



Akademie věd České republiky  
Ústav teorie informace a automatizace, v.v.i.

Academy of Sciences of the Czech Republic  
Institute of Information Theory and Automation

## RESEARCH REPORT

MIROSLAV KÁRNÝ, FRANTIŠEK HŮLA

### **Balancing Exploitation and Exploration via Fully Probabilistic Design of Decision Policies**

No. 2376

October 2018

ÚTIA AVČR, v.v.i., P.O.Box 18, 182 08 Prague, Czech Republic

Fax: (+420)286890378

<http://www.utia.cas.cz>

E-mail: [utia@utia.cas.cz](mailto:utia@utia.cas.cz)

This report is an unrefereed manuscript which is intended to be submitted for publication. Any opinions and conclusions expressed in this report are those of the authors and do not necessarily represent the views of the institute.

This research has been supported by GAČR, grants GA16-09848S and 18-15970S.

**Abstract.** Adaptive decision making learns an environment model serving a design of a decision policy. The policy-generated actions influence both the acquired reward and the future knowledge. The optimal policy properly balances exploitation with exploration. The inherent dimensionality curse of decision making under incomplete knowledge prevents the realisation of the optimal design. This has stimulated repetitive attempts to reach this balance at least approximately. Usually, either: (a) the exploitative reward is enriched by a part reflecting the exploration quality and a feasible approximate certainty-equivalent design is made; or (b) an explorative random noise is added to the purely exploitative actions. This paper avoids the inauspicious (a) and improves (b) by employing the non-standard fully probabilistic design (FPD) of decision policies, which naturally generates random actions. Monte-Carlo experiments confirm its achieved quality. The quality stems from methodological contributions, which include: (i) an improvement of the relation between FPD and standard Markov Decision Processes; (ii) a design of an adaptive tuning of an FPD-parameter. The latter also suits for the tuning of the temperature in both simulated annealing and Boltzmann’s machine.

**Keywords:** Exploitation · Exploration · Bayesian estimation · Adaptive systems · Fully probabilistic design · Kullback-Leibler divergence · Decision policy · Markov decision process.

## 1 Introduction

The inspected decision making is close to the traditional Markov Decision Process (MDP, [23]). The next summary of known basic facts allows us to formulate and solve the addressed problem. In order to focus on the paper’s topic, we restrict ourselves to a finite amount of possible agent’s actions<sup>1</sup>  $a_t \in \mathbf{a} = \{1, \dots, |\mathbf{a}|\}$ ,  $|\mathbf{a}| < \infty$ . They are selected in a finite amount of epochs  $t \in \mathbf{t} = \{1, \dots, |\mathbf{t}|\}$ ,  $|\mathbf{t}| < \infty$ . The agent’s environment responds to actions by discrete-valued observable states  $s_t \in \mathbf{s} = \{1, \dots, |\mathbf{s}|\}$ ,  $|\mathbf{s}| < \infty$ . A given real-valued reward  $r = (r_t(\tilde{s}, a, s), \tilde{s}, s \in \mathbf{s}, a \in \mathbf{a})_{t \in \mathbf{t}}$  quantifies the agent’s preferences. The sequence of transition probabilities

$$\mathbf{p} = (\mathbf{p}_t(\tilde{s}|a, s), \tilde{s}, s \in \mathbf{s}, a \in \mathbf{a})_{t \in \mathbf{t}}, \quad (1)$$

models the assumed Markov random environment. The sequence of probabilities  $\pi = (\pi_t(a|s), a \in \mathbf{a}, s \in \mathbf{s})_{t \in \mathbf{t}}$  describes the agent’s optional, randomised and

<sup>1</sup> Throughout,  $\mathbf{x}$  denotes set of  $x$ ’s of cardinality  $|\mathbf{x}|$ .

Markov policy. The MDP-optimal policy  $\pi^{\text{MDP}}$  maximises the expected cumulative reward

$$\pi^{\text{MDP}} \in \text{Arg max}_{\pi \in \boldsymbol{\pi}} \mathbf{E}^\pi \left[ \sum_{t \in \boldsymbol{t}} r_t(s_t, a_t, s_{t-1}) \right]. \quad (2)$$

The strategy-dependent expectation  $\mathbf{E}^\pi$  is implicitly conditioned on a known initial state. The optimisation runs over the set  $\boldsymbol{\pi}$  of Markov policies

$$\boldsymbol{\pi} = \left\{ \left( \pi_t(a|s) \geq 0, \sum_{a \in \boldsymbol{a}} \pi_t(a|s) = 1, \forall s \in \boldsymbol{s} \right)_{t \in \boldsymbol{t}} \right\}. \quad (3)$$

Dynamic programming (DP) provides the MDP-optimal policy consisting of *deterministic* decision rules  $(\pi_t(a|s))_{t \in \boldsymbol{t}}$  selecting the maximisers  $a_t^{\text{MDP}}(s)$  in

$$\mathbf{v}_{t-1}(s) = \max_{a \in \boldsymbol{a}} \mathbf{E}^\pi [r_t(\tilde{s}, a, s) + \mathbf{v}_t(\tilde{s})|a, s], \quad s \in \boldsymbol{s}, t \in \boldsymbol{t}. \quad (4)$$

The functional equation (4) evolves the value functions  $\mathbf{v}_t(s)$ ,  $s \in \boldsymbol{s}$ , and provides the used maximising arguments  $a_t^{\text{MDP}}(s)$ ,  $s \in \boldsymbol{s}$ . It is solved backwards starting with  $\mathbf{v}_{|\boldsymbol{t}|}(s) = 0$ ,  $\forall s \in \boldsymbol{s}$ . This standard solution extends to the case with the incompletely known environment model parameterised by the transition probability values

$$\mathbf{p}_t(\tilde{s}|a, s, \boldsymbol{\theta}) = \boldsymbol{\theta}(\tilde{s}|a, s), \quad \boldsymbol{\theta} \in \boldsymbol{\theta}. \quad (5)$$

The set  $\boldsymbol{\theta}$  is delimited by the meaning of the parameter

$$\boldsymbol{\theta} = \left\{ \theta(\tilde{s}|a, s) \geq 0, \sum_{\tilde{s} \in \boldsymbol{s}} \theta(\tilde{s}|a, s) = 1, \forall a \in \boldsymbol{a}, \forall s \in \boldsymbol{s} \right\}. \quad (6)$$

The parametric model (5) belongs to exponential family [3] and possesses Dirichlet's distribution  $\mathcal{D}_\theta(V_0)$ , given by the finite-dimensional occurrence array

$$V_0 = (V_0(\tilde{s}|a, s))_{\tilde{s}, s \in \boldsymbol{s}, a \in \boldsymbol{a}}, \quad V_0(\tilde{s}|a, s) > 0,$$

as its conjugate prior. With the chosen  $\mathcal{D}_\theta(V_0)$ , Bayesian learning [5] reproduces Dirichlet's form. It reduces to the updating of the occurrence array

$$V_t(s_t|a_t, s_{t-1}) = V_{t-1}(s_t|a_t, s_{t-1}) + 1, \quad (7)$$

where  $(s_t, a_t, s_{t-1})$  is the realised triple. This recursion, together with the predictive probabilities

$$\mathbf{p}(\tilde{s}|a, s, V_{t-1}) = \frac{V_{t-1}(\tilde{s}|a, s)}{\sum_{\tilde{s} \in \boldsymbol{s}} V_{t-1}(\tilde{s}|a, s)} = \hat{\boldsymbol{\theta}}_{t-1}(\tilde{s}|a, s), \quad \forall \tilde{s}, s \in \boldsymbol{s}, \forall a \in \boldsymbol{a}, \quad (8)$$

provides the Markov transition probability of the *information state*  $(s_t, V_t)$ . Thus, the MDP-optimal policy can *formally* be computed via DP (4) where  $s$  is replaced by  $(s, V)$ . Such an MDP-optimal policy  $(\pi_t(a|s, V))_{t \in \boldsymbol{t}}$  inevitably optimally balances the explorative effort, regarding the evolution of  $s_t$ , and the exploitative effort, regarding the evolution of  $V_t$ , cf. [11]. The number of possible

information states however, blows up exponentially. This prevents the evaluation and storing of the value functions  $(v_t(s, V))_{t \in \mathcal{T}}$ .

The common remedy is the use of the frozen point estimate  $\hat{\theta}_{t-1}$  instead of  $\theta$  in DP. This certainty-equivalent approximation diminishes the curse of dimensionality [4]. This approximation, however, gives up the care about the intentional exploration. It provably diverges from the optimal policy with a positive probability [21]. This experimentally well-confirmed fact has led to a range of attempts to recover the intentional exploration. The active exploration is mostly reached by introducing a random constituent into actions [9, 10, 31]. Good results are often achieved but the proper balance between exploration and exploitation is hard to find. This manifests itself in, repeatedly admitted, sensitivity to the choice of parameters determining the noise added to the exploitative actions.

This paper introduces the proper exploration by employing the fully probabilistic design of decision policies (FPD, [16, 18]). FPD is closely related to the Kullback-Leibler control [12, 13, 15]. In the given context, it is important that FPD leads to the *randomised*, and thus *explorative policy* unlike the usual MDP.

Methodologically, the paper relates MDP and FPD in a better way than the axiomatisation [18]. It proposes the adaptation of an optional FPD-parameter, similar to the temperature in simulated annealing [26] or Boltzmann's machine [30]. Practically, it presents Monte Carlo experiments, which show that the certainty-equivalent version of FPD is indeed adequately explorative.

*Layout:* Section 2 recalls basic facts about the ingredients of the advocated decision policy. It formalises and solves the addressed problem. Section 3 summarises the results of extensive simulations reflecting the properties of the proposed policy. Section 4 adds concluding remarks. Appendix contains data used in simulations so that our results can be reproduced.

## 2 FPD and its Relation to MDP

The environment model  $\mathfrak{p}$  (1) and any fixed policy  $\pi$  in (3) determine the joint probability  $\mathbf{c}^\pi$  of states and actions (implicitly conditioned on the initial state)

$$\mathbf{c}^\pi(\overbrace{s_{|t|}, a_{|t|}, s_{|t|-1}, a_{|t|-1} \dots, s_1, a_1}^{\text{behaviour } b \in \mathbf{b} = \mathbf{X}_{t \in \mathcal{T}}(\mathbf{s} \times \mathbf{a})}) = \prod_{t \in \mathcal{T}} p_t(s_t | a_t, s_{t-1}) \pi_t(a_t | s_{t-1}). \quad (9)$$

This closed-loop model  $\mathbf{c}^\pi(b)$  *completely* describes (closed-loop) behaviours  $b \in \mathbf{b}$  (9) consisting of observed and opted variables. Thus, all design ways, e.g. MDPs with different rewards, leading to the same  $\mathbf{c}^\pi$  are equivalent. This observation [27] implies that decision objectives can generally be expressed via an ideal (desired) closed-loop model  $\mathbf{c}^i(b)$ ,  $b \in \mathbf{b}$ . Informally, the ideal assigns high values to desired behaviours and small values to undesired behaviours. With the ideal closed-loop model chosen, the FPD-optimal policy  $\pi^{\text{FPD}}$  makes  $\mathbf{c}^{\pi^{\text{FPD}}}$  closest to  $\mathbf{c}^i$ . The FPD axiomatisation [18] specifies widely-acceptable conditions under

which the Kullback-Leibler divergence  $D(\mathbf{c}^\pi \|\mathbf{c}^i)$ , [20], is the adequate proximity measure. The FPD-optimal policy  $\pi^{\text{FPD}}$  is thus

$$\pi^{\text{FPD}} \in \text{Arg min}_{\pi \in \boldsymbol{\pi}} D(\mathbf{c}^\pi \|\mathbf{c}^i) = \text{Arg min}_{\pi \in \boldsymbol{\pi}} \sum_{b \in \mathbf{b}} c^\pi(b) \ln \left( \frac{c^\pi(b)}{c^i(b)} \right). \quad (10)$$

Proposition 1 presented below describes the FPD-optimal decision rules. The proposition is a direct counterpart of stochastic DP [1, 7]. It uses the chain-rule factorisation of  $\mathbf{c}^i$ , which delimits: (a) the ideal environment model  $\mathbf{p}_t^i(\tilde{s}|a, s)$ , which are the ideal counterparts of the transition probabilities  $\mathbf{p}_t(\tilde{s}|a, s)$ , and (b) the ideal decision rules  $\pi_t^i(a|s)$  forming the ideal policy.

Proof of Proposition 1 is, for instance, in [29]. The general FPD with the state estimation, corresponding to the partially observable MDP, is in [16].

**Proposition 1 (FPD-Optimal Policy)** *Decision rules  $\pi_t^{\text{FPD}}(a|s)$ ,  $t \in \mathbf{t}$ , forming the FPD-optimal policy (10) result from the following backward recursion, initiated by  $w_{|t}(s) = 1$ ,  $\forall s \in \mathbf{s}$ ,*

$$\pi_t^{\text{FPD}}(a|s) = \frac{\pi_t^i(a|s) \exp[-\omega_t(a, s)]}{\underbrace{\sum_{a \in \mathbf{a}} \pi_t^i(a|s) \exp[-\omega_t(a, s)]}_{w_{t-1}(s)}}, \quad a \in \mathbf{a}, \quad s \in \mathbf{s}, \quad (11)$$

$$\omega_t(a, s) = \sum_{\tilde{s} \in \mathbf{s}} \mathbf{p}_t(\tilde{s}|a, s) \ln \left( \frac{\mathbf{p}_t(\tilde{s}|a, s)}{\mathbf{p}_t^i(\tilde{s}|a, s) w_t(\tilde{s})} \right).$$

The work [18] containing axiomatisation of FPD also proved that: (i) any Bayesian decision making can be arbitrarily well approximated by the FPD formulation (10); (ii) there are FPD tasks having no standard counterpart. In other words, FPD tasks represent the proper dense extension of Bayesian decision making. Here, we modify the constructive way in which this result was shown. The construction explicitly relates the standard MDP to the less usual FPD. Importantly, it serves the purpose of this paper. It shows how the MDP-optimal deterministic policy is arbitrarily-well approximated by the naturally explorative, FPD-optimal, randomised policy. The construction uses the standard notion of entropy  $H^\pi$  [8] of the closed-loop model  $\mathbf{c}^\pi$  and the given cumulative reward  $\mathbf{R}$

$$H^\pi = - \sum_{b \in \mathbf{b}} c^\pi(b) \ln(c^\pi(b)), \quad \mathbf{R}(b) = \sum_{t \in \mathbf{t}} r_t(s_t, a_t, s_{t-1}), \quad b \in \mathbf{b}. \quad (12)$$

**Proposition 2 (FPD from MDP)** *The optimisation (2) over policies  $\pi \in \boldsymbol{\pi}$  (3), restricted by the additional requirement, determined by an optional  $h > 0$ ,*

$$H^\pi \geq h > H^{\pi^{\text{MDP}}}, \quad (13)$$

*leads to the FPD-optimal policy (10) with respect to the ideal closed-loop model<sup>2</sup>*

$$c^i(b) \propto \exp[\mathbf{R}(b)/\lambda]. \quad (14)$$

<sup>2</sup>  $\propto$  means proportionality.

The corresponding ideal environment model and the ideal decision rules are

$$\mathbf{p}_t^i(\tilde{s}|a, s) \propto \exp[r_t(\tilde{s}, a, s)/\lambda], \quad \pi_t^i(a|s) \propto \sum_{\tilde{s} \in \mathbf{s}} \exp[r_t(\tilde{s}, a, s)/\lambda]. \quad (15)$$

The optional bound  $h$  in (13) determines the scalar parameter  $\lambda = \lambda(h) > 0$  and

$$\lim_{h \rightarrow \mathbf{H}^{\pi^{\text{MDP}}}} \lambda(h) = 0. \quad (16)$$

*Proof.* It can be directly verified that any policy, which replaces some deterministic rules of the policy  $\pi^{\text{MDP}}$  by randomised ones has a higher entropy. Thus, when maximising the expected accumulated reward (2), under the inequality constraint (13), the constraint becomes active. The maximisation, equivalent to the negative-reward minimisation, reduces to the unconstrained minimisation of the Kuhn-Tucker functional [19], given by the multiplier  $\lambda = \lambda(h) > 0$ ,

$$\begin{aligned} \pi^{\text{FPD}} &\in \text{Arg min}_{\pi \in \boldsymbol{\pi}} \sum_{b \in \mathbf{b}} c^\pi(b) [-R(b) + \lambda \ln(c^\pi(b))] \\ &= \text{Arg min}_{\pi \in \boldsymbol{\pi}} \sum_{b \in \mathbf{b}} c^\pi(b) \ln \left( \frac{c^\pi(b)}{\exp[R(b)/\lambda]} \right) = \text{Arg min}_{\pi \in \boldsymbol{\pi}} D(c^\pi \| c^i). \end{aligned}$$

The additive reward form (12), standard conditioning and marginalisation imply the forms of the ideal factors (15). The limiting property (16) corresponds with the gradual relaxation of the constraint (13).  $\square$

#### Remarks

- ✓ The role of the *ideal* decision rule (15) differs from the closely-related Boltzmann's machine, which uses the decision rules

$$\pi_t(a|s) \propto \exp \left( \sum_{\tilde{s} \in \mathbf{s}} r_t(\tilde{s}|a, s) \mathbf{p}_t(\tilde{s}|a, s) / \lambda \right), \quad \lambda > 0. \quad (17)$$

- ✓ The original, less general, relation of FPD and MDP [18] led to the ideal closed-loop model  $c^{\text{iorig}}$  that also exploited the environment model  $\mathbf{p} = \prod_{t \in \mathbf{t}} \mathbf{p}_t$  (1)

$$c^{\text{iorig}}(b) \propto \mathbf{p}(b) \exp[R(b)/\lambda] \stackrel{(14)}{=} \mathbf{p}(b) c^i(b), \quad b \in \mathbf{b}. \quad (18)$$

Recovering the explorative nature of the certainty-equivalent MDP-optimal policy is the main reason for employing the constraint (13). The following accounting of the influence of the incomplete knowledge on resulting policy brings an insight into the exploration problem. Primarily, it guides the adaptive choice of  $\lambda = \lambda(h) > 0$  parameterising the ideal closed-loop model (14).

The policy  $\pi^{\text{MDP}}(\theta)$ , which maximises the expected cumulative reward while using a given parameter  $\theta \in \boldsymbol{\theta}$ , consists of the MDP-optimal deterministic rules

$$\pi_t^{\text{MDP}}(a|s, \theta) = 1 \text{ if } a = a_t^{\text{MDP}}(s, \theta), \quad \pi_t^{\text{MDP}}(a|s, \theta) = 0 \text{ otherwise.} \quad (19)$$

There  $a_t^{\text{MDP}}(s, \theta)$  is the maximising argument in the  $t$ th step of DP (4) with the explicit conditioning on  $\theta \in \boldsymbol{\theta}$ .

The decision rules<sup>3</sup>  $\pi_t^{\text{FPD}}(a|s, V, \lambda)$  of the constructed approximation of the FPD-optimal policy should approximate the policy  $\pi^{\text{MDP}}(\theta)$  made of the rules (19) exploiting the knowledge of the parameter  $\theta$

$$\pi^{\text{MDP}}(\theta) = (\pi_t^{\text{MDP}}(a_t|s_{t-1}, \theta))_{t \in \mathbf{t}}, \quad a_t \in \mathbf{a}, \quad s_{t-1} \in \mathbf{s}.$$

The approximate policy has the posterior probability  $\mathfrak{p}(\theta|s, V)$  as the only information about  $\theta \in \boldsymbol{\theta}$ .

Works [6, 17] imply that the expected Kullback-Leibler divergence of  $\pi^{\text{MDP}}(\theta)$  from  $\pi^{\text{FPD}}$  is the adequate proximity measure to be minimised by

$$\pi^{\text{FPD}} = (\pi_t^{\text{MDP}}(a_t|s_{t-1}, V_{t-1}))_{t \in \mathbf{t}}, \quad a_t \in \mathbf{a}, \quad V_{t-1} \text{ given by (7)}$$

via the adequately chosen  $\lambda^{\text{FPD}} = \lambda^{\text{FPD}}(s_{t-1}, V_{t-1})_t$ . This dictates the selection

$$\lambda^{\text{FPD}}(s_{t-1}, V_{t-1}) \in \text{Arg min}_{\lambda > 0} \quad (20)$$

$$\int_{\boldsymbol{\theta}} \sum_{a \in \mathbf{a}} \pi_t^{\text{MDP}}(a|s_{t-1}, \theta) \ln \left( \frac{\pi_t^{\text{MDP}}(a|s_{t-1}, \theta)}{\pi_t^{\text{FPD}}(a|s_{t-1}, V_{t-1}, \lambda)} \right) \mathfrak{p}(\theta|s_{t-1}, V_{t-1}) \, d\theta.$$

The optimal actions  $a_t^{\text{MDP}}(s, \theta)$  depend on the parameter  $\theta \in \boldsymbol{\theta}$  in a quite complex way. This makes us to solve (20) for greedy (one-stage-ahead) FPD. Importantly, the resulting ideal factors (15) are used in the multi-step policy design. Thus, the dynamic nature of the policy design is not compromised unlike in the wide-spread solutions of the exploration problem [31].

For choosing  $\lambda^{\text{FPD}}(s, V)$ , at the observed  $s = s_{t-1}$  and given  $V = V_{t-1}$ , let us define, cf. (2), (6),

$$\boldsymbol{\theta}_a = \left\{ \theta \in \boldsymbol{\theta} : \sum_{\tilde{s} \in \mathbf{s}} r_t(\tilde{s}, a, s) \theta(\tilde{s}|a, s) \geq \sum_{\tilde{s} \in \mathbf{s}} r_t(\tilde{s}, \tilde{a}, s) \theta(\tilde{s}|\tilde{a}, s), \quad \forall \tilde{a} \in \mathbf{a} \right\}, \quad \forall a \in \mathbf{a}. \quad (21)$$

<sup>3</sup> The dependence on  $\lambda$  is stressed by the extended condition.

On  $\theta_a$ , the action is optimal,  $a = a^{\text{MDP}}(\theta)$ . For the FPD-optimal greedy decision rule (11) and the ideal factors (15), the optimisation (20) reads<sup>4</sup>

$$\begin{aligned} \lambda^{\text{FPD}}(s, V) &\in \text{Arg min}_{\lambda > 0} \sum_{a \in \mathbf{a}} \mathbf{p}(\theta_a | s, V) \left[ -\bar{r}(a, s, V)/\lambda - \mathbf{H}(a, s, V) \right. \\ &\quad \left. + \ln \left( \sum_{\tilde{a} \in \mathbf{a}} \exp \left( +\bar{r}_t(\tilde{a}, s, V)/\lambda + \mathbf{H}(\tilde{a}, s, V) \right) \right) \right], \quad \text{where} \\ \bar{r}_t(a, s, V) &= \sum_{\tilde{s} \in \mathbf{s}} r_t(\tilde{s}, a, s) \mathbf{p}(\tilde{s} | a, s, V), \\ \mathbf{H}(a, s, V) &= - \sum_{\tilde{s} \in \mathbf{s}} \mathbf{p}_t(\tilde{s} | a, s, V) \ln(\mathbf{p}_t(\tilde{s} | a, s, V)), \\ \mathbf{p}(\theta_a | s, V) &= \int_{\theta_a} \mathbf{p}(\theta | s, V) d\theta. \end{aligned} \quad (22)$$

Numerical solution of the scalar minimisation (22) is simple and can be done by any off-the-shelf software. The evaluation of probabilities  $\mathbf{p}(\theta_a | s, V)$  (22) of the sets  $\theta_a$ ,  $a \in \mathbf{a}$  (21) is the only more involved step. Even it can be made by a straightforward Monte Carlo integration without excessive demands on its precision.

### 3 Experiments

This part provides a representative sample of Monte Carlo studies from [14].

*The simulated environment* corresponded to MDP with  $|\mathbf{s}| = 10$  possible states,  $|\mathbf{a}| = 5$  possible actions. These options balanced the wish to deal with a non-trivial example and to perform extensive Monte Carlo experiments within a reasonable time even in the experimental Matlab implementation. Numerical values of the time-invariant simulated environment model  $\mathbf{p}$  and of the time-invariant reward  $\mathbf{r}$  are in Appendix.

The considered number of epochs was  $|\mathbf{t}| = 10 \ll |\boldsymbol{\theta}| \approx |\mathbf{s}|^2 \times |\mathbf{a}| = 500$ . As already said, the proper balancing of exploration with exploitation is vital under the conditions of this type.

*The compared policies* are summarised in Table 1, which provides their labels, under which they are referred to in the figures. The table briefly characterises them and refers to their descriptions.

Policies depending on a fixed  $\lambda$  were judged on the uniform grid

$$\lambda \in \{0.15, 0.20, 0.25, \dots, 3.60\}. \quad (23)$$

<sup>4</sup> In experiments,  $\lambda^{\text{FPD}}$  was also optimised for the original ideal closed loop model (18). Then,  $\lambda^{\text{FPD}}$  minimises an analogy of (22).



**Table 1.** Compared Policies (CE is certainty equivalent version; model means environment model).

Label	Characterisation	Reference
DPknownPar	MDP, known model	(1), (2), (4)
DP	MPD, learnt model, CE	(2), (4), (7), (8)
FPD	FPD, learnt model, CE, former ideal given $\lambda$	(11), (18), (7), (8)
FPDAdaptive	FPD, learnt model, CE, former ideal adapting $\lambda$	(11), (18), (7), (8), (22)
FPDExp	FPD, learnt model, CE, proposed ideal given $\lambda$	(11), (14), (7), (8)
FPDExpAdaptive	FPD, learnt model, CE, proposed ideal adapting $\lambda$	(11), (14), (7), (8), (22)
Boltzmann	Greedy MDP, the learnt model, CE Boltzmann’s machine, learnt model, given $\lambda$	(7), (8), (17)
eps-Greedy	Greedy MDP, learnt model, CE uniform noise	[28]
UCB1	injected with probability $\varepsilon = 0.3$ Greedy MDP, learnt model, CE, noise tuned according upper confidence bound	[2, 25]

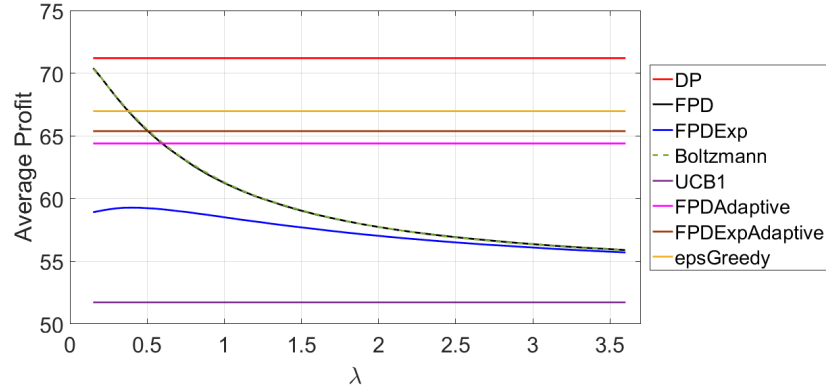
*The policy quality* was quantified by the sample mean (referred to as the average profit) of cumulative rewards  $R$  (12) evaluated for  $10^5$  Monte Carlo runs. Preliminary experiments verified that this number is more than sufficient to guarantee the representability of the results.

*Results* showing that the exploration is not necessarily helpful are in Figure 1 with abbreviations referring to labels in Table 1. They were obtained within *the first experiment* where the corresponding environment is described in Table 2. The policy DPknownPar, designed under the complete knowledge, reached the average profit of 72.11. Its variance  $\sigma = 51.72$  quantifies its volatility. Straight lines correspond to policies independent of  $\lambda$  varied on the considered grid.

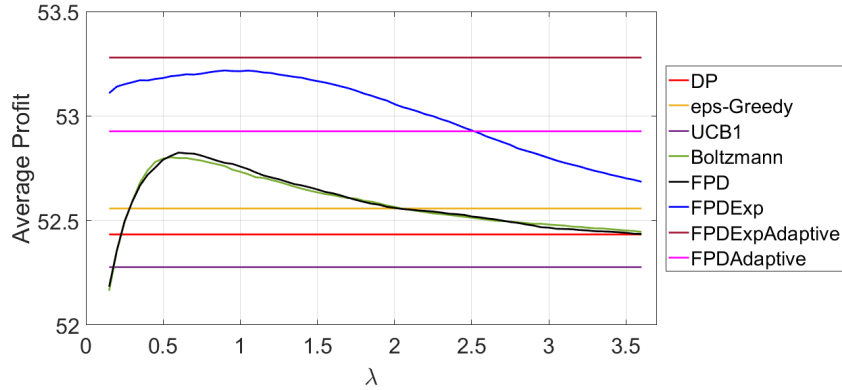
The results in which exploration was significant, were gained within *the second experiment*. They are summarised in Figure 2 with abbreviations again referring to labels in Table 1. The corresponding environment is described in Table 3. The policy DPknownPar, designed under the complete knowledge, reached the average profit of 62.84 and variance  $\sigma = 149.78$ .

*Discussion* starts with stressing that the inspected small number of epochs  $|t|$  respects that the exploration-exploitation balance is vital in this case. Otherwise, even a rare adding of random steps, whose non-optimal character with respect to the exploitation has negligible influence, guarantees convergence of learning and thus the policy optimality. This distinguishes our experiments from usual tests, e.g. [22].

The experiments dealt with structurally same static DM. The numerical choice of their parameters was based on the following, qualitatively obvious



**Fig. 1.** The results of the first experiment. The average profit is the sample mean of cumulative rewards (12) for the compared policies, Table 1, and different  $\lambda$  values on the grid (23).



**Fig. 2.** The results of the second experiment. The average profit is the sample mean of cumulative rewards (12) for compared policies, Table 1, and different  $\lambda$  values on the grid (23).

fact. The need for exploration (within the considered short-horizon scenario) depends on the mutual relation of the prior probability  $p(\theta|V_0)$  and the parameter  $\theta_{simulated}$  of the simulated environment model, see Table 4. The influence of this relation is enhanced or attenuated by the considered reward  $r$ . Note that

The first experiment, reflected in Figure 1, in which the DP policy is the best one warns that exploration need not be always helpful. Notably, FPD and Boltzmann's machine with sufficiently small  $\lambda$  can be arbitrarily close to this best behaviour. Due to the lack of exploration significance no other conclusions concerning the quality of the tested policies can be made. But it calls for an

improvement of  $\lambda$  tuning, which should converge to zero if the exploration is superfluous.

The second experiment, reflected in Figure 2, is more informative. The policy based on a newly proposed relation of FPD with MDP and an adaptive choice of  $\lambda$  (FPDExpAdaptive) brings the highest improvement (about 2%). A similar performance can be reached for a fixed but properly chosen  $\lambda$  (FPDExp). The adaptive FPD is worse (FPDAdaptive) but still outperforms the remaining competitors. The similarity of the results for the  $\lambda$ -dependent FPD and Boltzmann's machine supports the conjecture that the performance of Boltzmann's machine can be improved by adapting  $\lambda$ . This may be important in its other applications.

## 4 Concluding Remarks

The paper has arisen from inspecting the conjecture that the certainty-equivalent version of non-traditional fully probabilistic design (FPD) of decision policies properly balances exploitation with exploration. The achieved results supported it. Moreover the paper: (a) established a better relation of FPD to the widespread Markov decision processes; (b) proposed an adaptive tuning of the involved parameter, which can be used in the closely-related simulated annealing and Boltzmann's machine; (c) provided a sample of extensive experiments, which confirmed that standard exploration techniques are outperformed by the FPD-based policies.

The future work will concern: (i) an algorithmic recognition of cases in which exploration is unnecessary; (ii) inspection of a tuning mechanism based on extremum-seeking control; (iii) an efficient implementation of  $\lambda$ -tuning; (iv) application of the proposed ideas to continuous-valued MDP; (v) real-life problems, especially those in which a short, but non-unit, decision horizon is vital as in environmental decision making [24].

## References

1. Åström, K.: Introduction to Stochastic Control. Academic Press, NY (1970)
2. Auer, P., Cesa-Bianchi, N., Fischer, P.: Finite-time analysis of the multiarmed bandit problem. *Machine Learning* **47**(2-3), 235–256 (2002)
3. Barndorff-Nielsen, O.: Information and Exponential Families in Statistical Theory. Wiley, NY (1978)
4. Bellman, R.: Adaptive Control Processes. Princeton U. Press, NJ (1961)
5. Berger, J.: Statistical Decision Theory and Bayesian Analysis. Springer, NY (1985)
6. Bernardo, J.: Expected information as expected utility. *The An. of Stat.* **7**(3), 686–690 (1979)
7. Bertsekas, D.: Dynamic Programming and Optimal Control. Athena Scientific, US (2001)
8. Cover, T., Thomas, J.: Elements of Information Theory. Wiley (1991), 2nd edition
9. Črepinšek, M., Liu, S., Mernik, M.: Exploration and exploitation in evolutionary algorithms: A survey. *ACM Computing Survey* **45**(3), 37–44 (2013)

10. Duff, M.O.: Optimal Learning; Computational Procedures for Bayes-Adaptive Markov Decision Processes. Ph.D. thesis, University of Massachusetts Amherst (2002)
11. Feldbaum, A.: Theory of dual control. *Autom. Remote Control* **21,22**(9,2) (1960,61)
12. Gómez, A.G.V., Kappen, H.: Dynamic policy programming. *The J. of Machine Learning Research* **30**, 3207–3245 (2012)
13. Guan, P., Raginsky, M., Willett, R.: Online Markov decision processes with Kullback-Leibler control cost. In: *American Control Conference*. pp. 1388–1393. IEEE (2012)
14. Hůla, F.: Explorative Fully Probabilistic Design of Adaptive Decision Strategies. Master's thesis, FJFI, Czech Technical University, Prague (2018)
15. Kappen, H.: Linear theory for control of nonlinear stochastic systems. *Physical review letters* **95**(20), 200201 (2005)
16. Kárný, M., Guy, T.V.: Fully probabilistic control design. *Systems & Control Letters* **55**(4), 259–265 (2006)
17. Kárný, M., Guy, T.: On support of imperfect Bayesian participants. In: Guy, T., et al (eds.) *Decision Making with Imperfect Decision Makers*, vol. 28. Springer, Berlin (2012), *Intelligent Systems Reference Library*
18. Kárný, M., Kroupa, T.: Axiomatisation of fully probabilistic design. *Infor. Sciences* **186**(1), 105–113 (2012)
19. Kuhn, H., Tucker, A.: Nonlinear programming. In: *Proc. of 2nd Berkeley Symposium*, pp. 481–492. Univ. of California Press (1951)
20. Kullback, S., Leibler, R.: On information and sufficiency. *Annals of Mathematical Statistics* **22**, 79–87 (1951)
21. Kumar, P.: A survey on some results in stochastic adaptive control. *SIAM J. Control and Applications* **23**, 399–409 (1985)
22. Ouyang, Y., Gagrani, M., Nayyar, A., Jain, R.: Learning unknown Markov decision processes: A Thompson sampling approach. In: et al, I.G. (ed.) *Advances in Neural Information Processing Systems* 30, pp. 1333–1342. Curran Associates, Inc. (2017), <http://papers.nips.cc/paper/6732-learning-unknown-markov-decision-processes-a-thompson-sampling-approach.pdf>
23. Puterman, M.: *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley (2005)
24. Springborn, M.: Risk aversion and adaptive management: Insights from a multi-armed bandit model of invasive species risk. *Journal of Environmental Economics and Management* **68**, 226–242 (2014)
25. Tang, H., et al: #Exploration: A study of count-based exploration for deep reinforcement learning. In: et al, I.G. (ed.) *Advances in Neural Information Processing Systems* 30, pp. 2753–2762. Curran Associates, Inc. (2017), <http://papers.nips.cc/paper/6868-exploration-a-study-of-count-based-exploration-for-deep-reinforcement-learning.pdf>
26. Tanner, M.: *Tools for statistical inference*. Springer Verlag, NY (1993)
27. Ullrich, M.: Optimum control of some stochastic systems. In: *Preprints of the VIII-th conference ETAN*. Beograd (1964)
28. Vermorel, J., Mohri, M.: Multi-armed bandit algorithms and empirical evaluation. In: *European conference on machine learning*. pp. 437–448. Springer (2005)
29. Šindelář, J., Vajda, I., Kárný, M.: Stochastic control optimal in the Kullback sense. *Kybernetika* **44**(1), 53–60 (2008)

30. Witten, I., Frank, E., Hall, M., Pal, C.: Data Mining: Practical machine learning tools and techniques. 4th edition, Elsevier (2017)
31. Wu, H., Guo, X., Liu, X.: Adaptive exploration-exploitation trade off for opportunistic bandits (2017), preprint arXiv:1709.04004

## Appendix

This section provides considered rewards and transition probabilities used in experiments, Section 3. Static, time-invariant cases are considered. Their transition probabilities  $\mathbf{p}(\tilde{s}|a, s) = \mathbf{p}(\tilde{s}|a)$  are the same  $\forall s \in \mathbf{s}$ .

**Table 2.** The data used in the first experiment. Explicit values of the reward  $r_t(\tilde{s}, a, s) = r(\tilde{s}, a)$ , on the left side and of the transition probabilities  $\mathbf{p}_t(\tilde{s}|a, s) = r(\tilde{s}, a)$  on the right side. They are constant  $\forall s \in \mathbf{s}$ ,  $t \in \mathbf{t}$  and  $|\mathbf{s}| = 10$ ,  $|\mathbf{a}| = 5$ . Rows and columns correspond to states  $\tilde{s} \in \mathbf{s}$  and actions  $a \in \mathbf{a}$ , respectively.

	The reward $r_t$ actions $a \in \mathbf{a}$					The transition probability $\mathbf{p}$ actions $a \in \mathbf{a}$				
states $\tilde{s} \in \mathbf{s}$	5	7	6	5	10	0.12	0.16	0.12	0.12	0.08
	1	6	1	3	6	0.02	0.13	0.08	0.02	0.02
	6	2	5	7	9	0.08	0.16	0.06	0.14	0.15
	5	6	1	5	4	0.18	0.04	0.08	0.08	0.13
	5	2	2	6	6	0.10	0.06	0.18	0.10	0.06
	4	8	6	4	5	0.02	0.10	0.16	0.10	0.09
	3	9	3	8	5	0.06	0.07	0.08	0.08	0.13
	7	5	2	6	8	0.02	0.02	0.02	0.12	0.13
	3	9	3	2	6	0.20	0.18	0.16	0.20	0.04
	3	1	4	8	10	0.20	0.09	0.06	0.04	0.17

**Table 3.** The data used in the second experiment. Explicit values of the reward  $r_t(\tilde{s}, a, s) = r(\tilde{s}, a)$ , on the left side and of the transition probabilities  $\mathbf{p}_t(\tilde{s}|a, s) = r(\tilde{s}, a)$  on the right side. They are constant  $\forall s \in \mathbf{s}$ ,  $t \in \mathbf{t}$  and  $|\mathbf{s}| = 10$ ,  $|\mathbf{a}| = 5$ . Rows and columns correspond to states  $\tilde{s} \in \mathbf{s}$  and actions  $a \in \mathbf{a}$ , respectively.

	The reward $r_t$ actions $a \in \mathbf{a}$					The transition probability $\mathbf{p}$ actions $a \in \mathbf{a}$				
states $\tilde{s} \in \mathbf{s}$	1	1	1	1	1	0.03	0.05	0.03	0.02	0.08
	2	2	2	2	2	0.05	0.07	0.09	0.05	0.05
	3	3	3	3	3	0.08	0.12	0.14	0.07	0.08
	3	3	3	3	3	0.08	0.07	0.09	0.07	0.05
	5	5	5	5	5	0.11	0.17	0.11	0.12	0.11
	6	6	6	6	6	0.29	0.31	0.20	0.15	0.30
	12	12	12	12	12	0.13	0.07	0.09	0.29	0.16
	4	4	4	4	4	0.11	0.05	0.11	0.10	0.08
	3	3	3	3	3	0.08	0.07	0.09	0.07	0.05
	2	2	2	2	2	0.05	0.02	0.06	0.05	0.03

**Table 4.** Prior probability  $\mathbf{p}(\theta|V_0)$  for both the first experiment on the left side and the second experiment on the right side. They are constant  $\forall s \in \mathbf{s}$ ,  $t \in \mathbf{t}$  and  $|\mathbf{s}| = 10$ ,  $|\mathbf{a}| = 5$ . Rows and columns correspond to states  $\tilde{s} \in \mathbf{s}$  and actions  $a \in \mathbf{a}$ , respectively.

	The first experiment actions $a \in \mathbf{a}$					The second experiment actions $a \in \mathbf{a}$				
states $\tilde{s} \in \mathbf{s}$	0.1	0.1	0.1	0.1	0.1	0.03	0.04	0.02	0.06	0.08
	0.1	0.1	0.1	0.1	0.1	0.05	0.06	0.05	0.09	0.05
	0.1	0.1	0.1	0.1	0.1	0.08	0.11	0.07	0.09	0.08
	0.1	0.1	0.1	0.1	0.1	0.08	0.06	0.07	0.09	0.05
	0.1	0.1	0.1	0.1	0.1	0.11	0.15	0.10	0.11	0.11
	0.1	0.1	0.1	0.1	0.1	0.29	0.19	0.36	0.26	0.30
	0.1	0.1	0.1	0.1	0.1	0.13	0.06	0.12	0.14	0.16
	0.1	0.1	0.1	0.1	0.1	0.11	0.11	0.10	0.06	0.08
	0.1	0.1	0.1	0.1	0.1	0.08	0.09	0.07	0.06	0.05
	0.1	0.1	0.1	0.1	0.1	0.05	0.13	0.05	0.06	0.03