t-1}) = r^\circ(a_t|s_{t-1})$ and provides the value function $v(s_{t-1}) = -\ln(n(s_{t-1})) \geq 0$. ◻

2.2. Mutual relationship of FPD and MDP

Introduction claimed that the FPD tasks strictly extend the class of MDP tasks. This subsection proves this claim and primarily points to a “normalisation problem” similar to that faced when formulating the discounted FPD in Sec. 3.

Propositions 2, 3 below motivate the way to the discounted FPD, Sec. 3. They are of an independent interest. Proposition 2 was proved in [44]. It shows how FPD reduces to MDP. The novel observation warns that this reduction cannot be directly inverted. Proposition 3 slightly refines the claim of [14] that the set of FPD tasks densely extends the set of MDP tasks. Thus, an FPD task can be found to each MDP task providing the decision rules arbitrarily close to the MDP-optimal ones.

FPD uses the ideal pd c^i , which is the product (2) of the ideal environment model m^i and the ideal decision rules r^i . Often, the agent is unwilling or unable to specify some factors forming them. In this case, it is wise to let FPD to identify these factors with their counterparts in c^i . When applied to a decision rule, then this *leave-to-the-fate* option (LTF, [44]) reads

$$r^i := r \Leftrightarrow \text{the ideal decision rule equals to the optimised one.} \quad (13)$$

Proposition 2 (FPD under LTF (13) becomes MDP). Under LTF, FPD (4) reduces to MDP with the total, decision-rules-independent,⁵ loss

$$L_D^h(b) := \sum_{i \in \mathcal{I}} L_D(s_t, a_t, s_{t-1}) := \sum_{i \in \mathcal{I}} \ln \left(\frac{m(s_t | a_t, s_{t-1})}{m^i(s_t | a_t, s_{t-1})} \right). \quad (14)$$

The MDP-optimal decision rules r^o , assigned to the sum of arbitrary decision-rules-independent partial losses $L(s_t, a_t, s_{t-1})$ (including L_D), are deterministic. They generate the minimising arguments $a^o(s_{t-1})$ in

$$v(s_{t-1}) := \min_{a_t \in \mathcal{A}} \int_s [L(s_t, a_t, s_{t-1}) + v(s_t)] m(s_t | a_t, s_{t-1}) ds_t. \quad (15)$$

The backward functional recursion (15) coincides with the usual dynamic programming [4] and $v(s_h) = 0$ initiates it.

Proof. Under LTF, decision rules r and ideal decision rules r^i cancel in the logarithm used in the definition (4). The minimised KL becomes the expectation (given by the closed-loop model $c^i = m r$ (1)) of the total loss $L_D^h(b)$ (14). The formulation coincides with the MDP in which the optimised functional (8) is linear in the optimised r_t . This implies the standard deterministic solution (15), [4]. \square

Proposition 2 seemingly shows that any MDP, determined by the environment model $m(b)$ and by its given total loss $L^h(b)$, can be seen as a special case of FPD that:

- ▶ applies LTF (13), and;
- ▶ uses the ideal environment model m^i of the form

$$m^i(b) \propto m(b) \exp [-L^h(b)]. \quad (16)$$

This conclusion is *invalid* as the normalisation factor N in (16), guaranteeing the normalisation (3), depends on actions

$$N(a_h, \dots, a_1) = \int_{s^h} m(b) \exp [-L^h(b)] d(s_h, \dots, s_1).$$

Thus, the choice of the ideal model (16) for LTF (13) leads to the minimisation of

$$D(c^i || m^i r) = E^r [L^h] + E^r [\ln(N)], \quad (17)$$

instead of the desired minimisation of $E^r [L^h]$ made within MDP. The good news is that any MDP task, given by an environment model $m(b)$ and a total loss $L^h(b)$, $b \in \mathcal{B}$, can be arbitrarily well approximated by a specific FPD task.

Proposition 3 (FPD Strictly and Densely Extends MDP).

1. There are FPD tasks having no MDP equivalent.
2. Any MDP with a decision-rules-independent total loss L^h and stabilising decision rules \bar{r} making $E^{\bar{r}} [L^h] < \infty$ can be approximated to an arbitrary precision by an FPD task with the same environment model $m(b)$ and the ideal pd

$$c^{i,\lambda}(b) := \frac{\bar{c}(b) \exp[-L^h(b)/\lambda]}{\int_{\mathcal{B}} \bar{c}(b) \exp[-L^h(b)/\lambda] db}. \quad (18)$$

It is given by $\lambda > 0$ and a pd $\bar{c}(b) > 0$ on \mathcal{B} making

⁵ It explicitly depends on the behaviour realisation, not on the decision rules that influenced it.

$$\begin{aligned}
 & i) \text{ the pd } c^{i\lambda}(b) \text{ (18) a proper one with } \int_b \tilde{c}(b) \exp[-L^h(b)/\lambda] db < \infty \\
 & ii) D(\text{mr}||\tilde{c}) < \infty \text{ for deterministic decision rules } r \text{ with } E^r[L^h] < \infty.
 \end{aligned} \tag{19}$$

Proof. ad 1. The randomised nature of the FPD optimal decision rules, see Proposition 1, implies that they cannot be optimal for any MDP with a unique minimiser $a^o(s_{t-1})$ in dynamic programming, see Proposition 2.

ad 2. The following decision rules $r^o, r^{o\lambda} \in \mathbf{r}$ are well-defined, for $\lambda > 0$ and the pd $\tilde{c}(b) > 0$ used in (18),

$$\begin{aligned}
 & r^o \in \text{Arg min}_{r \in \mathbf{r}} E^r[L^h], \text{ which guarantees } E^{r^o}[L^h] \leq E^r[L^h] < \infty \\
 & r^{o\lambda} \in \text{Arg min}_{r \in \mathbf{r}} [E^r[L^h] + \lambda D(\text{mr}||\tilde{c})] \underbrace{=}_{(18)} \text{Arg min}_{r \in \mathbf{r}} D(\text{mr}||c^{i\lambda}).
 \end{aligned} \tag{20}$$

The joint pd $c^{i\lambda}(b)$ (18) exists due to the assumption (19) i). The definitions (20) of r^o and $r^{o\lambda}$ imply the next inequalities for the induced expected losses

$$\begin{aligned}
 & 0 \leq \underbrace{E^{r^{o\lambda}}[L^h] - E^{r^o}[L^h]}_{(20)} \leq \underbrace{E^{r^{o\lambda}}[L^h] + \lambda D(\text{mr}^{o\lambda}||\tilde{c}) - E^{r^o}[L^h]}_{\lambda D(\bullet||\bullet) \geq 0} \\
 & \underbrace{\leq}_{(20)} E^{r^o}[L^h] + \lambda D(\text{mr}^o||\tilde{c}) - E^{r^o}[L^h] = \lambda D(\text{mr}^o||\tilde{c}) \rightarrow_{\lambda \rightarrow 0^+} 0.
 \end{aligned}$$

The assumption ii) in (19) yields $D(\text{mr}^o||\tilde{c}) < \infty$. Thus, the expected loss $E^{r^{o\lambda}}[L^h]$, for the FPD-optimal decision rules $r^{o\lambda}$ (20), is arbitrarily close (for $\lambda \rightarrow 0^+$) to the expected loss $E^{r^o}[L^h]$ with the optimal decision rules r^o (20). \square

Remark 3 (On Generalisation of Proposition 3). The proof of Proposition 3 adds the positive term $\lambda D(\bullet||\bullet)$. λ -multiple of another f-divergence [42], say α -divergence [45,46], could be used to prove that the gained FPD generalisation, see Remark 2, densely extends MDP. The extension will be strict due to the non-linearity of the optimised criterion with respect to the decision rules making the optimal decision rules randomised. \square

3. Main result: discounted FPD

First, the quest for the discounted FPD is shown to be nontrivial. Then, the presentation provides the DM formulation leading to the discounted FPD.

3.1. Unsuccessful attempt

The additive expression of the optimised KL, employing the expectation functional $E^r[\bullet]$ introduced in (6),

$$D(c^r||c^i) = E^r \left[\sum_{t \in \mathcal{I}} \ln \left(\frac{m(s_t|a_t, s_{t-1})r(a_t|s_{t-1})}{m^i(s_t|a_t, s_{t-1})r^i(a_t|s_{t-1})} \right) \right] \tag{21}$$

together with LTF (13), hint a way to the discounted FPD. It seems that it suffices to apply a *partial* LTF to factors forming the ideal closed-loop model and to use the weighted ideal pd

$$\begin{aligned}
 c^{iw}(b) & := \prod_{t \in \mathcal{I}} n^{rw}(s_{t-1})^{-1} [m^i(s_t|a_t, s_{t-1})r^i(a_t|s_{t-1})]^{w(s_{t-1})} \\
 & \quad \times [m(s_t|a_t, s_{t-1})r(a_t|s_{t-1})]^{1-w(s_{t-1})} \\
 n^{rw}(s_{t-1}) & := \int_{(s,a)} [m^i(s_t|a_t, s_{t-1})r^i(a_t|s_{t-1})]^{w(s_{t-1})} \\
 & \quad \times [m(s_t|a_t, s_{t-1})r(a_t|s_{t-1})]^{1-w(s_{t-1})} d(s_t, a_t)
 \end{aligned} \tag{22}$$

with weights $w := (w(s_{t-1}) \in [0, 1], s_{t-1} \in \mathcal{S})_{t \in \mathcal{I}}$. It uses the ideal pds $m^i(b), r^i(b)$ from (2). Inserting (22) into (21) in the place of $c^i(b)$ yields

$$D(c^r||c^{iw}) = E^r \left[\sum_{t \in \mathcal{I}} w(s_{t-1}) \ln \left(\frac{m(s_t|a_t, s_{t-1})r(a_t|s_{t-1})}{m^i(s_t|a_t, s_{t-1})r^i(a_t|s_{t-1})} \right) + \ln(n^{rw}(s_{t-1})) \right]. \tag{23}$$

Above, the first addend after the summation sign has the desired form of the weighted partial loss but the unwanted logarithm of the normalisation $n^{rw}(s_{t-1})$ appears as the second addend. Thus, the same problem as that met in connection with Proposition 2 (17), appears: *careless handling of the ideal pd changes the optimised functional in an undesired way.*

3.2. Successful approach

This part provides the remedy of the above problem. It formulates and solves the relevant FPD.

The adopted approach stems from the recognition that discounting partially *gives up the full optimisation* in order to avoid amplifying errors caused either by an imprecise model [31] and (or) imprecisely quantified aims [29].

Here, this “giving-up the full optimisation” is achieved by introducing a random optional pointer. In each optimisation epoch, it decides whether the selected decision rule should take the current partial loss into account or not. The LTF option allows to neglect the partial loss, see Sec. 2.2. It will be shown that the probability of the employed optimisation acts as the discounting rate. The found optimal solution exhibits the expected properties of DM under discounting.

Let us proceed in the outlined way. Let us introduce optional random pointers $p_t \in \mathbf{p} := \{0, 1\}$, $t \in t$, influencing the factors of the constructed ideal pd $c^i(b)$ (2) in the way leading to the discounted FPD

$$m^i(s_t|a_t, p_t, s_{t-1})r^i(a_t|p_t, s_{t-1}) := \begin{cases} m^i(s_t|a_t, s_{t-1})r^i(a_t|s_{t-1}) & \text{if } p_t = 1 \\ m(s_t|a_t, s_{t-1})r(a_t|p_t, s_{t-1}) & \text{if } p_t = 0 \end{cases} \quad (24)$$

The value $p_t = 1$ says that the t -th addend in (21) is to be optimised and $p_t = 0$ omits the addend from the optimisation.

Pointers $(p_t)_{t \in t}$ extend the behaviour $b \in \mathbf{b}$ to $b := (s_t, a_t, p_t)_{t \in t}$. The closed-loop model of the extended behaviour becomes

$$c^{r^w}(b) := \prod_{t \in t} m(s_t|a_t, s_{t-1})r(a_t|p_t, s_{t-1})w^{p_t}(s_{t-1})(1 - w(s_{t-1}))^{1-p_t}, \text{ with} \\ w(s_{t-1}) := r(p_t = 1|s_{t-1}) \in [0, 1]. \quad (25)$$

The form (25) uses the assumption that s_{t-1} is the state, i.e. the pointer p_t brings no information regarding the state transition

$$m(s_t|a_t, p_t, s_{t-1}) := m(s_t|a_t, s_{t-1}). \quad (26)$$

The assumption (26) of no influence of p_t on s_t , if a_t, s_{t-1} are given, is a version of natural conditions of control [39].

The pointer $p_t \in \{1, 0\}$ should ideally behave according to a chosen ideal pd $w^i(s_{t-1}) := r^i(p_t = 1|s_{t-1}) \in [0, 1]$, $r^i(p_t = 0|s_{t-1}) = 1 - w^i(s_{t-1})$. The ideal pd on the *extended* behaviour $b \in \mathbf{b}$ then gets the form implied by (24)

$$c^{iw^i}(b) := \prod_{t \in t} m^i(s_t|a_t, p_t, s_{t-1})r^i(a_t|p_t, s_{t-1})r^i(p_t|s_{t-1}) = \prod_{t \in t} [m^i(s_t|a_t, s_{t-1}) \\ \times r^i(a_t|s_{t-1})w^i(s_{t-1})]^{p_t} [m(s_t|a_t, s_{t-1})r(a_t|p_t, s_{t-1})(1 - w^i(s_{t-1}))]^{1-p_t}. \quad (27)$$

Using (25), (27) and the definitions of $w(s_{t-1}), w^i(s_{t-1})$, the optimised KL reads

$$D(c^{r^w}||c^{iw^i}) = E^r \left[\sum_{t \in t} \ln \left(\frac{m(s_t|a_t, s_{t-1})r(a_t|p_t, s_{t-1})r(p_t|s_{t-1})}{m^i(s_t|a_t, p_t, s_{t-1})r^i(a_t|p_t, s_{t-1})r^i(p_t|s_{t-1})} \right) \right] \\ = E^r \left[\sum_{t \in t} \ln \left(\frac{m(s_t|a_t, s_{t-1})r(a_t|p_t = 1, s_{t-1})w(s_{t-1})}{m^i(s_t|a_t, s_{t-1})r^i(a_t|s_{t-1})w^i(s_{t-1})} \right) w(s_{t-1}) \right. \\ \left. + \ln \left(\frac{1 - w(s_{t-1})}{1 - w^i(s_{t-1})} \right) (1 - w(s_{t-1})) \right].$$

By separating the weights from the first logarithm the KL of the weights-vector to its ideal counterpart appears. The overall KL gets the discounted form

$$D(c^{r^w}||c^{iw^i}) = E^r \left[\sum_{t \in t} w(s_{t-1}) \ln \left(\frac{m(s_t|a_t, s_{t-1})r(a_t|p_t = 1, s_{t-1})}{m^i(s_t|a_t, s_{t-1})r^i(a_t|s_{t-1})} \right) \right. \\ \left. + D([w(s_{t-1}), 1 - w(s_{t-1})] || [w^i(s_{t-1}), 1 - w^i(s_{t-1})]) \right]. \quad (28)$$

The possible choice $r^i(p_t = 1|s_{t-1}) := w^i(s_{t-1}) := w(s_{t-1}) = r(p_t = 1|s_{t-1})$ mimics LTF introduced in (13). It eliminates KL divergence of the rule generating the pointer $r(p_t|s_{t-1})$ to $r^i(p_t|s_{t-1})$. Then, (28) gives the discounted FPD with the weight $w(s_{t-1}) := r(p_t = 1|s_{t-1}) = w^i(s_{t-1})$.

Generally, the pointer-generating rule $r(p_t|s_{t-1}) := w(s_{t-1})^{p_t}(1 - w(s_{t-1}))^{1-p_t}$ can be optimised, as well. It suffices to apply Proposition 2 for finding the optimal rules $r^o(a_t, p_t|s_{t-1}) = r^o(a_t|p_t, s_{t-1})r^o(p_t|s_{t-1})$ for the extended actions (a_t, p_t) . The optimal weight $w^o(s_{t-1}) = r^o(p_t = 1|s_{t-1})$ is obtained by integrating $r^o(a_t, p_t|s_{t-1})$ over the actions $a_t \in \mathbf{a}$.

Proposition 4 (Main Result: Discounted FPD). *Let the closed-loop model*

$c^{r^w}(b)$ (25) *operate on the extended behaviour* b *made of the preserved states* $s_t \in s$ (26) *and of the extended actions* $(a_t, p_t) \in (\mathbf{a}, \{0, 1\})$, $t \in t$. *Let the ideal pd* $c^{iw^i}(b)$ *have the form* (27). *Then* (28) *holds. It implies that:*

- ▶ *the term* $-E^r[\ln(n^{r^w}(s_{t-1}))]$, *as seen in* (23), *does not appear;*
- ▶ *the extra KL in* (28) *expresses the proximity of* $w(s_{t-1})$ *to* $w^i(s_{t-1})$;

► the optimal rule $r^o(a_t|p_t, s_{t-1})$ and the weight (the pointer-generating rule) $r^o(p_t = 1|s_{t-1}) := w^o(s_{t-1})$ evaluate as follows

$$r^o(a_t|p_t = 1, s_{t-1}) = \frac{r^i(a_t|s_{t-1}) \exp[-d(a_t, s_{t-1})]}{n^i(s_{t-1})}$$

$$n^i(s_{t-1}) := \int_a r^i(a_t|s_{t-1}) \exp[-d(a_t, s_{t-1})] da_t$$

$$d(a_t, s_{t-1}) := \int_s m(s_t|a_t, s_{t-1}) \ln \left[\frac{m(s_t|a_t, s_{t-1})}{n(s_t)m^i(s_t|a_t, s_{t-1})} \right] ds_t.$$

$r^o(a_t|p_t = 0, s_{t-1})$ selects $a^o(s_{t-1}) \in \text{Arg min}_{a_t \in a} K(a_t, s_{t-1})$ with

$$K(a_t, s_{t-1}) := \int_s m(s_t|a_t, s_{t-1}) \ln \left[\frac{1}{n(s_t)} \right] ds_t$$

$$n(s_{t-1}) = n^i(s_{t-1})w^i(s_{t-1}) + k(s_{t-1})(1 - w^i(s_{t-1})),$$

$$k(s_{t-1}) := \exp[-K(a^o(s_{t-1}), s_{t-1})],$$

$$w^o(s_{t-1}) = \frac{n^i(s_{t-1})w^i(s_{t-1})}{n^i(s_{t-1})w^i(s_{t-1}) + k(s_{t-1})(1 - w^i(s_{t-1}))}.$$

This backward recursion, $t = h, h - 1, \dots, 1$, is initiated by $n(s_h) = n^i(s_h) = 1$.⁶

Proof. The optimal decision rules, dealing with extended actions a_t, p_t , have to be found. Proposition 1 provides the optimal rules

$$d(a_t, p_t, s_{t-1}) \stackrel{(11)}{=} \int_s m(s_t|a_t, p_t, s_{t-1}) \ln \left[\frac{m(s_t|a_t, p_t, s_{t-1})}{n(s_t)m^i(s_t|a_t, p_t, s_{t-1})} \right] ds_t$$

$$\stackrel{(24),(26)}{=} p_t \underbrace{\int_s m(s_t|a_t, s_{t-1}) \ln \left[\frac{m(s_t|a_t, s_{t-1})}{n(s_t)m^i(s_t|a_t, s_{t-1})} \right] ds_t}_{:= d(a_t, s_{t-1})}$$

$$+ (1 - p_t) \underbrace{\int_s m(s_t|a_t, s_{t-1}) \ln \left[\frac{1}{n(s_t)} \right] ds_t}_{K(a_t, s_{t-1}) :=} \tag{29}$$

The weight $w^i(s_{t-1})$ influences the value-function $v(s_{t-1}) = -\ln(n(s_{t-1}))$ (11) via

$$n(s_{t-1}) = \int_{(a,p)} r^i(a_t, p_t|s_{t-1}) \exp[-d(a_t, p_t, s_{t-1})] d(a_t, p_t)$$

$$\stackrel{(24),(29)}{=} w^i(s_{t-1}) \underbrace{\int_a r^i(a_t|s_{t-1}) \exp[-d(a_t, s_{t-1})] da_t}_{n^i(s_{t-1}) :=}$$

$$+ (1 - w^i(s_{t-1})) \int_a r^o(a_t|p_t = 0, s_{t-1}) \exp[-K(a_t, s_{t-1})] da_t.$$

LTF applied for $p_t = 0$ implies that the optimal decision rule is deterministic and selects

$$a^o(s_{t-1}) \in \text{Arg min}_{a_t \in a} K(a_t, s_{t-1}) = \text{Arg min}_{a_t \in a} \int_s m(s_t|a_t, s_{t-1}) \ln \left[\frac{1}{n(s_t)} \right] ds_t.$$

The recursion for the function $n(s_{t-1})$ then gets its final form

⁶ K resembles Kerridge's inaccuracy [47] and k is its transformation.

$n(s_{t-1}) := w^i(s_{t-1})n^i(s_{t-1}) + (1 - w^i(s_{t-1}))k(s_{t-1})$ with

$k(s_{t-1}) := \exp[-K(a^o(s_{t-1}), s_{t-1})]$.

The optimal rule generating the extended actions (a_t, p_t) reads

$$\begin{aligned} r^o(a_t, p_t | s_{t-1}) &\stackrel{(12)}{=} \frac{r^i(a_t, p_t | s_{t-1}) \exp[-d(a_t, p_t, s_{t-1})]}{n(s_{t-1})} \\ &= \frac{p_t r^i(a_t | s_{t-1}) \exp[-d(a_t, s_{t-1})] w^i(s_{t-1})}{n(s_{t-1})} \\ &\quad + \frac{(1 - p_t) r^o(a_t | p_t = 0, s_{t-1}) \exp[-K(a_t, s_{t-1})] (1 - w^i(s_{t-1}))}{n(s_{t-1})}. \end{aligned} \quad (30)$$

The second equality uses the chain rule $r^i(a_t, p_t | s_{t-1}) = r^i(a_t | p_t, s_{t-1}) r^i(p_t | s_{t-1})$ and the definition $w^i(s_{t-1}) = r^i(p_t = 1 | s_{t-1})$.

Let us consider $p_t = 1$. The marginalisation of (30) yields

$$\begin{aligned} r^o(p_t = 1 | s_{t-1}) &:= w^o(s_{t-1}) = \frac{n^i(s_{t-1}) w^i(s_{t-1})}{n(s_{t-1})} \\ &= \frac{n^i(s_{t-1}) w^i(s_{t-1})}{n^i(s_{t-1}) w^i(s_{t-1}) + k(s_{t-1}) (1 - w^i(s_{t-1}))}, \quad \text{while} \end{aligned}$$

$r^o(a_t | p_t = 1, s_{t-1}) \propto r^o(a_t, p_t = 1 | s_{t-1})$ implies

$$r^o(a_t | p_t = 1, s_{t-1}) = \frac{r^i(a_t | s_{t-1}) \exp[-d(a_t, s_{t-1})]}{n^i(s_{t-1})}.$$

It remains to find explicitly the deterministic optimal decision rule for $p_t = 0$. A direct check confirms that the equality

$$r^o(a_t | p_t = 0, s_{t-1}) = \frac{r^o(a_t, p_t = 0 | s_{t-1})}{r(p_t = 0 | s_{t-1})}$$

is met for any $a^o(s_{t-1})$ minimising $K(a_t, s_{t-1})$. \square

Remark 4 (On Proposition 4). \blacktriangleright For $w^i(s_{t-1}) = 1$, the last addend in (28) is finite only for

$$w^o(s_{t-1}) = w^i(s_{t-1}) = 1 \text{ as otherwise } D\left([w(s_{t-1}), 1 - w(s_{t-1})] \parallel [1, 0]\right)$$

$= w(s_{t-1}) \ln(w(s_{t-1})) + (1 - w(s_{t-1})) \ln((1 - w(s_{t-1}))/0) = \infty$. Thus, the non-discounted FPD, see Proposition 1, recovers for $w^i(s_{t-1}) = 1$.

\blacktriangleright The optimal rule $r^o(a_t | p_t = 1, s_{t-1})$ coincides with the rule in the non-discounted FPD, see Proposition 1. Discounting influences the value function via $n(s_{t-1})$ as expected. \square

4. Illustrative experiment

The example illustrates that the designed, randomly switching off, the optimisation leads to the discounted FPD that counteracts adverse consequences of imprecise modelling and acts as expected and desirable.

The simulation had three state values $s \in s := \{1, 2, 3\}$ and started from $s_0 = 1$. The environment was stimulated by actions with two values $a \in a := \{1, 2\}$ generated by stationary discounted FPD-optimal decision rules applied up to the horizon $h = 200$. Table 1 describes the used ideal environment model and the ideal decision rule, both time-invariant. The FPD ran with the environment model $m(s_t | a_t, s_{t-1}) \propto \sqrt{m^i(s_t | a_t, s_{t-1})}$. The modelling error was introduced by simulating the environment with the transition pd derived from $m(s_t | a_t, s_{t-1})$ by swapping the pd values with entries $(s_t | a_t, s_{t-1}) = (3|2, 2)$ and $(s_t | a_t, s_{t-1}) = (2|2, 3)$.

Fig. 1 shows the expected influence of the discounting by comparing the cases with $w^i(s_{t-1}) = 0.4$ and $w^i(s_{t-1}) = 1.0$ (no discounting). The strongly preferred state $s_t = 3$ is reached more often with the discounting rate $w^i(s_{t-1}) = 0.4$.

All simulation options are presented to allow a cross-check of our results. Their choice has no deeper meaning. It just fits the illustrative purpose.

The practical choice of the discounting rate is non-trivial and the discounted FPD is expected to help in this respect. It is necessary to conduct a further applied research. Thesis [35] reflects its current state. At present, it is too soon for a meaningful comparative study. It will be done and published elsewhere after combining the derived decision rules with on-line learning [48] and the related quantification of the ideal closed-loop model [21].

Table 1
The Ideal Environment Model and the Ideal Decision Rule.

$m^i(s_t a_t, s_{t-1})$	$a_t = 1$	$a_t = 2$		
$s_t = 1, s_{t-1} = 1$	0.316	0.043		
$s_t = 1, s_{t-1} = 2$	0.307	0.080		
$s_t = 1, s_{t-1} = 3$	0.377	0.016		
$s_t = 2, s_{t-1} = 1$	0.263	0.053		
$s_t = 2, s_{t-1} = 2$	0.373	0.069		
$s_t = 2, s_{t-1} = 3$	0.330	0.027		
$s_t = 3, s_{t-1} = 1$	0.421	0.904		
$s_t = 3, s_{t-1} = 2$	0.320	0.851		
$s_t = 3, s_{t-1} = 3$	0.292	0.957		

$r^i(a_t s_{t-1})$	$a_t = 2$
$s_{t-1} = 1$	0.933
$s_{t-1} = 2$	0.930
$s_{t-1} = 3$	0.957

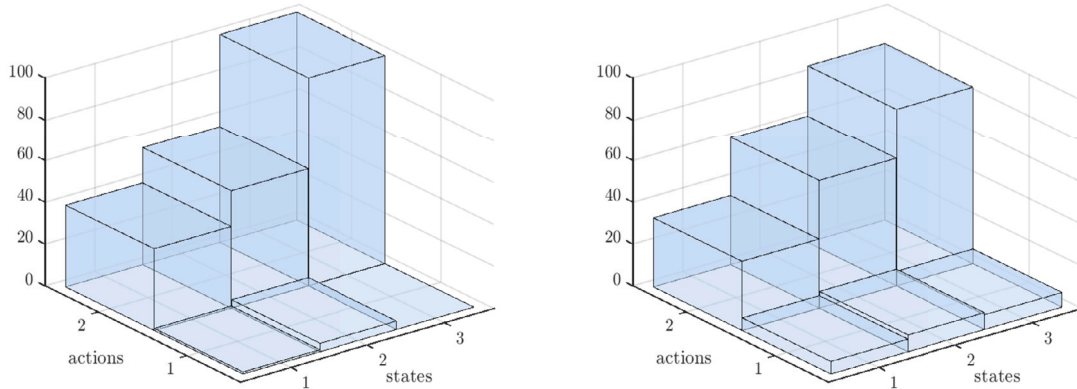


Fig. 1. The demonstration of the discounting influence under mismodelling. The left panel reflects the discounting with $w^i(s_{t-1}) = 0.4$, the right one reflects no discounting, $w^i(s_{t-1}) = 1.0$. Recall: the most desired state is $s = 3$ and the most desired action is $a = 2$.

5. Concluding remarks

The paper pushes further on the support of DM of agents exploiting FPD and thus MDP. It offers them its novel discounted version, Proposition 4. Its usefulness relies on an operational specification of the ideal weights $\{w^i(s_{t-1})\}_{t \in T}$. The unified, pd-based, formulation of FPD, opens a direct way of their choice.

The weights act as forgetting rates in estimation [32]. They can directly be estimated using Bayesian estimation paradigm, at least on discrete grid [33]. The probabilistic interpretation, so useful in preference elicitation [49,50], then allows us to use such estimates as data-dependent choices of the discounting weight.

The running research [35] tries to:

- ▶ relate tightly the discounting to forgetting;
- ▶ admit different weights to the environment models and decision rules;
- ▶ find the conditions under which the designed decision rules converge and stabilise closed loop [13] if the horizon is unbounded.

Authors will be happy if the paper will encourage readers to convert the presented advanced theory into a practically useful tool.

CRedit authorship contribution statement

Miroslav Kárný: Writing – original draft, Investigation. **Soňa Molnárová:** Investigation.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

No data was used for the research described in the article.

References

[1] L. Savage, *Foundations of Statistics*, Wiley, 1954.

- [2] A. Wald, *Statistical Decision Functions*, J. Wiley, 1950.
- [3] E. Feinberg, A. Shwartz, *Handbook of Markov Decision Processes: Methods & Applications*, Kluwer, 2002.
- [4] D. Bertsekas, *Dynamic Programming and Optimal Control*, Athena Scientific, 2017.
- [5] J. van Schuppen, *Stochastic Control Theory*, Springer International Publishing, Cham, 2021, pp. 617–624.
- [6] Y. Yang, G. Hu, C. Spanos, Stochastic optimal control of HVAC system for energy-efficient buildings, *IEEE Trans. Control Syst. Technol.* 30 (1) (2022) 376–383.
- [7] C. Gollier, Discounting an uncertain future, *J. Public Econ.* 85 (2) (2002) 149–166.
- [8] P. Dasgupta, Discounting climate change, *J. Risk Uncertain.* 37 (2008) 141–169.
- [9] A. Attema, W. Brouwer, K. Claxton, Discounting in economic evaluations, *Pharmacoeconomics* 36 (2018) 745–758.
- [10] Y. Gu, Y. Cheng, C. Chen, X. Wang, Proximal policy optimization with policy feedback, *IEEE Trans. Syst. Man Cybern. Syst.* 52 (7) (2022) 4600–4610.
- [11] L. Pan, Y. Deng, K. Cheong, Dynamical Markov decision-making model based on mass function to quantitatively predict interference effects, *Inf. Sci.* 648 (2023) 119482.
- [12] P. Yan, D. Wang, H. Li, D. Liu, Error bound analysis of Q-function for discounted optimal control problems with policy iteration, *IEEE Trans. Syst. Man Cybern. Syst.* 47 (7) (2017) 1207–1216.
- [13] D. Wang, J. Ren, M. Ha, Discounted linear Q-learning control with novel tracking cost and its stability, *Inf. Sci.* 626 (2023) 339–353.
- [14] M. Kárný, Axiomatisation of fully probabilistic design revisited, *Syst. Control Lett.* 141 (2020) 104719.
- [15] M. Kárný, T. Kroupa, Axiomatisation of FPD, *Inf. Sci.* 186 (1) (2012) 105–113.
- [16] I. Landau, A survey of MRAS techniques, *Automatica* 10 (4) (1974) 353–379.
- [17] M. Kárný, Towards fully probabilistic control design, *Automatica* 32 (12) (1996) 1719–1722.
- [18] M. Kárný, T. Guy, Fully probabilistic control design, *SCL* 55 (2006) 259–265.
- [19] A. Quinn, M. Kárný, T. Guy, Fully probabilistic design of hierarchical Bayesian models, *Inf. Sci.* 369 (2016) 532–547.
- [20] E. Garrabé, G. Russo, Probabilistic design of optimal sequential decision-making algorithms in learning and control, *Annu. Rev. Control* 54 (2022) 81–102.
- [21] M. Kárný, T. Siváková, Model-based preference quantification, *Automatica* 156 (2023) 111185.
- [22] S. Kullback, R. Leibler, On information and sufficiency, *Ann. Math. Stat.* 22 (1951) 79–87.
- [23] H. Kappen, Linear theory for control of nonlinear stochastic systems, *Phys. Rev. Lett.* 95 (20) (2005) 200201.
- [24] P. Guan, M. Raginsky, R. Willett, Online Markov decision processes with Kullback Leibler control cost, *IEEE Trans. Automat. Control* 59 (6) (2014) 1423–1438.
- [25] Y. Abbasi-Yadkori, P. Bartlett, X. Chen, A. Malek, Large-Scale Markov Decision Problems with KL Control Cost and Its Application to Crowdsourcing, *Proc. of the 32nd Intern. Conf. on Machine Learning*, vol. 37, JMLR, Lille, France, 2015, pp. 19–28.
- [26] D. Palenicek, A survey on constraining policy updates using the KLD, in: B. Belousov, et al. (Eds.), *Reinforcement Learning Algorithms: Analysis and Applications*, Springer, 2021.
- [27] N. Vieillard, T. Kozuno, B. Scherrer, O. Pietquin, R. Munos, M. Geist, Leverage the average: an analysis of KL regularization in reinforcement learning, *Adv. Neural Inf. Process. Syst.* 33 (2020) 12163–12174.
- [28] K. Rana, V. Dasagi, J. Haviland, B. Talbot, M. Milford, N. Sünderhauf, Bayesian controller fusion: leveraging control priors in deep reinforcement learning for robotics, *Int. J. Robot. Res.* 42 (3) (2023) 123–146.
- [29] J. Doyle, Survey of time preference, delay discounting models, *Judg. Dec. Making* 8 (2013) 116–135.
- [30] V. Gaitsgory, L. Grüne, M. Höger, C. Kellett, S. Weller, Stabilization of strictly dissipative discrete time systems with discounted optimal control, *Automatica* 93 (2018) 311–320.
- [31] C. Rohrs, L. Valavani, M. Athans, G. Stein, Robustness of adaptive control algorithms in the presence of unmodeled dynamics, in: *IEEE Conf. on Decision and Control*, vol. 1, Orlando, FL, 1982, pp. 3–11.
- [32] R. Kulhavý, M.B. Zarrop, On a general concept of forgetting, *Int. J. Control* 58 (4) (1993) 905–924.
- [33] K. Dedecius, I. Nagy, M. Kárný, Parameter tracking with partial forgetting method, *Int. J. Adapt. Control Signal Process.* 26 (1) (2012) 1–12.
- [34] M. Aguayo, L. Bellido, C. Lentisco, E. Pastor, DASH adaptation algorithm based on adaptive forgetting factor estimation, *IEEE Trans. Multimed.* 20 (5) (2018) 1224–1232.
- [35] S. Molnárová, *Applicable Adaptive Discounted Fully Probabilistic Design of Decision Strategy*, FNSPE, CTU, Prague, 2024, <http://hdl.handle.net/10467/114570>.
- [36] M. Rao, *Measure Theory and Integration*, J. Wiley, 1987.
- [37] N. Cammardella, A. Bušić, S. Meyn, Kullback-Leibler-quadratic optimal control in a stochastic environment, in: *60th IEEE Conference on Decision and Control (CDC)*, IEEE, 2021, pp. 158–165.
- [38] D. Gagliardi, G. Russo, On a probabilistic approach to synthesize control policies from example datasets, *Automatica* 137 (2022) 110121.
- [39] V. Peterka, Bayesian system identification, in: P. Eykhoff (Ed.), *Trends & Progress in System Identification*, 1981, pp. 239–304.
- [40] A. Feldbaum, Theory of dual control, *Autom. Remote Control* 22 (1961) 3–19.
- [41] H. Hu, S. Song, G. Huang, Self-attention-based temporary curiosity in reinforcement learning exploration, *IEEE Trans. Syst. Man Cybern. Syst.* 51 (9) (2021) 5773–5784.
- [42] I. Sason, On f-divergences: integral representations, local behavior, and inequalities, *Entropy* 20 (5) (2018) 383.
- [43] M. Kárný, Rényi's extension of fully probabilistic design of decision-making rules, *IEEE Trans. Syst. Man Cybern. Syst.* (2024), submitted for publication.
- [44] M. Kárný, J. Böhm, T. Guy, L. Jirsa, I. Nagy, P. Nedoma, L. Tesaf, *Optimized Bayesian Dynamic Advising: Theory and Algorithms*, Springer, London, UK, 2006.
- [45] S. Amari, α -divergence is unique, belonging to both f-divergence and Bregman divergence classes, *IEEE Trans. Inf. Theory* 55 (11) (2009) 4925–4931.
- [46] T. van Erven, P. Harremoës, Rényi divergence and Kullback-Leibler divergence, *IEEE Trans. Inf. Theory* 60 (7) (2014) 3797–3820.
- [47] D. Kerridge, Inaccuracy and inference, *J. R. Stat. Soc. B* 23 (1961) 284–294.
- [48] M. Kárný, FPD of strategies with estimator, *Automatica* 141 (2022) 110269.
- [49] U. Chajewska, D. Koller, Utilities as random variables: density estimation and structure discovery, in: *Proc. UAI-00*, 2000, pp. 63–71.
- [50] M. Kárný, T. Guy, Preference elicitation within framework of fully probabilistic design of decision strategies, in: *IFAC Workshop ALCOS*, vol. 52, 2019, pp. 239–244.