

Fully Probabilistic Design of Strategies with Estimator[★]

Miroslav Kárný

The Czech Academy of Sciences, Institute of Information Theory and Automation, POB 18, 182 08 Prague 8, Czech Republic

Abstract

The axiomatic fully probabilistic design (FDP) of decision strategies strictly extends Bayesian decision making (DM) theory. FDP also models the closed decision loop by a joint probability density (pd) of all inspected random variables, referred as behaviour. FDP expresses DM aims via an ideal pd of behaviours, unlike the usual DM. Its optimal strategy minimises Kullback-Leibler divergence (KLD) of the joint, strategy-dependent, pd of behaviours to its ideal twin. A range of FDP results confirmed its theoretical and practical strength. Curiously, no guide exists how to select a specific ideal pd for an estimator design. The paper offers it. It advocates the use of the closed-loop state notion and generalises dynamic programming so that FDP is its special case. Primarily, it provides an explorative optimised feedback that “naturally” diminishes exploration (gained in learning) as the learning progresses.

Key words: Bayes methods, Closed loop systems, Decision making, Dynamic programming, Estimation

1 Introduction

On the paper context and main result This brief paper focuses on a technical problem related to a prescriptive theory of dynamic decision making (DM). The theory is dubbed fully probabilistic design (FDP) of decision strategies¹. It generalises methodologies developed in connection with (adaptive) control theory [1,4] and Markov decision processes [32]. Since its initial publication [13], it was broadly elaborated [17], axiomatised [14], applied [16,33] and used for supporting decision makers [9,15,18,42].

The paper deals with the evergreen known as dual control [7,22,28] or exploration-exploitation dichotomy [5]. It concerns the balance of random explorative actions, supporting parameter estimation, with actions moving the closed control loop to the desired state². The main contribution of the paper is an optimised feedback that “naturally” diminishes exploration (gained in learning) as the learning progresses.

On the addressed technical problem Any estimation serves to decision making seen as the aim-focused selection and use of actions. The agent — the decision maker or the action selector, referred as “it” — acts under uncertainty. The

inspected agent uses FDP. FDP models the closed-loop behaviour by the joint probability density (pd). The behaviour consists of all considered uncertain variables. The inspected estimation arises when the behaviour includes a parameter unknown to the agent. Its adopted handling as random variable coincides with bayesianism [3].

The FDP-optimal strategy minimises Kullback-Leibler divergence (KLD) [24] of the behaviours’ pd to its ideal, DM-aims expressing, twin. The estimation has the parameter estimates as (a part of) agent’s actions. The wish to obtain good estimates of the unknown parameter is the generic agent’s aim. The key question is: *What ideal pd expresses this wish?* A universal conversion of a usual loss into the ideal pd exists, Prop. 3. in [14]. It often violates the dictum [20]: *Select an ambitious but reachable ideal pd!* Our solution meets this dictum and leads to the mentioned main result.

Layout Sec. 2 recalls FDP, embeds it into a slightly generalised dynamic programming and advocates the use of *closed-loop* states. Core Sec. 3 solves FDP with an estimator. It proposes the relevant ideal pd and finds the FDP-optimal estimator. Sec. 4 summarises properties of the proposed strategy and outlines open problems.

Notation $\{x\}$ marks the set of x s defined if needed. Sanserif fonts denote mappings. The superscripts ⁱ, ^o refer to the ideal pd and optimality, respectively. The symbol $:=$ defines by assigning; \propto is proportionality; \sim marks interim objects. The time subscript $t \in \{t\}$ of a function f_t on $\{x\}$ drops if the function argument has it, $f(x_t) := f_t(x_t)$. The text prefers mnemonic identifiers.

[★] This paper was not presented at any IFAC meeting.

Email address: school@utia.cas.cz (Miroslav Kárný).

¹ It has overlaps with KL control [6,8,12,30,39], with minimum relative entropy principle [29,35], softmax, etc.

² The exploration counteracts the positive probability of non-optimality connected with certainty-equivalence-based strategies, see e.g. [25] the example on p. 347 and Theorem 7.1

2 Fully Probabilistic Design

FPD deals with the closed DM loop. An agent and its environment form it. The agent's actions $a_t \in \{a\} \neq \emptyset$, at time moments tagged by $t \in \{t\} := \{1, \dots, h\}$, $h < \infty$, influence transitions of states $s_{t-1} \in \{s\} \neq \emptyset$ to states $s_t \in \{s\}$. The inspected transition model depends on an unknown, time-invariant, parameter $p \in \{p\}$. The closed-loop states $(s_t)_{t \in \{t\}}$ are gradually observed and constructed. A fixed, known initial state s_0 implicitly conditions all used pds. The case with internal states is left aside to keep the paper simple. The fact that s_t is *closed-loop* state also simplifies the text. The state s_t includes the value w_t of the sufficient statistic w_t and has the structure

$$s_t := (o_t, w_t) := (\text{observation of the environment,} \quad (1) \\ \text{value of the sufficient statistic } w_t), \quad w_t := w(o_t, a_t, w_{t-1}).$$

The data $d := (d_t)_{t \in \{t\}} := (s_t, a_t)_{t \in \{t\}}$ consists of states, $s_t \in \{s\}$, and actions, $a_t \in \{a\}$, up to the horizon h . The data, $d \in \{d\}$, and the unknown parameter, $p \in \{p\}$, form the random behaviour $b := (d, p) \in \{b\} := (\{d\}, \{p\})$. The agent opts its actions via a randomised causal strategy³

$$s \in \{s\} := \left\{ s(b) := \prod_{t \in \{t\}} r(a_t | s_{t-1}) = s(d) \right\}. \quad (2)$$

Pds $(r(a_t | s_{t-1}))_{t \in \{t\}}$ describe the decision rules of the strategy s . They meet natural conditions of control (NCC) [31]

$$r(a_t | s_{t-1}) = r(a_t | s_{t-1}, p), \quad (3) \\ a_t \in \{a\}, \quad s_{t-1} \in \{s\}, \quad p \in \{p\}.$$

NCC express the adopted assumption that the parameter is unknown to the agent. The s -dependent joint pd $c^s(b) = c^s(d, p)$ fully models random behaviours. The chain rule for pds [31], the state notion and NCC (3) factorise this closed-loop model of the behaviours, $b \in \{b\} = \{b = (d, p)\}$,

$$c^s(b) := c^s(d|p)p(p) := \quad (4) \\ \prod_{t \in \{t\}} m(s_t | a_t, s_{t-1}, p) r(a_t | s_{t-1}) p(p) := m(d|p)s(d)p(p).$$

The parametric models $m(d|p)$, $p \in \{p\}$,

$$m(d|p) \in \{m(d|p)\} := \quad (5) \\ \left\{ m(d|p) := \prod_{t \in \{t\}} m(s_t | a_t, s_{t-1}, p), \quad d = (s_t, a_t)_{t \in \{t\}} \right\}$$

serve the agent to describe its environment. The factors $m(s_t | a_t, s_{t-1}, p)$ are pds of $s_t \in \{s\}$ conditioned on $a_t \in \{a\}$. They model the state transition from $s_{t-1} \in \{s\}$. They

³ A decision function [40], act [34], policy [32], etc.

do it for each $p \in \{p\}$ and all $t \in \{t\}$. The strategy s influences the agent's environment solely via the actions it generates. This explains the dropped superscript ^s in (4). The strategy s has no influence on the prior pd $p(p)$ modelling the unknown parameter $p \in \{p\}$. It drops ^s, too.

The value of the prior pd $p(p)$, $p \in \{p\}$, expresses agent's belief that $m(d|p)$ is the best projection of the objective environment description on the set of parametric models $\{m(d|p), d \in \{d\}\}_{p \in \{p\}}$, [15]. Under NCC (3), the value of the prior pd $p(p)$ also quantifies the belief that $m(d|p)s(d)$ is the best projection of the objective closed-loop model.

FPD quantifies agent's aims and constraints by an ideal closed-loop model c^i . It is a joint pd $c^i(b)$ of ideally distributed behaviours $b \in \{b\}$. The agent chooses the ideal pd $c^i(b)$. It assigns high values to preferred behaviours b , small values to unwanted ones and zero to forbidden bs . The pd $c^i(b)$ factorises as $c^s(b)$, $b \in \{b\}$,

$$c^i(b) := c^i(d|p)p^i(p) := \quad (6) \\ := \prod_{t \in \{t\}} m^i(s_t | a_t, s_{t-1}, p) r^i(a_t | s_{t-1}, p) p^i(p) \\ := m^i(d|p)s^i(d|p)p(p).$$

The factors of the ideal parametric environment model, $m^i(s_t | a_t, s_{t-1}, p)$, $s_t, s_{t-1} \in \{s\}$, $a_t \in \{a\}$, $t \in \{t\}$, describe the *desired* state transitions. The ideal strategy $s^i(d|p) := \prod_{t \in \{t\}} r^i(a_t | s_{t-1}, p)$ consists of the preferred decision rules. Unlike the optimised decision rules (3), they may depend on the unknown $p \in \{p\}$ as they model agent's wishes. The strategy s has no influence on the prior pd, which is left-to-its-fate [16] by setting $p^i(p) := p(p)$. FPD axiomatics [14] implies that the FPD-optimal strategy $s^o \in \{s\}$ minimises the KLD $D(c^s || c^i)$ of c^s to c^i

$$s^o \in \text{Arg min}_{s \in \{s\}} D(c^s || c^i) \\ := \text{Arg min}_{s \in \{s\}} \int_{\{b\}} c^s(b) \ln \left(\frac{c^s(b)}{c^i(b)} \right) db. \quad (7)$$

Remark 1 (On Knowledge Accumulation)

- ✓ Having a parametric model $\tilde{m}(o_t | a_t, s_{t-1}, p)$ relating the observations $o_t \in \{o\}$ (1) to the actions $a_t \in \{a\}$, to the past states $s_{t-1} \in \{s\}$, and to the parameters $p \in \{p\}$, Bayes' rule [31] gradually updates the prior pd $p(p) := p(p|s_0)$ into the posterior pds $p(p|s_t)$

$$p(p|s_t) = \frac{\tilde{m}(o_t | a_t, s_{t-1}, p) p(p|s_{t-1})}{\int_{\{p\}} \tilde{m}(o_t | a_t, s_{t-1}, p) p(p|s_{t-1}) dp} \\ \propto \tilde{m}(o_t | a_t, s_{t-1}, p) p(p|s_{t-1}). \quad (8)$$

This knowledge accumulation is correct under NCC (3) that have the equivalent expression

$$p(p|a_t, s_{t-1}) = p(p|s_{t-1}) \Leftrightarrow r(a_t | s_{t-1}) = r(a_t | s_{t-1}, p) \\ a_t \in \{a\}, \quad s_{t-1} \in \{s\}, \quad p \in \{p\}. \quad (9)$$

✓ Bayes' rule (8) updates the value w_t of the sufficient statistic w_t that is a part of the state (1). Thus, the parametric observation model \tilde{m}_t and Bayes' rule describe the parametric model m_t (5).

The full t th step of the knowledge collection starts with $(p(p|s_{t-1}), s_{t-1})$, uses the action $a_t \sim r(a_t|s_{t-1})$ and the observation o_t made on agent's environment in Bayes' rule (8) giving the statistic value $w_t = w(o_t, a_t, s_{t-1})$. This completes the state $s_t = (o_t, w_t)$, see (1).

✓ The use of the closed-loop state:

- ★ opens a way to a better estimation as the options like forgetting [23] or a trust weight in Bayes' rule [21] just extend the actions;
- ★ unifies inference with tasks influencing the environment;
- ★ allows to model human in DM cycle [38] and to respect agent's attitudes and emotions [11];
- ★ gives a specific view on (ir)rationality of DM [19];
- ★ fits conceptually to dual control [7] and to exploration-exploitation dichotomy [5,37].

3 FPD with Estimator

Part 3.1 constructs the FPD-optimal strategy. It relies on a slight extension of stochastic dynamic programming [4] that minimises⁴ the strategy-dependent expectation $E^s[L^s]$ of the s -dependent additive loss $L^s(b)$, $\{b\}$. The optimal strategy s^o minimises the expectation of the loss

$$L^s(b) := \sum_{t \in \{t\}} l^t(s_{t-1}), \quad l^t : \{s\} \rightarrow (-\infty, \infty)$$

$$s^o \in \text{Arg min}_{\{s\}} E^s[L^s] := \text{Arg min}_{\{s\}} \int_{\{b\}} L^s(b) c^s(b) db. \quad (10)$$

The dependence of the loss L^s on s makes the optimised functional (10) non-linear in the strategy $s \in \{s\}$. This leads to the non-standard optimisation giving the optimal randomised strategy. FPD is a special case of (10). It has an explicit form of the FPD-optimal strategy s^o .

Part 3.2 finds the FPD-optimal strategy s^o with the optimal estimator by using the results of Part 3.1.

3.1 Dynamic Programming on Strategy-Dependent Loss

The next proposition prepares the solution of (10).

Proposition 1 (Static Design) *Let us search for the optimal strategy s^o consisting of single decision rule r^o*

$$s^o := (r^o(a|s))_{a \in \{a\}, s \in \{s\}}$$

$$\in \text{Arg min}_{\{s\}} E^s[L^s] := \text{Arg min}_{\{r\}} \int_{\{s\}} l^r(s) \mu(s) ds,$$

⁴ The assumed generic existence of minima keeps the text simple.

where the marginal pd μ on $\{s\}$ does not depend on r . Let $l^r(s) : \{r\} := \{pds \ r(a|s) \text{ on } \{a\}\} \rightarrow (-\infty, \infty)$ have a minimiser r^* for μ -almost all $s \in \{s\}$.

Then, $s^o = r^o = r^*$. Symbolically,

$$(r^o(a|s))_{a \in \{a\}, s \in \{s\}} \in \text{Arg min}_{\{r(a|s)\}} \int_{\{s\}} l^r(s) \mu(s) ds$$

$$= \text{Arg min}_{\{r(a|s)\}} l^r(s).$$

The minimum is $\int_{\{s\}} l^{r^o}(s) \mu(s) ds$.

Proof By definition, $l^{r^*}(s) \leq l^r(s)$ for any $r \in \{r\}$ and for μ -almost all $s \in \{s\}$. Multiplication of this inequality by the positive values of $\mu(s)$ and integration over $\{s\}$ preserve it. This proves the claim. \square

The solution of (10) uses the optimal value functions, [4], $v_{t-1} : \{s\} \rightarrow (-\infty, \infty)$, $t \in \{t\}$,

$$v(s_{t-1}) := \min_{\{(r_\tau)_{\tau \geq t}\}} E^s \left[\sum_{\tau=t}^h l^{r_\tau}(s_{\tau-1}) \middle| s_{t-1} \right]. \quad (11)$$

Proposition 2 (Dynamic Programming) *The backward, $t = h, \dots, 1$, functional recursion, initiated by $v(s_h) = 0$, evolves the optimal value functions $(v_t)_{t \in \{t\}}$, $s_{t-1} \in \{s\}$,*

$$v(s_{t-1}) = \min_{\{r_t\}} \{l^{r_t}(s_{t-1}) + E^s[v(s_t)|s_{t-1}]\}. \quad (12)$$

The minimising arguments $(r_t^o)_{t \in \{t\}}$ in (12) form the optimal strategy s^o (10). The minimum of (10) is $v(s_0)$.

Proof It exploits the “backward” induction. For $t = h$ and $v(s_h) = 0$, $\forall s_h \in \{s\}$, (12) holds. In a generic step,

$$v(s_{t-1}) \stackrel{(11)}{=} \min_{\{(r_\tau)_{\tau \geq t}\}} E^s \left[\sum_{\tau=t}^h l^{r_\tau}(s_{\tau-1}) \middle| s_{t-1} \right] = \min_{\{r_t\}} \left\{ \begin{aligned} & \min_{\{(r_\tau)_{\tau \geq t+1}\}} \left\{ l^{r_t}(s_{t-1}) + E^s \left[E^s \left[\sum_{\tau=t+1}^h l^{r_\tau}(s_{\tau-1}) \middle| s_t \right] \middle| s_{t-1} \right] \right\} \\ & \stackrel{(11), Prop. 1}{=} \min_{\{r_t\}} \{l^{r_t}(s_{t-1}) + E^s[v(s_t)|s_{t-1}]\}. \end{aligned} \right.$$

The way to the final formula exploits: a) $l^{r_t}(s_{t-1})$ behaves as a constant with respect to conditional expectation $E^s[\bullet|s_{t-1}]$; b) the chain rule for expectation $E^s[\bullet|\diamond] = E^s[E^s[\bullet|\star, \diamond]|\diamond]$, valid for arbitrary \bullet, \star, \diamond ; c) equality $E^s[\bullet|s_t, s_{t-1}] = E^s[\bullet|s_t]$ holds as s_t is the state; d) both l^{r_t} and $c^s(s_t|s_{t-1})$ (used in the outer expectation) are independent of the rules r_τ , $\tau > t$. The final claim of the

proposition follows from the identity $\min_{\{s\}} \mathbb{E}^s[L^s] = v(s_0)$ for the given s_0 , cf. (11), (12). \square

The next proposition provides the FPD-optimal strategy s^o .

Proposition 3 (FPD with an Unknown Parameter) *The backward, $t = h, h-1, \dots, 1$, functional recursion for $n(s_{t-1}) \in [0, 1]$, $d(a_t, s_{t-1}) \geq 0$, $s_{t-1} \in \{s\}$, $a_t \in \{a\}$,*

$$\begin{aligned} d(a_t, s_{t-1}) &:= \int_{\{p\}} p(p|s_{t-1}) \int_{\{s\}} m(s_t|a_t, s_{t-1}, p) \\ &\times \ln \left(\frac{m(s_t|a_t, s_{t-1}, p)}{m^i(s_t|a_t, s_{t-1}, p)r^i(a_t|s_{t-1}, p)n(s_t)} \right) ds_t dp \\ n(s_{t-1}) &:= \int_{\{a\}} \exp[-d(a_t, s_{t-1})] da_t, \quad n(s_h) := 1, \end{aligned} \quad (13)$$

gives the minimum of (7) $v(s_0) := -\ln(n(s_0)) = D(c^{s^o}||c^i)$. The rules r_t^o of the strategy (7) in (2) are

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

Proof The definition of KLD $D(c^s||c^i)$ (7), the factorised forms of the involved pds (4), (6) and NCC (9) imply that KLD is of the additive form (10) $D(c^s||c^i) =$

$$\begin{aligned} \mathbb{E}^s[L^s] &= \mathbb{E}^s \left[\sum_{t \in \{t\}} \int_{\{\tilde{s}\}, \{a\}, \{p\}} m_t(\tilde{s}|a, s, p) r_t(a|s) p_{t-1}(p|s) \right. \\ &\times \ln \left(\frac{m_t(\tilde{s}|a, s, p) r_t(a|s)}{m_t^i(\tilde{s}|a, s, p) r_t^i(a|s, p)} \right) d(\tilde{s}, a, p) \left. \right]. \end{aligned}$$

The summands, denoted $l^t(s)$, are obviously DM-rule dependent. The related value function $v(s_{t-1})$ is non-negative as it is the sum of conditional KLD's. Thus, the transformation $v(s_{t-1}) := -\ln(n(s_{t-1}))$ implies that the function $n(s_{t-1}) \in [0, 1]$, $n(s_h) = 1 \Leftrightarrow v(s_h) = 0$. The minimised right-hand side of (12) reads

$$\begin{aligned} l^t(s) + \mathbb{E}^s[v_t(\tilde{s})|s] &= \int_{\{a\}} r_t(a|s) \left[\ln(r_t(a|s)) \right. \\ &+ \int_{\{p\}} p_{t-1}(p|s) \int_{\{s\}} m_t(\tilde{s}|a, s, p) \\ &\times \ln \left(\frac{m_t(\tilde{s}|a, s, p)}{m_t^i(\tilde{s}|a, s, p) r_t^i(a|s, p) n_t(\tilde{s})} \right) d\tilde{s} dp \left. \right] da \end{aligned}$$

$$\begin{aligned} &= \int_{\{a\}} r_t(a|s) \ln \left(\frac{r_t(a|s)}{\underbrace{\int_{\{a\}} \exp[-d_t(\tilde{a}, s)] d\tilde{a}}_{:= r_t^*(a|s)}} \right) da \\ &\quad - \ln \left(\int_{\{a\}} \exp[-d_t(\tilde{a}, s)] d\tilde{a} \right). \end{aligned}$$

The 1st summand after the last equality is the KLD (conditioned on s) of the decision rules' pair. It reaches its smallest zero value for $r_t = r_t^* \stackrel{\text{Prop.2}}{=} r_t^o$. The 2nd summand, $-\ln(\text{normaliser})$ is the reached minimum. \square

Remark 2 (On Generality and Constraints)

- ✓ The paper [17] deals with the general case containing internal states (time-varying unknown "parameters") but unlike Prop. 3 does not consider the admissible dependence of r^i (6) on them.
- ✓ The support of $m_t^i r_t^i$ includes the support of $m_t r_t^o$ as, otherwise infinite, $d(a_t, s_{t-1})$ (13) guarantees this. This shows that the ideal pd provides also hard constraints on the optimal closed-loop model c^{s^o} .

3.2 FPD with Estimation

DM with estimation deals with an action a_t that has a parameter estimate $\hat{p}_t \in \{p\}$ as its part. It means

$$a_t := (\alpha_t, \hat{p}_t) \in \{a\} := (\{\alpha\}, \{p\}). \quad (14)$$

The possibly empty part $\alpha_t \in \{\alpha\}$ serves to other DM aim than the estimation. This structures the used pds (the explanations of the choices follow their presentation in (15))

$$\begin{aligned} \text{(a): } m(s_t|a_t, s_{t-1}, p) &\stackrel{(14)}{=} m(s_t|\alpha_t, \hat{p}_t, s_{t-1}, p) \\ &:= m(s_t|\alpha_t, s_{t-1}, p) \\ \text{(b): } m^i(s_t|a_t, s_{t-1}, p) &\stackrel{(14)}{=} m^i(s_t|\alpha_t, \hat{p}_t, s_{t-1}, p) \\ &:= m^i(s_t|\alpha_t, s_{t-1}, \hat{p}_t) \\ m^i(s_t|\alpha_t, s_{t-1}, \hat{p}_t) &:= m(s_t|\alpha_t, s_{t-1}, \hat{p}_t) \text{ if } \alpha_t \text{ is free} \\ \text{(c): } r(a_t|s_{t-1}) &\stackrel{(14)}{=} r(\alpha_t|\hat{p}_t, s_{t-1}) r(\hat{p}_t|s_{t-1}) \\ r^i(a_t|s_{t-1}, p) &\stackrel{(14)}{=} r^i(\alpha_t|\hat{p}_t, s_{t-1}, p) r^i(\hat{p}_t|s_{t-1}, p) \\ &:= r^i(\alpha_t|s_{t-1}, \hat{p}_t) \times p(\hat{p}_t|s_{t-1}). \end{aligned} \quad (15)$$

The initial equalities follow from (14). The last equalities in (15) quantify conditions and aims of DM with estimation. They respectively express:

- (a): independence of the state of the parameter estimate: the environment does not care about agent's estimate and it also does not influence the sufficient statistic;

- (b): the wish to choose \hat{p}_t that ideally replaces the unknown p in $m^i(s_t|\alpha_t, s_{t-1}, p)$ expressing the α -driving aim; $m^i := m$ when α_t is not optimised (left-to-its-fate [16]);
- (c): the estimate \hat{p}_t ideally replaces the unknown $p \in \{p\}$ also in r^i and in $p(p|s_{t-1})$ collecting the knowledge on p .

Proposition 4 (FPD with Estimator) *The FPD-optimal strategy with an estimator, given by options (14), (15), results from the backward, $t = h, \dots, 1$, recursion for $n(s_{t-1}) \in [0, 1]$ and $\Delta(\alpha_t, \hat{p}_t, s_{t-1}) \geq 0$, $s_{t-1} \in \{s\}$, $\alpha_t \in \{\alpha\}$, $\hat{p}_t \in \{p\}$, $t \in \{t\}$, with $n(s_h) := 1$,*

$$\begin{aligned} \Delta(\alpha_t, \hat{p}_t, s_{t-1}) &:= \int_{\{p\}} p(p|s_{t-1}) \\ &\times \int_{\{s\}} m(s_t|\alpha_t, s_{t-1}, p) \ln \left[\frac{m(s_t|\alpha_t, s_{t-1}, p)}{m^i(s_t|\alpha_t, s_{t-1}, \hat{p}_t)n(s_t)} \right] ds_t dp \\ n(s_{t-1}) &:= \int_{\{\alpha\}, \{s\}} r^i(\alpha_t|s_{t-1}, \hat{p}_t) p(\hat{p}_t|s_{t-1}) \\ &\times \exp[-\Delta(\alpha_t, \hat{p}_t, s_{t-1})] d(\hat{p}_t, \alpha_t). \end{aligned} \quad (16)$$

It gives the minimum $v(s_0) := -\ln(n(s_0)) = D(c^{r^o}||c^i)$. The rules r^o of the optimal strategy are

$$r^o(\alpha_t, \hat{p}_t|s_{t-1}) \propto r^i(\alpha_t|s_{t-1}, \hat{p}_t) p(\hat{p}_t|s_{t-1}) \exp[-\Delta(\alpha_t, \hat{p}_t, s_{t-1})]. \quad (17)$$

Proof It uses Prop. 3 specialised by (14), (15). The exponent $d(a_t, s_{t-1}) \stackrel{(14)}{=} d(\alpha_t, \hat{p}_t, s_{t-1})$ in (13) under (15) gets the form, with $n(s_h) = 1$ and $\Delta(\alpha_t, \hat{p}_t, s_{t-1})$ defined in (16),

$$\begin{aligned} d(\alpha_t, \hat{p}_t, s_{t-1}) &:= \int_{\{p\}} p(p|s_{t-1}) \int_{\{s\}} m(s_t|\alpha_t, s_{t-1}, p) \times \\ &\ln \left(\frac{m(s_t|\alpha_t, s_{t-1}, p)}{m^i(s_t|\alpha_t, s_{t-1}, \hat{p}_t) r^i(\alpha_t|s_{t-1}, \hat{p}_t) p(\hat{p}_t|s_{t-1}) n(s_t)} \right) ds_t dp \\ &= -\ln \left(r^i(\alpha_t|s_{t-1}, \hat{p}_t) p(\hat{p}_t|s_{t-1}) \times \right. \\ &\quad \left. \exp \left[\underbrace{- \int_{\{p\}} p(p|s_{t-1}) \int_{\{s\}} m(s_t|\alpha_t, s_{t-1}, p) \right.}_{:= -\Delta(\alpha_t, \hat{p}_t, s_{t-1})} \right. \right. \\ &\quad \left. \left. \times \ln \left(\frac{m(s_t|\alpha_t, s_{t-1}, p)}{m^i(s_t|\alpha_t, s_{t-1}, \hat{p}_t) n(s_t)} \right) ds_t dp \right] \right) \\ n(s_{t-1}) &:= \int_{\{\alpha\}, \{s\}} r^i(\alpha_t|s_{t-1}, \hat{p}_t) p(\hat{p}_t|s_{t-1}) \\ &\quad \times \exp[-\Delta(\alpha_t, \hat{p}_t, s_{t-1})] d(\alpha_t, \hat{p}_t) \\ r^o(\alpha_t, \hat{p}_t|s_{t-1}) &= r^o(\alpha_t|s_{t-1}, \hat{p}_t) r^o(\hat{p}_t|s_{t-1}) \\ &\propto r^i(\alpha_t|s_{t-1}, \hat{p}_t) p(\hat{p}_t|s_{t-1}) \exp[-\Delta(\alpha_t, \hat{p}_t, s_{t-1})]. \quad \square \end{aligned}$$

Remark 3 (Complexity and Propositions 3 & 4)

- ✓ *The dynamic programming generally needs approximations [36]. This fact is behind need for model-based predictive control [26]. The proposed strategy has a potential to enhance such approximate strategies as it supports exploration (dual control [7]) in a novel powerful way.*

- ✓ *Propositions 3 and 4 deal with an unknown parameter but only Prop. 4 formulates its estimation as a part of the DM aim. It allows to solve sole estimation task. Primarily, it makes the optimal strategy more amenable to approximation as seen from the discussion in Sec. 4.*
- ✓ *The posterior pd $p(\hat{p}_t|s_{t-1})$ influences complexity of any strategy with estimation. This implies a wide use of models from the exponential family (EF) [2] for which the functional Bayes' rule reduces to algebraic operations.*
- ✓ *Linear Gaussian autoregressive-regressive model is a prominent member of EF. Its estimation reduces to recursive least squares [31]. Markov decision process [32] with discrete-valued states and actions is the other EF member with the estimation that simply counts the number of observed configurations (s_t, a_t, s_{t-1}) . Many extensions rely on these EF members. Their mixtures [16] are typical and serve as universal approximators [27].*

4 On the Proposed Strategy

The novel FPD-optimal strategy with estimator:

- ✓ respects both the knowledge collected in the posterior pd p_{t-1} (8) and influence of the parameter estimate \hat{p} on α -driven DM via the function $\Delta_t(\alpha, \hat{p}, s)$ (16), which is the expected (weighted) divergence of the environment model $m_t(\tilde{s}|\alpha, s, p)$ to its ideal twin $m_t^i(\tilde{s}|\alpha, s, \hat{p})$;
 - ✓ correlates, due to the previous property, usual actions α_t with estimates \hat{p}_t more deeply⁵ than the common ‘‘certainty equivalence’’ that uses an estimate \hat{p}_t based only on the past knowledge accumulated in the posterior pd p_{t-1} and unrelated with the DM dynamics;
 - ✓ shows that the estimation is dynamic DM task influenced by the value function even when the ‘‘ordinary’’ action α_t is un-optimised.
- Moreover, Prop. 4 opens ways to:
- ✓ the combination of the classical estimation with sequential stopping [41] with a dynamic DM;
 - ✓ the radically novel approach to feasible exploration: non-explorative, certainty-equivalent, α_t -driven DM will be stimulated by random sampling of the parameter estimate \hat{p}_t from the optimal rule (17): it is influenced by the posterior pd $p(\hat{p}_t|s_{t-1})$ (8), which ‘‘naturally’’ shrinks.

The presented result reflects a part of an open-ended research, which surely requires to deal with:

- ✓ the analysis, almost surely simulation-based, of the proposed strategy and its model-based predictive variants;
- ✓ the opened research ways outlined above;
- ✓ specific cases and important real-life problems;
- ✓ the potential indicated in connection with the use of closed-loop states;
- ✓ the general FPD with unobserved states in the ideal pd.

⁵ The paper [10], indicates, within a quite different set up, how useful such a dependence can be.

The foreseeable positive consequences make this research worthwhile. You are invited to join it.

Acknowledgements

MŠMT ČR LTC18075 and EU-COST Action CA16228 support this research.

References

- [1] K.J. Åström and B. Wittenmark. *Adaptive Control*. Addison-Wesley, 1994.
- [2] O. Barndorff-Nielsen. *Information and Exponential Families in Statistical Theory*. Wiley, N.Y., 1978.
- [3] J.O. Berger. *Statistical Decision Theory and Bayesian Analysis*. Springer, 1985.
- [4] D.P. Bertsekas. *Dynamic Programming and Optimal Control*. Athena Scientific, 2017.
- [5] O. Besbes, Y. Gur, and A. Zeevi. Optimal exploration – exploitation in a multi-armed bandit problem with nonstationary rewards. *Stochastic Systems*, 9(4):319–337, 2019.
- [6] A. Bušič and S. Meyn. Action-constrained Markov decision processes with Kullback-Leibler cost. In S. Bubeck, V. Perchet, and P. Rigollet, editors, *Proc. of Machine Learning Research*, volume 75, pages 1–14. MLR Press, 2018.
- [7] A.A. Feldbaum. Theory of dual control. *Autom. Remote Control*, 22:3–19, 1961.
- [8] P. Guan, M. Raginsky, and R. Willett. Online Markov decision processes with Kullback-Leibler control cost. In *Am. Control Conference*, pages 1388–1393. IEEE, June 2012.
- [9] T.V. Guy, S. Fakhimi Derakhshan, and J. Štěch. Lazy fully probabilistic design: Application potential. In F. Belardinelli, editor, *Multi-Agent Systems and Agreement Technologies*. Springer, 2018.
- [10] T.A.N. Heirung, B.E. Ydstie, and B. Foss. Dual adaptive model predictive control. *Automatica*, 80:340–348, 2017.
- [11] J. Homolová, A. Černěcká, T.V. Guy, and M. Kárný. Affective decision-making in ultimatum game: Responder. In M. Slavkovik, editor, *Multi-Agent Systems, EUMAS 2018*, volume LNAI 11450, pages 127–139. Springer Nature Switzerland AG, 2019.
- [12] H.J. Kappen. Linear theory for control of nonlinear stochastic systems. *Physical review letters*, 95(20):200201, 2005.
- [13] M. Kárný. Towards fully probabilistic control design. *Automatica*, 32(12):1719–1722, 1996.
- [14] M. Kárný. Axiomatisation of fully probabilistic design revisited. *SCL*, 104719, 2020.
- [15] M. Kárný. Towards on-line tuning of adaptive-agent’s multivariate meta-parameter. *Int. J. of Machine Learning and Cybernetics*, 12(9):2717–2731, 2021.
- [16] M. Kárný, J. Böhm, T.V. Guy, L. Jirsa, I. Nagy, P. Nedoma, and L. Tesaf. *Optimized Bayesian Dynamic Advising: Theory and Algorithms*. Springer, London, UK, 2006.
- [17] M. Kárný and T.V. Guy. Fully probabilistic control design. *SCL*, 55:259–265, 2006.
- [18] M. Kárný and T.V. Guy. On support of imperfect Bayesian participants. In T.V. Guy and et al, editors, *Decision Making with Imperfect Decision Makers*, volume 28, pages 29–56. Springer Int. Syst. Ref. Lib., 2012.
- [19] M. Kárný and T.V. Guy. On the origins of imperfection and apparent non-rationality. In T.V. Guy and et al, editors, *Decision Making: Uncertainty, Imperfection, Deliberation and Scalability*, Studies in Comput. Intelligence, pages 57–91. Springer, Switzerland, 2015.
- [20] M. Kárný and T.V. Guy. Preference elicitation within framework of fully probabilistic design of decision strategies. In *IFAC Int. Workshop on Adaptive and Learning Control Systems*, volume 52, pages 239–244, 2019.
- [21] M. Kárný and F. Hůla. Fusion of probabilistic unreliable indirect information into estimation serving to decision making. *Int. Journal of Machine Learning and Cybernetics*, 12:33673378, 2021.
- [22] E.D. Klenke and P. Hennig. Dual control for approximate Bayesian reinforcement learning. *JMLR*, 17:1–30, 2016.
- [23] R. Kulhavý and M. B. Zarrop. On a general concept of forgetting. *Int. J. of Control*, 58(4):905–924, 1993.
- [24] S. Kullback and R. Leibler. On information and sufficiency. *Ann Math Stat*, 22:79–87, 1951.
- [25] P.R. Kumar. A survey on some results in stochastic adaptive control. *SIAM J. Control and Applications*, 23:399–409, 1985.
- [26] D.Q. Mayne. Model predictive control: Recent developments and future promise. *Automatica*, pages 2967–2986, 2014.
- [27] G.J. McLachlan and D. Peel. *Finite Mixture Models*. Wiley Series in Probab. & Stat. Wiley, N.Y., 2000.
- [28] Ali Mesbah. Stochastic model predictive control with active uncertainty learning: A survey on dual control. *Annual Reviews in Control*, 45:107 – 117, 2018.
- [29] P.A. Ortega and D.A. Braun. A minimum relative entropy principle for learning and acting. *Journal of Artificial Intelligence Research*, 38:475–511, 2010.
- [30] D. Palenicek. A survey on constraining policy updates using the KL divergence. In B. Belousov, H. Abdulsamad, P. Klink, S. Parisi, and J. Peters, editors, *Reinforcement Learning Algorithms: Analysis and Applications*. Springer, Cham, 2021.
- [31] V. Peterka. Bayesian system identification. In P. Eykhoff, editor, *Trends & Progress in System Identification*, pages 239–304. Perg. Press, 1981.
- [32] M.L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley, 2005.
- [33] A. Quinn, P. Ettlér, L. Jirsa, I. Nagy, and P. Nedoma. Probabilistic advisory systems for data-intensive applications. *Int. J. Adapt Control Signal Process.*, 17(2):133–148, 2003.
- [34] L.J. Savage. *Foundations of Statistics*. Wiley, 1954.
- [35] J. Shore and R. Johnson. Axiomatic derivation of the principle of maximum entropy & the principle of minimum cross-entropy. *IEEE Tran. on Inf. Th.*, 26(1):26–37, 1980.
- [36] J. Si, A.G. Barto, W.B. Powell, and D. Wunsch, editors. *Handbook of Learning and Approximate Dynamic Programming*. Wiley, 2004.
- [37] H. Tang and et al. #Exploration: A study of count-based exploration for deep reinforcement learning. In I. Guyon et al, editor, *Adv. in Neur. Inf. Proc.* 30, pages 2753–2762. Curran Associates, Inc., 2017.
- [38] N. Tishby and D. Polani. Information theory of decisions and actions. In V. Cutsuridis and et al, editors, *Perception-Action Cycle*, pages 601–636. Springer, N.Y., 2011.
- [39] E. Todorov. Linearly-solvable Markov decision problems. In B. Schölkopf and et al, editors, *Adv. in Neur. Inf. Proc.*, pages 1369 – 1376. MIT Press, 2006.
- [40] A. Wald. *Statistical Decision Functions*. J. Wiley, 1950.
- [41] A. Wald. *Sequential Analysis*. Dover Publications, 2013.
- [42] E. Zugarová and T.V. Guy. Similarity-based transfer learning of decision policies. In *Proc. of the IEEE International Conference on Systems, Man, and Cybernetics*, pages 37–44. IEEE, 2020.