

Chapter 3

Automated Preference Elicitation for Decision Making

Miroslav Kárný

Abstract. In the contemporary complex world decisions are made by an imperfect participant devoting limited deliberation resources to any decision-making task. A normative decision-making (DM) theory should provide support systems allowing such a participant to make rational decisions in spite of the limited resources. Efficiency of the support systems depends on the interfaces enabling a participant to benefit from the support while exploiting the gradually accumulating knowledge about DM environment and respecting incomplete, possibly changing, participant's DM preferences. The insufficiently elaborated preference elicitation makes even the best DM supports of a limited use. This chapter proposes a methodology of automatic eliciting of a quantitative DM preference description, discusses the options made and sketches open research problems. The proposed elicitation serves to fully probabilistic design, which includes a standard Bayesian decision making.

Keywords: Bayesian decision making, fully probabilistic design, DM preference elicitation, support of imperfect participants.

3.1 Introduction

This chapter concerns of an imperfect participant¹, which solves a real-life decision-making problem under uncertainty, which is worth of its optimising effort. The topic has arisen from the recognition that a real participant often cannot benefit from sophisticated normative DM theories due to an excessive deliberation effort needed

Miroslav Kárný
Institute of Information Theory and Automation,
Academy of Sciences of the Czech Republic, Pod vodárenskou věží 4, 182 08 Prague 8,
Czech Republic
e-mail: school@utia.cas.cz

¹ A participant is also known as user, decision maker, agent. A participant can be human, an artificial object or a group of both. We refer the participant as "it".

for their mastering and for feeding them by DM elements² they need. This observation has stimulated a long-term research, which aims to equip a participant with automatic tools (intelligent interfaces) mapping its knowledge, DM preferences and constraints on DM elements while respecting its *imperfection*, i.e. ability to devote only a limited deliberation effort to a particular DM task. The research as well as this chapter concentrates on the Bayesian DM theory because of its exceptional, axiomatically justified, role in DM under uncertainty, e.g. [49].

The adopted *concept of the ultimate solution* considers creating an automated supporting system, which covers the complete design and use of a decision-generating DM strategy. It has to preserve the theoretically reachable DM quality and free the participant's cognitive resources to tasks specific to its application domain. This concept induces:

Requirement 1: The supporting system uses a consistent and complete DM theory.

Requirement 2: To model the environment³ the supporting system fully exploits both participant's knowledge and information brought by the data observed during the use of the DM strategy.

Requirement 3: The supporting system respects participant's DM preferences and refines their description by the information gained from the observed data.

This chapter represents a further step to the ultimate solution. It complements the results of the chapter [28] devoted to DM of imperfect participants. The tools needed for a conceptual solution are based on a generalised Bayesian DM, called *fully probabilistic design* (FPD) [23, 29], see also Section 3.2. The FPD minimises Kullback-Leibler divergence [38] of the optimised, strategy-dependent probabilistic model of the closed DM loop on its ideal counterpart, which describes the desired behaviour of the closed decision loop. The design replaces the maximisation of an expected utility over a set of admissible decision strategies [4] for the FPD densely extending all standard Bayesian DM formulations [31]. The richness, intuitive plausibility, practical advantages and axiomatic basis of the FPD motivate its acceptance as a unified theoretical DM basis, which meets Requirement 1.

Requirement 2 concerns the description of the environment with which the participant interacts during DM. Traditionally, its construction splits into a structural and semi-quantitative modelling of the environment and *knowledge elicitation* understood as a quantitative description of unknown variables entering the environment model. Both activities transform domain-specific knowledge into the model-related part of DM elements.

The environment modelling is an extremely wide field that exploits first principles (white and grey box models, e.g. [6, 22]), application-field traditions, e.g. [9],

² This is a common label for all formal, qualitatively and quantitatively specified, objects needed for an exploitation of the selected normative DM theory.

³ The environment, also called system, is an open part of the World considered by the participant and with which it interacts within the solved DM task.

universal approximation (black box models, e.g. [17, 50]) and their combinations. An automated mapping of these models on probabilistic DM elements of the FPD is the expected service of the supporting DM system. The tools summarised in [28] are conjectured to be sufficient to this purpose. Knowledge elicitation in the mentioned narrow sense is well surveyed in [14, 45] and automated versions related to this chapter are in [24, 25]. The ordinary Bayesian framework [4, 47] adds the required ability to learn from the observed data.

Requirement 3 reflects the fact that a feasible and effective solution of *preference elicitation* problem decides on the efficiency of any intelligent system supporting DM. This extracting the information about the participant's DM preferences has been recognised as a vital problem and repeatedly addressed within artificial intelligence, game theory, operation research. Many sophisticated approaches have been proposed [10, 11, 13, 16], often in connection with applied sciences like economy, social science, clinical decision making, transportation, see, for instance, [21, 41].

Various assumptions on the structure of DM preferences have been adopted in order to ensure feasibility and practical applicability of the resulting decision support. Complexity of the elicitation problem has yet prevented to get a satisfactory widely-applicable solution. For instance, the conversion of DM preferences on individual observable decision-quality-reflecting *attributes* into the overall DM preferences is often done by assuming their additive independence [33]. The DM preferences on attributes are dependent in majority of applications and the enforced independence assumption significantly worsens the elicitation results⁴. This example indicates a deeper drawback of the standard Bayesian DM, namely, the lack of unambiguous rules how to combine low-dimensional description of DM preferences into a global one.

The inability of the participant to completely specify its DM preferences is another problem faced. In this common case, the DM preferences should be learned from either domain-specific information (technological requirements and knowledge, physical laws, etc.) or the observed data.

Eliciting the needed information itself is an inherently difficult task, which success depends on experience and skills of an elicitation expert. The process of eliciting of the domain-specific information is difficult, time-consuming and error-prone activity⁵. Domain experts provide subjective opinions, typically expressed in different forms. Their processing requires a significant cognitive and computational effort of the elicitation expert. Even if the cost of this effort⁶ is negligible, the elicitation

⁴ The assumption can be weakened by a introducing a conditional preferential independence, [8].

⁵ It should be mentioned that practical solutions mostly use a laborious and unreliable process of manual tuning a number of parameters of the pre-selected utility function. The high number of parameters makes this solution unfeasible and enforces attempts to decrease the number to recover feasibility.

⁶ This effort is usually very high and many sophisticated approaches aim at optimising a trade-off between elicitation cost and value of information it provides (often, a decision quality is considered), see for instance [7].

result is always limited by the expert's imperfection, i.e. his/her inability to devote an unlimited deliberation effort to eliciting. Unlike the imperfection of experts providing the domain-specific information, the imperfection of elicitation experts can be eliminated by preparing a feasible *automated* support of the preference elicitation, that does not rely on any elicitation expert.

The dynamic decision making strengthens the dependence of the DM quality on the preference elicitation. Typically, the participant acting within a dynamically changing environment with evolving parameters gradually changes its DM preferences. The change may depend on the expected future behaviour or other circumstances. The overall task is getting even harder when the participant dynamically interacts with other imperfect participants within a common environment. When DM preferences evolve, their observed-data-based learning becomes vital.

The formal disparity of modelling language (probabilities) and the DM preference description (utilities) makes Bayesian learning of DM preferences difficult. It needs a non-trivial "measurement" of participant's satisfaction of the decision results, which often puts an extreme deliberation load on the participant. Moreover, the degree of satisfaction must be related to conditions under which it has been reached. This requires a non-standard and non-trivial modelling. Even, if these learning obstacles are overcome, the space of possible behaviour is mostly larger than that the observed data cover. Then, the initial DM preferences for the remaining part of the behaviour should be properly assigned and exploration has to care about making the DM preference description more precise. Altogether, a weak support of the preference elicitation (neglecting of Requirement 3) is a significant gap to be filled. Within the adopted FPD, an ideal probability density (pd^7) is to be elicited⁸. The ideal pd describes the closed-loop behaviour, when the participant's DM strategy is an optimal one and the FPD searches for the optimal randomised strategy minimising the divergence from the current closed-loop description to the ideal one.

Strengthening the support with respect to Requirement 3 forms the core of this chapter. The focus on the preference elicitation for the FPD brings immediate methodological advantages. For instance, the common probabilistic language for knowledge and DM preference descriptions simplifies an automated elicitation as the ideal pd provides a standardised form of quantitatively expressed DM preferences. Moreover, the raw elicitation results reflect inevitably incomplete, competitive or complementing opinions with respect to the same collection of DM preference-expressing multivariate attributes. Due to their automated mapping on probabilities, their logically consistent merging is possible with the tools described in [28]. Besides, domain experts having domain-specific information are often

⁷ Radon-Nikodým derivative [48] of the strategy-dependent measure describing closed DM loop with respect to a dominating, strategy-independent measure. The use of this notion helps us to keep a unified notation that covers cases with mixed – discrete and continuous valued – variables.

⁸ Let us stress that no standard Bayesian DM is omitted due to the discussed fact that the FPD densely covers all standard Bayesian DM tasks.

unable to provide their opinion on a part of behaviour due to either limited knowledge of the phenomena behind or the indifference towards the possible instances of behaviour. Then, the DM preference description has to be extended to the part of behaviour not “covered” by the domain-specific information. This extension is necessary as the search for the optimal strategy heavily depends on the full DM preference description. It is again conceptually enabled by the tools from [28]. The usual Bayesian learning is applicable whenever the DM preferences are related to the observed data [27].

In summary, the chapter concerns a construction of a probabilistic description of the participant’s DM preferences based on the available information. Decision making under uncertainty is considered from the perspective of an imperfect participant. It solves a DM task with respect to its environment and indirectly provides a finite description of the DM preferences in a non-unique way⁹ and leaves uncertainty about the DM preferences on a part of closed-loop behaviour. To design an optimal strategy, the participant employs the fully probabilistic design of DM strategies [23, 29] whose DM elements are probability densities used for the environment modelling, DM preference description and description of the observed data.

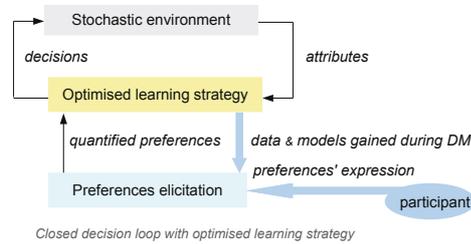
The explanations prefer discussion of the solution aspects over seemingly definite results. After a brief summary of common tools Section 3.2, they start with a problem formalisation that includes the basic adopted assumptions, Section 3.3. The conceptual solution summarised in Section 3.4 serves as a guide in the subsequent extensive discussion of its steps in Section 3.5. Section 3.6 provides illustrative simulations and Section 3.7 contains concluding remarks.

Concept of the proposed preference elicitation is reflected in Figure¹⁰ 3.1. The usual decision loop formed by a stochastic environment and a decision strategy complemented by a preference elicitation block is expanded to the proposed solution. The considered strategy consists of the standard Bayesian learning of the environment model and of a standard fully probabilistic design (FPD). Its explicit structuring reveals the need of the ideal closed-loop model of the desired closed-loop behaviour. The designed strategy makes the closed decision loop closest to this ideal model, which is generated by the elicitation block as follows. The observed data is censored¹¹ to data, which contains an information about the optimal strategy and serves for its Bayesian learning. The already learnt environment model is combined with the gained model of the optimal strategy into the model of the DM loop closed by it. Within the set of closed-loop models, which comply with the participant’s DM preferences and are believed to be reachable, the ideal closed-loop model is selected as the nearest one to the learnt optimal closed-loop model.

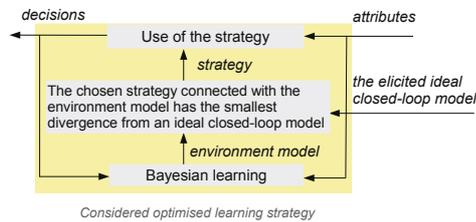
⁹ Even, when we identify instances of behaviour that cannot be preferentially distinguished.

¹⁰ A block performing the inevitable knowledge elicitation is suppressed to let the reader focus on the proposed preference elicitation.

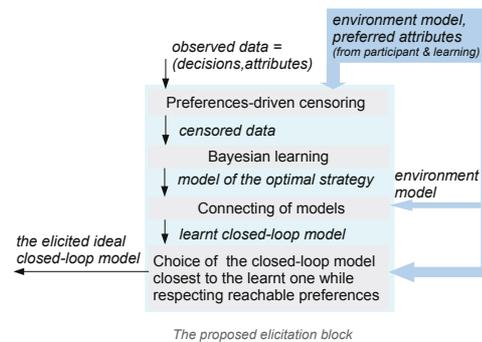
¹¹ Such data processing is also called filtering. This term is avoided as it has also another meaning.



(a)



(b)



(c)

Fig. 3.1 The figure 3.1a displays a closed decision loop with an optimised learning strategy. The figure 3.1b expands the considered optimised learning strategy that uses the FPD and Bayesian learning. The figure 3.1c shows the *proposed elicitation block*. The observed data is censored to reveal information about an unknown optimal strategy. The Bayesian learning on the censored data provides the model of the optimal strategy, which together with the learnt environment model describes the optimally closed loop. The elicitation block selects the ideal closed-loop model as the model, which: (i) complies with participant's DM preferences; (ii) is reachable by an available strategy; (iii) is the nearest one to the model of the optimal closed loop.

Notation*General conventions*

\mathbf{x}	is a set of x -values having cardinality $ \mathbf{x} $
$d \in \mathbf{d}, \mathbf{d} \neq \emptyset$	are decisions taken from a finite-dimensional set \mathbf{d}
$a_i \in \mathbf{a}_i, i \in \mathbf{i} = \{1, \dots, \mathbf{i} \}$	are attribute entries in finite-dimensional sets \mathbf{a}_i
$\mathbf{a} \in \mathbf{A}$	is a collection of all attributes in the set $\mathbf{a} = \times_{i \in \mathbf{i}} \mathbf{a}_i$,
\mathbf{X}	denotes the Cartesian product
$\boldsymbol{\alpha} \subseteq \mathbf{a}, \boldsymbol{\alpha} \neq \emptyset$	is the set of the most desirable attribute values specified entry-wise $\boldsymbol{\alpha} = \times_{i \in \mathbf{i}} \alpha_i$
$t \in \mathbf{t} = \{1, \dots, \mathbf{t} \}$	is discrete time
$(x_t)_{t \in \mathbf{t}}$	is a sequence of x_t indexed by discrete time $t \in \mathbf{t}$.

Probability densities

$g(\cdot), h(\cdot)$	are probability densities (pds): Radon-Nikodým derivatives with respect to a dominating measure denoted d
$M_t(a d), M(a d, \Theta)$	are the environment model and its parametric version with an unknown parameter $\Theta \in \boldsymbol{\Theta}$
$F_t(\Theta), t \in \mathbf{t} \cup \{0\}$	is the pd quantifying knowledge available at time t about the unknown parameter Θ of the environment model
$S_t(d)$	describes the randomised decision strategy to be selected
$s_t(d), s(d \theta)$	are the model of the optimal strategy and its parametric version with an unknown parameter $\theta \in \boldsymbol{\theta}$
$f_t(\theta), t \in \mathbf{t} \cup \{0\}$	is the pd quantifying knowledge available at time t about the unknown parameter $\theta \in \boldsymbol{\theta}$ of the optimal strategy
$l_t(a, d)$	is the ideal pd quantifying the elicited participant's DM preferences
$P_t(a, d)$	is the pd modelling the decision loop with the optimal strategy
$M_t(a_i d), l_t(a_i), i \in \mathbf{i}$	are marginal pds of a_i derived from pds $M_t(a d)$ and $l_t(a, d)$

Convention on time indexing

$F_{t-1}(\Theta), f_{t-1}(\theta)$	quantify knowledge accumulated before time t
$M_t(a d), s_t(d), S_t(d)$	serve to the t th DM task and exploit
$l_t(a, d), P_t(a, d)$	the knowledge accumulated before time t .

Frequent symbols

$\bar{d} \in \mathbf{d}$	is a decision leading to a high probability of the set $\boldsymbol{\alpha}$
$\mathbb{D}(h g)$	is the Kullback-Leibler divergence (KLD, [38]) of a pd h from a pd g
$\mathbb{E}[\cdot]$	denotes expectation
V	is a sufficient statistic of the exponential family (EF), which becomes the occurrence table in Markov-chain case
$\phi \in [0, 1]$	is a forgetting factor
\propto	denotes an equality up to normalisation.

3.2 Preliminaries

Introduction repeatedly refers to the tools summarised in [28]. Here, we briefly recall its sub-selection used within this chapter.

1. The *Kullback-Leibler divergence* (KLD, [38]) $\mathbb{D}(\mathbf{g}||\mathbf{h})$ of a pd \mathbf{g} from a pd \mathbf{h} , both defined on a set \mathbf{x} and determined by a dominating strategy-independent measure dx , is defined by the formula

$$\mathbb{D}(\mathbf{g}||\mathbf{h}) = \int_{\mathbf{x}} \mathbf{g}(x) \ln \left(\frac{\mathbf{g}(x)}{\mathbf{h}(x)} \right) dx. \quad (3.1)$$

The KLD is a convex functional in the pd \mathbf{g} , which reaches its smallest zero value iff $\mathbf{g} = \mathbf{h}$ dx -almost everywhere.

$\mathbb{D}(\mathbf{g}||\mathbf{h}) \neq \mathbb{D}(\mathbf{h}||\mathbf{g})$ and a correct pd should be used as its first argument when measuring (di)similarity of pds by the KLD. A pd is called *correct*¹² if it fully exploits the knowledge about the random variable it models. Its existence is assumed.

2. Under the commonly met conditions [5], the *optimal Bayesian approximation* $\mathbf{h}^o \in \mathbf{h}$ of a correct pd \mathbf{g} by a pd $\mathbf{h} \in \mathbf{h}$ should be defined

$$\mathbf{h}^o \in \underset{\mathbf{h} \in \mathbf{h}}{\text{Arg min}} \mathbb{D}(\mathbf{g}||\mathbf{h}). \quad (3.2)$$

3. The *minimum KLD principle* [28, 51] recommends to select a pd $\mathbf{h}^e \in \mathbf{h}$

$$\mathbf{h}^e \in \underset{\mathbf{h} \in \mathbf{h}}{\text{Arg min}} \mathbb{D}(\mathbf{h}||\mathbf{g}) \quad (3.3)$$

as an extension of the available information about the correct pd \mathbf{h} . The assumed available information consists of a given set \mathbf{h} and of a rough (typically flat) estimate \mathbf{g} of the pd \mathbf{h} .

The minimum KLD principle provides such an extension of the available information that the pd \mathbf{h}^e deviates from its estimate \mathbf{g} only to the degree enforced by the constraint $\mathbf{h} \in \mathbf{h}$. It reduces to the celebrated maximum entropy principle [20] for the uniform pd \mathbf{g} .

The paper [51] axiomatically justifies the minimum KLD principle for sets \mathbf{h} delimited by values of \mathbf{h} moments. The generalisation in [28] admits a richer collection of the sets \mathbf{h} . For instance, the set \mathbf{h} can be of the form

$$\mathbf{h} = \{ \mathbf{h} : \mathbb{D}(\mathbf{h}||\hat{\mathbf{h}}) \leq k < \infty \}, \quad (3.4)$$

determined by a given pd $\hat{\mathbf{h}}$ and by a positive constant k .

For the set (3.4), the pd \mathbf{h}^e (3.3) can be found by using the Kuhn-Tucker optimality conditions [35]. The solution reads

$$\mathbf{h}^e \propto \hat{\mathbf{h}}^\phi \mathbf{g}^{1-\phi}, \quad \phi \in [0, 1], \quad (3.5)$$

¹² This is an operational notion unlike often used adjectives “true” or “underlying”.

where \propto denotes an equality up to a normalisation factor and ϕ is to be chosen to respect the constraint (3.4). The solution formally coincides with the so-called stabilised forgetting [37] and ϕ is referred as *forgetting factor*.

3.3 Problem Formalisation

The considered participant repeatedly solves a sequence of static¹³ DM tasks indexed by (*discrete*) time $t \in \mathbf{t} = \{1, 2, \dots, |\mathbf{t}|\}$. DM concerns a stochastic, incompletely-known, time-invariant static environment. The decision d influencing the environment is selected from a finite-dimensional set \mathbf{d} . The participant judges DM quality according to a multivariate attribute $a \in \mathbf{a}$, which is a participant-specified image of the observed environment response to the applied decision. The attribute has $|\mathbf{a}| < \infty$, possibly vectorial, entries a_i . Thus, $a = (a_i)_{i=1}^{|\mathbf{a}|}$, $a_i \in \mathbf{a}_i$, $i \in \mathbf{i} = \{1, \dots, |\mathbf{a}|\}$, and \mathbf{a} is the Cartesian product $\mathbf{a} = \times_{i \in \mathbf{i}} \mathbf{a}_i$.

The solution of a sequence of static DM tasks consists of the choice and use of an *admissible randomised strategy*, which is formed by a sequence $(\mathbf{S}_t)_{t \in \mathbf{t}}$ of the randomised causal mappings

$$\mathbf{S}_t \in \mathbf{S}_t \subset \{\text{knowledge available at time } (t-1) \rightarrow d_t \in \mathbf{d}\}, t \in \mathbf{t}. \quad (3.6)$$

We accept the following basic non-restrictive assumptions.

Agreement 1 (Knowledge Available). *The knowledge available at time $(t-1)$ (3.6), $t \in \mathbf{t}$, includes*

- *the data observed up to time $(t-1)$ inclusive, i.e. decisions made (d_1, \dots, d_{t-1}) and the corresponding realisations of attributes (a_1, \dots, a_{t-1}) ;*
- *a time-invariant parametric environment model $M(a|d, \Theta) > 0$, which is a conditional pd known up to a finite-dimensional parameter $\Theta \in \Theta$;*
- *a prior pd $F_0(\Theta) > 0$ on the unknown parameter $\Theta \in \Theta$.*

The standard Bayesian learning and prediction [47] require availability of the knowledge described in Agreement 1 in order to provide the predictive pds $(M_t(a|d))_{t \in \mathbf{t}}$. They model the environment in the way needed for the design of the admissible strategy $(\mathbf{S}_t)_{t \in \mathbf{t}}$.

Agreement 2 (Optimality in the FPD Sense). *The following optimal strategy $(\mathbf{S}_t^o)_{t \in \mathbf{t}}$ in the FPD sense [31] is chosen*

$$(\mathbf{S}_t^o)_{t \in \mathbf{t}} \in \text{Arg} \min_{(\mathbf{S}_t \in \mathbf{S}_t)_{t \in \mathbf{t}}} \frac{1}{|\mathbf{t}|} \sum_{t \in \mathbf{t}} \mathbb{E}[\mathbb{D}(M_t \mathbf{S}_t || I_t)], \quad (3.7)$$

¹³ The restriction to the static case allows us to avoid technical details making understanding of the conceptual solution difficult. All results are extendable to the dynamic DM with a mixed (discrete and continuous) observed data and considered but unobserved internal variables.

where the participant's DM preferences in t th DM task are quantified by an ideal pd $l_t(a, d)$ assigning a high probability to the desirable pairs $(a, d) \in (\mathbf{a}, \mathbf{d})$ and a low probability to undesirable ones. The expectation $\mathbb{E}[\bullet]$ is taken over conditions of the individual summands in (3.7)¹⁴.

The strategy $(\mathbf{S}_t^o)_{t \in \mathbf{i}}$ minimises an average Kullback-Leibler divergence $\mathbb{D}(\mathbf{M}_t \mathbf{S}_t \| l_t)$ of the strategy-dependent closed-loop model $\mathbf{M}_t(a|d) \mathbf{S}_t(d)$ from the participant's DM preferences-expressing ideal pd $l_t(a, d)$.

Assumption 1 (Preference Specification). *The participant provides the time-invariant sets α_i , $i \in \mathbf{i}$, of the most desirable values of individual attribute entries a_i*

$$\alpha_i \subset a_i, \alpha_i \neq \emptyset, i \in \mathbf{i}. \quad (3.8)$$

These sets define the set of the most desirable attributes' values α

$$\alpha = \bigcup_{i \in \mathbf{i}} \alpha_i \subseteq \mathbf{a}, \alpha \neq \emptyset. \quad (3.9)$$

The participant can also assign *importance weights* $w \in \mathbf{w} = \{w = (w_1, \dots, w_{|\mathbf{i}|}), w_i \geq 0, \sum_{i \in \mathbf{i}} w_i = 1\}$ ¹⁵ to particular attribute entries but the availability of w is rarely realistic.

Generally, the participant may specify a number of not necessarily embedded sets α^μ , $\mu \in \boldsymbol{\mu} = \{1, \dots, |\boldsymbol{\mu}|\}$, $|\boldsymbol{\mu}| > 1$, of the desirable attribute values with the desirability decreasing with μ . The participant may also specify similar information about the possible decisions. The chosen version of the partially specified DM preferences suffices for presenting an essence of the proposed approach.

Preferences are elicited under the following non-standard assumption.

Assumption 2 (Modelling of the Unknown Optimal Strategy). *A parametric model $s(d|\theta)$ of an unknown optimal randomised strategy and a prior pd $f_0(\theta)$ of an unknown finite-dimensional $\theta \in \boldsymbol{\theta}$ parameterising this model are available.*

The feasibility of Assumption 2 follows from the time-invariance of the parametric model of the environment and from the assumed invariance of the (partially specified) participant's DM preferences¹⁶. Neither the environment model nor the complete DM preferences are known and the parameter. The only source of knowledge is observed closed-loop data. Therefore the model of the optimal strategy can be learnt from it during application of non-optimal strategy. Having this non-standard learning problem solved, the standard Bayesian prediction [47] provides the model of the optimal strategy as the predictive pd $s_t(d)$. The chain rule [47] for pds and the already learnt environment model $\mathbf{M}_t(a|d)$ imply the availability of the closed-loop model with the estimated optimal strategy $\mathbf{S}_t(d)$

¹⁴ The considered KLD measures divergence between the conditional pds. The environment model $\mathbf{M}_t(a|d)$, the optimised mapping $\mathbf{S}_t(a)$ as well as the ideal pd $l_t(a, d)$ depend on the random knowledge available at time $(t-1)$, see Agreement 1.

¹⁵ The set is referred as probabilistic simplex.

¹⁶ The proposed preference elicitation with time-invariant sets α_i can be extended to time-varying cases.

$$P_t(a, d) = M_t(a|d)s_t(d), \quad t \in \mathbf{t}. \quad (3.10)$$

Problem Formulation. Under Assumptions 1, 2 describing the available information about the environment and partially specified DM preferences¹⁷, design a well-justified automated construction of the ideal pds $(l_t)_{t \in \mathbf{t}}$ quantifying the given participant's DM preferences.

The ideal-pds construction is a novel part in the following formalised description of the closed loop depicted in Figure 3.1

$$\begin{aligned} & \text{given DM elements} \\ & \left\{ \begin{array}{c} \text{observed data} \\ \alpha \subseteq \mathbf{a}, \mathbf{d}, a_1, \dots, a_{t-1}, d_1, \dots, d_{t-1} \\ \theta, M(a|d, \theta), F_{t-1}(\theta), \theta, s(d|\theta), f_{t-1}(\theta) \end{array} \right\} \Rightarrow \left\{ \begin{array}{c} M_t(a|d) \\ s_t(d) \\ l_t(a, d) \end{array} \right\} \Rightarrow S_t^o \\ & \Rightarrow d_t \in \mathbf{d} \quad \text{environment} \quad a_t \in \mathbf{a} \Rightarrow \left\{ \begin{array}{c} F_t(\theta) \\ f_t(\theta) \end{array} \right\}, \quad t \in \mathbf{t}. \end{aligned} \quad (3.11)$$

3.4 Conceptual Solution of Preference Elicitation

The proposed preference elicitation and the optimised learning strategy form a unity described by the following conceptual algorithm. Section 3.5 contains discussion providing details on the options made.

Algorithm 1

1. Delimit the time-invariant DM elements listed in Agreements 1, 2 and Assumption 2:
 - a. Specify the set \mathbf{d} of available decisions and the set $\mathbf{a} = \times_{i \in \mathbf{i}} \mathbf{a}_i$ of the multivariate attributes $a = a = (a_i)_{i=1}^{|\mathbf{i}|}$.
 - b. Select the parametric models of the environment $M(a|d, \theta)$ and of the optimal strategy $s(d|\theta)$.
 - c. Specify the prior pds $F_0(\theta)$, $f_0(\theta)$ of the parameters $\theta \in \Theta$, $\theta \in \theta$.

Further on, the algorithm runs for the increasing time $t \in \mathbf{t}$.

2. Specify DM preferences on attributes $(a_i)_{i \in \mathbf{i}}$ via the sets α_i , (3.8), as required by Assumption 1. If possible, specify their relative importance by assigning them weights w_i in probabilistic simplex \mathbf{w} or set $w_i = 1/|\mathbf{i}|$, $i \in \mathbf{i}$. Change the obsolete pd $f_{t-1}(\theta)$ from the previous time if the participant's DM preferences have been changed at this time instance.

This step completes specification of DM elements coinciding with the collection of formal objects in the first curly brackets before the rightwards double arrow, see (3.11).

¹⁷ The partial specification of the DM preferences via Assumptions 1, 2 is much easier than a direct specification of the DM-aims-expressing ideal pds.

3. Evaluate predictive pds, [47], $M_t(a|d)$, $s_t(d)$ serving as the environment model and the optimal-strategy model.

The models $M_t(a|d)$, $s_t(d)$, serve for a design of S_t (3.6) generating the decision d_t . Thus, they can exploit data measured up to and including time $t - 1$, cf. Agreement 1.

4. Select a decision $\bar{d} = \bar{d}(w)$ (it depends on the weights w assigned)

$$\bar{d} = \bar{d}(w) \in \text{Arg max}_{d \in \mathbf{d}} \sum_{i \in \mathbf{i}} w_i \int_{\alpha_i} M_t(a_i|d) da_i \quad (3.12)$$

and define the set \mathbf{I}_t of the reachable ideal pds expressing the participant's DM preferences

$$\mathbf{I}_t = \{I_t(a, d) : I_t(a_i) = M_t(a_i|\bar{d}), \forall a_i \in \alpha_i, i \in \mathbf{i}\}. \quad (3.13)$$

The decision $\bar{d} \in \mathbf{d}$ provides the Pareto-optimised probabilities¹⁸

$$\left(\int_{\alpha_1} M_t(a_1|\bar{d}) da_1, \dots, \int_{\alpha_{|i|}} M_t(a_{|i|}|\bar{d}) da_{|i|} \right)$$

of the desirable attribute-entries sets (3.8). The weight w with constant entries $w_i = 1/|i|$ can be replaced by the weight w^o maximising the probability of the set $\alpha = X_{i \in \mathbf{i}} \alpha_i$ of the most desirable attribute values

$$w^o \in \text{Arg max}_{w \in \mathbf{w}} \int_{\alpha} M_t(a|\bar{d}(w)) da,$$

see (3.12).

5. Extend the partial specification $I_t \in \mathbf{I}_t$, (3.13) to the pd $I_t^e(a, d)$ via the following application of the minimum KLD principle

$$I_t^e \in \text{Arg min}_{I_t \in \mathbf{I}_t} \mathbb{D}(I_t \| P_t) \text{ with } P_t(a, d) = M_t(a|d)s_t(d). \quad (3.14)$$

The set \mathbf{I}_t created in Step 4 reflects the participant's DM preferences. The extension to the ideal pd $I_t^e(a, d)$ supposes that s_t is a good guess of the optimal strategy.

This step finishes the specification of the mapping marked by the first rightwards double arrow in (3.11).

6. Perform the FPD (3.7) with the environment model $M_t(a|d)$ and the ideal pd $I_t^e(a, d)$. Then generate d_t according to the mapping S_t^o optimal in the FPD sense (3.7), apply it, and observe a_t .

Enriching the knowledge available makes the solved DM task dynamic one even for the time-invariant parametric environment model. The dynamics is enhanced by the dependence of the used ideal pd on data and time. The unfeasible

¹⁸ A vector function, dependent on a decision, becomes Pareto-optimised if an improvement of any of its entries leads to a deterioration of another one, [46].

optimal design arising for this complex dynamic DM task has to be solved approximately and the approximation should be revised at each time moment.

This step finishes the specification of the mappings symbolised by the second and marked by the third rightwards double arrow in (3.11).

7. Update the pd $F_{t-1}(\Theta) \rightarrow^{(a_t, d_t)} F_t(\Theta)$ in the Bayesian way, i.e. enrich the knowledge about the parameter Θ of the environment model $M(a|d, \Theta)$ by a_t, d_t .

This step is inevitable even when making decisions without the preference elicitation. The updating may include forgetting [37] if the parameter Θ varies.

8. Update the information about the parameter of the model of the optimal strategy, i.e. update $f_{t-1}(\theta) \rightarrow^{(a_t, d_t)} f_t(\theta)$ according to the following weighted version of the Bayes rule¹⁹

$$f_t(\theta) \propto \mathbf{s}^{\phi_t}(d|\theta)f_{t-1}(\theta), \quad \phi_t = \int_{\mathbf{a}} M_{t+1}(a|d = d_t) da. \quad (3.15)$$

This data censoring is inevitable for learning the optimal strategy.

The step finishes the specification of the mapping expressed by the last rightwards double arrow in (3.11).

9. Increase time t and go to Step 2 or to Step 3, if the DM preferences have not been changed.

3.5 Individual Steps of the Conceptual Solution

This section provides details and discussion of the solution steps. The following sections correspond to the individual steps of the conceptual solution summarised in Section 3.4. The third digit of a section coincides with the corresponding number of the discussed step. Steps 2, 4, 5, 6 and 8 are the key ones, the remaining are presented for completeness.

The general solution is specialised to the important parametric models from the exponential family [2] used in the majority of practically feasible solutions.

3.5.1 Specification of Decision Making Elements

This step transforms a real-life DM task formulated in domain-oriented terms into the formal one.

- ad 1a Specification of the sets of available decisions \mathbf{d} and observable attributes \mathbf{a}

Often, these sets are uniquely determined by the domain-specific conditions of the solved DM task. Possible ambiguities can be resolved by Bayesian testing of hypothesis, e.g. [32], about informativeness of the prospective attributes and about influence of the available decisions.

¹⁹ The environment model $M_{t+1}(a|d = d_t)$ used in (3.15) exploits data measured up to and including time t and will also serve for the choice of d_{t+1} .

ad 1b Choice of the parametric models $M(a|d, \theta)$, $\theta \in \Theta$, $s(d|\theta)$, $\theta \in \theta$

A full modelling art, leading to grey- or black-box models, can be applied here. The modelling mostly provides deterministic but approximate models, which should be extended to the needed probabilistic models via the minimum KLD principle, Section 3.2.

Illustrative simulations, Section 3.6, use zero-order Markov chains that relate discrete-valued attributes and decisions. Markov chains belong to a dynamic *exponential family* (EF)²⁰, [2, 26],

$$M(a_t|a_1, \dots, a_{t-1}, d_1, \dots, d_t, \theta) = M(a_t|d_t, \theta) = \exp \langle \mathbf{B}(a_t, d_t), \mathbf{C}(\theta) \rangle, \quad (3.16)$$

determined by a scalar product $\langle \mathbf{B}, \mathbf{C} \rangle$ of compatible values of multivariate functions $\mathbf{B}(a, d)$ ²¹ and $\mathbf{C}(\theta)$.

The following formula provides the considered most general parametrisation of Markov-chain models and its correspondence with the EF

$$\begin{aligned} M(a_t|d_t, \theta) &= \theta(a_t|d_t) = \exp \left(\sum_{a \in \mathbf{a}, d \in \mathbf{d}} \delta(a, a_t) \delta(d, d_t) \ln(\theta(a|d)) \right) \\ &= \exp \langle \mathbf{B}(a_t, d_t), \mathbf{C}(\theta) \rangle, \quad \theta \in \Theta, \\ s(d_t|\theta) &= \exp \left(\sum_{d \in \mathbf{d}} \delta(d, d_t) \ln(\theta(d)) \right) = \exp \langle \mathbf{b}(d_t), \mathbf{c}(\theta) \rangle, \quad \theta \in \theta, \end{aligned} \quad (3.17)$$

where \mathbf{b} , \mathbf{c} have the same meaning as \mathbf{B} , \mathbf{C} in (3.16). Θ , θ are in appropriate probabilistic simplex sets and *Kronecker delta* $\delta(x, \tilde{x})$ equals to one for $x = \tilde{x}$ and it is zero otherwise.

ad 1c Specification of prior pds $F_0(\theta)$, $f_0(\theta)$, $\theta \in \Theta$, $\theta \in \theta$

The specification of the prior pds is known as knowledge elicitation. A relatively complete elaboration of the automated knowledge elicitation (focused on the dynamic EF) is in [24]. Comparing to a general case, the knowledge elicitation problem within the EF is much simpler as the EF admits a finite-dimensional sufficient statistic [34] and possesses a conjugate prior pd [4]. It has the following self-reproducing functional form

$$F_0(\theta) \propto \exp \langle V_0, \mathbf{C}(\theta) \rangle \chi_{\Theta}(\theta), \quad f_0(\theta) \propto \exp \langle v_0, \mathbf{c}(\theta) \rangle \chi_{\theta}(\theta) \quad (3.18)$$

where an *indicator function* $\chi_x(x)$ equals one for $x \in \mathbf{x}$ and is zero otherwise. In the EF, the knowledge elicitation reduces to a selection of finite-dimensional tables V_0, v_0 parameterising these prior pds.

²⁰ Models from the EF dominate in practice.

²¹ The dependence only on d_t formalises the assumed static case. In a fully dynamic case, the function \mathbf{B} acts on values of a finite-dimensional function of the observed data and of the current decision. Its dimension is fixed even for growing number of collected data and its values can be updated recursively.

For Markov chains, the conjugate priors are Dirichlet pds, cf. (3.17). Their functional forms (chosen for $t = 0$) are preserved for all $t \in \mathbf{t}$

$$\begin{aligned} F_t(\boldsymbol{\theta}) &= \prod_{d \in \mathbf{d}} \frac{\prod_{a \in \mathbf{a}} \Theta(a|d)^{V_t(a|d)-1}}{\beta(V_t(\cdot|d))} \chi_{\boldsymbol{\theta}}(\boldsymbol{\theta}) = \prod_{d \in \mathbf{d}} \text{Di}_{\Theta(\cdot|d)}(V_t(\cdot|d)) \\ V_t(\cdot|d) &= (V_t(a_1|d), \dots, V_t(a_{|\mathbf{i}}|d)), \quad \boldsymbol{\theta}(\cdot|d) = (\theta(a_1|d), \dots, \theta(a_{|\mathbf{i}}|d)), \quad a_i \in \mathbf{a}_i, \\ \mathbf{f}_t(\boldsymbol{\theta}) &= \text{Di}_{\theta}(v_t), \quad V_0(a|d) > 0, \quad v_0(d) > 0 \quad \text{on } \mathbf{a}, \mathbf{d}. \end{aligned} \quad (3.19)$$

The used multivariate beta function

$$\beta(x) = \frac{\prod_{l \in \mathbf{l}} \Gamma(x_l)}{\Gamma(\sum_{l \in \mathbf{l}} x_l)}$$

is defined for a positive $|\mathbf{l}|$ -vector x . $\Gamma(\cdot)$ is Euler gamma function [1]. V_0 can be interpreted as an *occurrence table*: $V_0(a|d)$ means the (thought) number of occurrences of the value a following the decision d . v_0 has a similar interpretation.

3.5.2 Specification of Preferences and Their Changes

The domain-specific description of the participant's DM preferences via the set $(\boldsymbol{\alpha}_i)_{i \in \mathbf{i}}$ of the most desirable attribute values (3.8) is a relatively straightforward task. The considered entry-wise specification respects limits of the human being, who can rarely go beyond pair-wise comparison. The DM preferences specified by (3.8) mean

$$\boldsymbol{\alpha}_i \text{ is the more desirable than } (\mathbf{a}_i \setminus \boldsymbol{\alpha}_i), \quad i \in \mathbf{i}, \quad (3.20)$$

where $\boldsymbol{\alpha}_i \subsetneq \mathbf{a}_i$ is a set of desirable values of the i th attribute entries and $\mathbf{a}_i \setminus \boldsymbol{\alpha}_i$ is its complement to \mathbf{a}_i . The specification of the DM preferences (3.20) is mostly straightforward. However, the participant can change them at some time t , for instance, by changing the selection of the attribute entries non-trivially constrained by the set $\boldsymbol{\alpha} = \times_{i \in \mathbf{i}} \boldsymbol{\alpha}_i$. Let us discuss how to cope with a DM preference change from $\boldsymbol{\alpha}$ to $\tilde{\boldsymbol{\alpha}} \neq \boldsymbol{\alpha}$.

The discussed change modifies the set of candidates of the optimal strategy and makes the learnt pd $\mathbf{f}_{t-1}(\boldsymbol{\theta})$ inadequate. It is possible to construct a new pd $\tilde{\mathbf{f}}_{t-1}(\boldsymbol{\theta})$ from scratch if the observed data is stored. It requires a sufficient supply of deliberation resources for performing a completely new censoring of the observed data, see Section 3.5.8, respecting the new DM preferences given by the set $\tilde{\boldsymbol{\alpha}} \subsetneq \boldsymbol{\alpha}$.

The situation is more complex if the pd $\mathbf{f}_{t-1}(\boldsymbol{\theta})$, $\boldsymbol{\theta} \in \boldsymbol{\theta}$, reflecting the obsolete DM preferences is solely stored. Then, the prior pd $\mathbf{f}_0(\boldsymbol{\theta})$ is the only safe guess of the parameter $\boldsymbol{\theta} \in \boldsymbol{\theta}$, which should describe a new optimal strategy. We hope that the divergence of the correct pd $\tilde{\mathbf{f}}_{t-1}(\boldsymbol{\theta})$ (describing the strategy optimal with respect to $\tilde{\boldsymbol{\alpha}}$) on $\mathbf{f}_{t-1}(\boldsymbol{\theta})$ is bounded. This motivates the choice of the pd $\tilde{\mathbf{f}}_{t-1}(\boldsymbol{\theta})$ via the following version of the minimum KLD principle

$$\tilde{\mathbf{f}}_{t-1}(\cdot) = \mathbf{f}_{t-1}(\cdot|\phi) \in \text{Arg min}_{\tilde{\mathbf{f}}_{t-1}} \mathbb{D}(\tilde{\mathbf{f}}|\mathbf{f}_0), \quad \tilde{\mathbf{f}}_{t-1} = \left\{ \text{pds } \tilde{\mathbf{f}}(\cdot) \text{ on } \boldsymbol{\theta}, \mathbb{D}(\tilde{\mathbf{f}}|\mathbf{f}_{t-1}) \leq k \right\} \quad (3.21)$$

for some $k \geq 0$. The solution of (3.21), recalled in Section 3.2, provides the following rule for tracking of DM preference changes.

The change of the most desirable set α to $\tilde{\alpha} \neq \alpha$ is respected by the change ²²

$$f_{t-1}(\theta) \rightarrow f_{t-1}(\theta|\phi) \propto f_{t-1}^\phi(\theta) f_0^{1-\phi}(\theta). \quad (3.22)$$

The adequate forgetting factor $\phi \in [0, 1]$ is unknown due to the lack of the knowledge of k , which depends on the sets α and $\tilde{\alpha}$ in too complex way. Each specific choice of the forgetting factor ϕ provides a model

$$s_t(d|\phi) = \int_{\Theta} s(d|\theta) f_{t-1}(\theta|\phi) d\theta \quad (3.23)$$

of the optimal strategy. Each model (3.23) determines the probability of a new set $\tilde{\alpha}$ of the most desirable attribute values, which should be maximised by the optimal strategy. This leads to the choice of the best forgetting factor ϕ^o as the maximiser of this probability

$$\phi^o \in \text{Arg max}_{\phi \in [0,1]} \int_{\tilde{\alpha}} M_t(a|d) s_t(d|\phi) dd. \quad (3.24)$$

Qualitatively this solution is plausible due to a smooth dependence of $f_{t-1}(\theta|\phi)$ (3.22) on the forgetting factor ϕ . The pd $f_{t-1}(\theta|\phi)$ has also desirable limit versions, which for: i) $\phi \approx 0$ describe that the optimal strategies corresponding to α and $\tilde{\alpha}$ are unrelated; ii) $\phi \approx 1$ express that the DM preference change of α to $\tilde{\alpha}$ has a negligible influence on the optimal strategy.

Quantitative experience with this tracking of the DM preference changes is still limited but no conceptual problems are expected. Unlike other options within the overall solution, the choice (3.24) is a bit of ad-hoc nature.

3.5.3 Evaluation of Environment and Optimal Strategy Models

The admissible strategies can at most use the knowledge available, Agreement 1. They cannot use correct values of parameters Θ, θ , i.e. they have to work with predictive pds serving as the environment model $M_t(a|d)$ and the model of the optimal strategy $s_t(d)$ ²³. If the DM preferences have changed, $s_t(d)$ should be replaced by the pd $s_t(d|\phi)$ reflecting this change (3.23) with the best forgetting factor $\phi = \phi^o$ (3.24)

²² For pds $f_{t-1}(\theta), f_0(\theta)$ conjugated to an EF member, the pd $f_{t-1}(\theta|\phi)$ (3.22) is also conjugated to it.

²³ Let us recall that all DM tasks work with the same time-invariant parametric models of the environment and the optimal strategy $M(a|d, \Theta)$ and $s(d|\theta)$. The predictive pds $M_t(a|d), s_t(d)$, serving the t th DM task, exploit the knowledge accumulated before time t quantified by $F_{t-1}(\Theta)$ and $f_{t-1}(\theta)$.

$$\begin{aligned} M_t(a|d) &= \int_{\Theta} M(a|d, \Theta) F_{t-1}(\Theta) d\Theta > 0 \text{ on } (\mathbf{a}, \mathbf{d}) \\ s_t(d) &= \int_{\Theta} s(d|\theta) f_{t-1}(\theta) d\theta \text{ or } s_t(d|\phi) = \frac{\int_{\Theta} s(d|\theta) f_{t-1}^{\phi}(\theta) f_0^{1-\phi}(\theta) d\theta}{\int_{\Theta} f_{t-1}^{\phi}(\theta) f_0^{1-\phi}(\theta) d\theta}. \end{aligned} \quad (3.25)$$

The formulae (3.25) can be made more specific for the exponential family (3.16) and for the corresponding conjugate pds (3.18). The self-reproducing property of the pd

$$F_{t-1}(\Theta) = F(\Theta|V_{t-1}) \propto \exp\langle V_{t-1}, \mathbf{C}(\Theta) \rangle \quad (3.26)$$

conjugated to a parametric environment model $M(a|d, \Theta) = \exp\langle \mathbf{B}(a, d), \mathbf{C}(\Theta) \rangle$ and the parametric model of the optimal strategy $s(d|\theta) = \exp\langle \mathbf{b}(d), \mathbf{c}(\theta) \rangle$ imply

$$\begin{aligned} M_t(a|d) &= \frac{J(V_{t-1} + \mathbf{B}(a, d))}{J(V_{t-1})}, \quad J(V) = \int_{\Theta} \exp\langle V, \mathbf{C}(\Theta) \rangle F(\Theta|V) d\Theta \\ s_t(d) &= \frac{j(v_{t-1} + \mathbf{b}(d))}{j(v_{t-1})}, \quad j(v) = \int_{\Theta} \exp\langle v, \mathbf{c}(\theta) \rangle f(\theta|v) d\theta. \end{aligned} \quad (3.27)$$

The stabilised forgetting (3.5), suitable also for tracking of the varying parameter of the environment model, replaces the sufficient statistics V_{t-1}, v_{t-1} by the convex combinations

$$\tilde{V}_{t-1} = \phi V_{t-1} + (1 - \phi)V_0, \quad \tilde{v}_{t-1} = \phi v_{t-1} + (1 - \phi)v_0.$$

For Markov chains and the parametrisation (3.17), the conjugate Dirichlet pds (3.19)

$$F_{t-1}(\Theta) = F(\Theta|V_{t-1}) = \prod_{d \in \mathbf{d}} \text{Di}_{\Theta(\cdot|d)}(V_{t-1}(\cdot|d))$$

and

$$f_{t-1}(\theta) = f(\theta|v_{t-1}) = \text{Di}_{\theta}(v_{t-1})$$

reproduce. They depend on the occurrence tables V_{t-1}, v_{t-1} , which sufficiently compress the observed (censored) data. The corresponding explicit forms of predictive pds serving as the environment model and the model of the optimal strategy are

$$M_t(a|d) = \frac{V_{t-1}(a|d)}{\sum_{\tilde{a} \in \mathbf{a}} V_{t-1}(\tilde{a}|d)}, \quad a \in \mathbf{a}, \quad s_t(d) = \frac{v_{t-1}(d)}{\sum_{\tilde{d} \in \mathbf{d}} v_{t-1}(\tilde{d})}, \quad d \in \mathbf{d}. \quad (3.28)$$

Up to the important influence of initial values of occurrence tables, these formulae coincide with the relative frequency of occurrences of the realised configurations of a specific attribute a after making a specific decision d and the relative frequency of occurrences of the decision value d . The formulas (3.28) follow from the known property $\Gamma(x+1) = x\Gamma(x)$ [1], which also implies that the environment model coincides with conditional expectations $\hat{\Theta}_{t-1}(a|d) = \mathbb{E}[\Theta(a|d)|V_{t-1}]$ of $\Theta(a|d)$

$$M_t(a|d) = \hat{\Theta}_{t-1}(a|d) = \frac{V_{t-1}(a|d)}{\sum_{\tilde{a} \in \mathbf{a}} V_{t-1}(\tilde{a}|d)}, \quad a \in \mathbf{a}, d \in \mathbf{d}. \quad (3.29)$$

It suffices to store the point estimates $\hat{\Theta}_{t-1}(a|d)$ of the unknown parameter $\Theta(a|d)$ and the *imprecision vector* κ_{t-1} with entries

$$\kappa_{t-1}(d) = \frac{1}{\sum_{a \in \mathbf{a}} V_{t-1}(a|d)}, \quad d \in \mathbf{d}, \quad (3.30)$$

instead of V_{t-1} as the following transformations are bijective

$$V \leftrightarrow (\hat{\Theta}, \kappa) \text{ and } v \leftrightarrow (\hat{\theta}, \tau) = \left(\frac{v_d}{\sum_{d \in \mathbf{d}} v_d}, \frac{1}{\sum_{d \in \mathbf{d}} v_d} \right). \quad (3.31)$$

These transformations of sufficient statistics suit for the approximate design of the optimal strategy, see Section 3.5.6.

3.5.4 Expressing of Preferences by Ideal Probability Densities

Within the set of pds on (\mathbf{a}, \mathbf{d}) , we need to select an ideal pd, which respects the participant's partially specified DM preferences. We use the following way.

Find pds assigning the highest probability to the most desirable attribute values according to the environment model when using a proper decision $\bar{d} \in \mathbf{d}$. Then choose the required singleton among them via the minimum KLD principle.

This verbal description delimits the set of ideal-pd candidates both with respect to their functional form and numerical values but does not determine them uniquely. Here, we discuss considered variants of their determination and reasons, which led us to the final choice presented at the end of this section.

Throughout several subsequent sections the time index t is fixed and suppressed as uninformative.

Independent Choice of Marginal Probability Densities of the Ideal Probability Density

Let us consider the i th attribute entry a_i and find the decision $\bar{d}^i \in \mathbf{d}$

$$\bar{d}^i \in \text{Arg max}_{d \in \mathbf{d}} \int_{\alpha_i} M(a_i|d) da_i, \quad i \in \mathbf{i}.$$

The maximised probability of the set α_i of the most desirable values of the i th attribute is given by the i th marginal pd $M(a_i|d)$ of the joint pd $M(a|d)$. The decision \bar{d}^i guarantees the highest probability of having the i th attribute entry in the set α_i at the price that a full decision effort concentrates on it. This motivates to consider a set of ideal pds \mathbf{l} respecting the participant's DM preferences (3.8) as the pds $l(a, d)$ having the marginal pds

$$l(a_i) = M(a_i|\bar{d}^i), \quad i \in \mathbf{i}. \quad (3.32)$$

Originally, we have focused on this option, which respects entry-wise specification of the DM preferences (3.8). It also allows a simple change of the number of attribute entries that actively delimit the set of the most desirable attribute values. Moreover, this specification of \mathbf{I} uses the marginal pds $M(a_i|d)$ of the pd $M(a|d)$, which are more reliable than the full environment model $M(a|d)$. Indeed, the learning of mutual dependencies of respective attribute entries is data-intensive as the number of “dependencies” of $|\mathbf{i}|$ -dimensional discrete-valued vector a grows very quickly.

A closer look, however, reveals that such ideal pds are unrealistic as generally $\bar{d}^i \neq \bar{d}^j$ for $i \neq j$. Then, the ideal pd $l(a)$ with the marginal pds (3.32) cannot be attained. This feature made us to abandon this option.

Joint Choice of the Marginal Probability Densities of the Ideal PD

The weakness of the independent choice has guided us to consider the following option. A single decision $\bar{d} \in \mathbf{d}$ is found, which maximises the probability of reaching the set of the most desirable attribute values α

$$\bar{d} \in \text{Arg max}_{d \in \mathbf{d}} \int_{\alpha} M(a|d) da. \quad (3.33)$$

It serves for constraining the set of ideal pds \mathbf{I} to those having the marginal pds

$$l(a_i) = M(a_i|\bar{d}) = \int_{\mathbf{a}_{\setminus i}} M(a|\bar{d}) da_{\setminus i}, \quad (3.34)$$

where subscript $\setminus i$ indexes a vector created from a by omitting the i th entry

$$\mathbf{a}_{\setminus i} = (a_1, \dots, a_{i-1}, a_{i+1}, \dots, a_{|\mathbf{i}|}).$$

This variant eliminates drawback of the independent choice at the price of using joint pd in (3.33). Otherwise, it seemingly meets all requirements on the DM preferences-expressing ideal pds. It may, however, be rather bad with respect to the entry-wise specified DM preferences (3.8). Indeed, the marginalisation (3.34) may provide an ideal pd with marginal probabilities of the sets α_i (i.e. of the most desirable attribute values) are strictly smaller than their complements

$$\int_{\alpha_i} l(a_i) da_i < \int_{\mathbf{a}_i \setminus \alpha_i} l(a_i) da_i, \quad \forall i \in \mathbf{i}. \quad (3.35)$$

This contradicts the participant’s wish (3.20). The following example shows that this danger is real one.

Example 1

Let us consider a two-dimensional attribute $a = (a_1, a_2)$ and a scalar decision $d = d_1$ with the sets of possible values $\mathbf{a}_1 = \mathbf{a}_2 = \mathbf{d} = \{1, 2\}$. The environment model $M(a|d)$ is a table, explicitly given in Table 3.1 together with the marginal pds $M_1(a_1|d)$, $M_2(a_2|d)$.

Table 3.1 The left table describes the discussed environment model $M(a|d)$, the right one provides its marginal pds $M(a_i|d)$. The parameters in this table have to meet constraints $\sigma_i, \rho_i, \zeta_i > 0$, $1 - \sigma_i - \rho_i - \zeta_i > 0$, $i \in \mathbf{i} = \{1, 2\}$ guaranteeing that $M(a|d)$ is a conditional pd.

	$a=(1,1)$	$a=(1,2)$	$a=(2,1)$	$a=(2,2)$		$a_1=1$	$a_1=2$	$a_2=1$	$a_2=2$
$d=1$	σ_1	ρ_1	ζ_1	$1 - \sigma_1 - \rho_1 - \zeta_1$	$d=1$	$\sigma_1 + \rho_1$	$1 - \sigma_1 - \rho_1$	$\sigma_1 + \zeta_1$	$1 - \sigma_1 - \zeta_1$
$d=2$	σ_2	ρ_2	ζ_2	$1 - \sigma_2 - \rho_2 - \zeta_2$	$d=2$	$\sigma_2 + \rho_2$	$1 - \sigma_2 - \rho_2$	$\sigma_2 + \zeta_2$	$1 - \sigma_2 - \zeta_2$

For an entry-wise-specified set of the most desirable attribute values $\alpha_1 = \{1\}$, $\alpha_2 = \{1\}$, we get the set α of the most desirable attribute values as the singleton $\alpha = \{(1, 1)\}$. Considering $\sigma_1 > \sigma_2$ in Table 3.1, the decision $\bar{d} = 1$ maximises $\int_{\alpha} M(a|d) da = M(1, 1|d)$. There is an infinite amount of the parameter values $\sigma_i, \rho_i, \zeta_i$, $i \in \mathbf{i}$, for which $M(a|\bar{d})$ has the marginal pds (3.34) with the adverse property (3.35). A possible choice of this type is in Table 3.2.

Table 3.2 The left table contains specific numerical values of the discussed environment model $M(a|d)$, the right one provides its marginal pds $M(a_i|d)$, $i \in \mathbf{i} = \{1, 2\}$

	$a=(1,1)$	$a=(1,2)$	$a=(2,1)$	$a=(2,2)$		$a_1=1$	$a_1=2$	$a_2=1$	$a_2=2$
$d=1$	0.40	0.05	0.05	0.50	$d=1$	0.45	0.55	0.45	0.55
$d=2$	0.30	0.30	0.30	0.10	$d=2$	0.60	0.40	0.60	0.40

Table 3.2 shows that the decision $\bar{d} = 1$, maximising $\int_{\alpha} M(a|d) da = M((1, 1)|d)$, gives the marginal pds $M(a_i = 1|\bar{d} = 1) = 0.45 < 0.55 = M(a_i = 2|\bar{d} = 1)$. The other decision $d^o = 2$ leads to $M(a_i = 1|d^o = 2) = 0.6 > 0.4 = M(a_i = 2|d^o = 2)$, both for $i = 1, 2$. This property disqualifies the choice (3.33), (3.34).

Pareto-Optimal Marginal Probability Densities of the Ideal PD

The adverse property (3.35) means that the solution discussed in the previous section can be dominated in Pareto sense [46]: the marginal pds $l(a_i)$ (3.34) may lead to probabilities $\int_{\alpha_i} l(a_i) da_i$, $i \in \mathbf{i}$, which are smaller than those achievable by other decision $d^o \in \mathbf{d}$ used in (3.34) instead of \bar{d} . This makes us to search directly for a non-dominated, Pareto optimal, solution reachable by a $\bar{d} \in \mathbf{d}$.

Taking an $|\mathbf{i}|$ -dimensional vector w of arbitrary positive probabilistic weights $w \in \mathbf{w} = \{w = (w_1, \dots, w_{|\mathbf{i}|}), w_i > 0, \sum_{i \in \mathbf{i}} w_i = 1\}$ and defining the w -dependent decision

$$\bar{d} \in \text{Arg max}_{d \in \mathbf{d}} \sum_{i \in \mathbf{i}} w_i \int_{\alpha_i} M(a_i|d) da_i \text{ and } l(a_i) = M(a_i|\bar{d}), \quad i \in \mathbf{i}, \quad (3.36)$$

ensures the found solution be non-dominated.

Indeed, let us assume that there is another $d^o \in \mathbf{d}$ such that $\int_{\alpha_i} M(a_i|d^o) da_i \geq \int_{\alpha_i} M(a_i|\bar{d}) da_i, \forall i \in \mathbf{i}$, with some inequality being strict. Multiplying these inequalities by the positive weights w_i and summing them over $i \in \mathbf{i}$ we get the sharp inequality contradicting the definition of \bar{d} as the maximiser in (3.36).

The possible weights w : i) are either determined by the participant if it is able to distinguish the importance of individual attribute entries; ii) or are fixed to the constant $w_i = 1/|\mathbf{i}|$ if the participant is indifferent with respect to individual attribute entries; iii) or can be chosen as maximiser of the probability of a set of the most desirable attribute values $\alpha = X_{i \in \mathbf{i}} \alpha_i$ by selecting

$$w^o \in \text{Arg max}_{w \in \mathbf{w}} M(X_{i \in \mathbf{i}} \alpha_i | \bar{d}), \text{ with } w\text{-dependent } \bar{d} \text{ given by (3.36)}. \quad (3.37)$$

The primary option i) is rarely realistic and the option iii) is to be better than ii).

In summary, the most adequate set of ideal pds respecting the participant's DM preferences (3.8) in a reachable way reads

$$\mathbf{I} = \{l(a, d) > 0 \text{ on } (\mathbf{a}, \mathbf{d}) : l(a_i) = M(a_i | \bar{d}), \text{ with } \bar{d} \text{ given by (3.36)}\}. \quad (3.38)$$

Note that in the example of the previous section the optimisation (3.37) is unnecessary. General case has not been analysed yet but no problems are foreseen.

A natural question arises: Why the decision \bar{d} is not directly used as the optimal one? The negative answer follows primarily from heavily exploitation-oriented nature of \bar{d} . It does not care about exploration, which is vital for a gradual improvement of the used strategy, which depends on the improvement of the environment and optimal strategy models. In the considered static case, the known fact that the repetitive use of \bar{d} may completely destroy learning and consequently decision making [39] manifests extremely strongly. The example of such an adverse behaviour presented in Section 3.6 appeared without any special effort.

Moreover in dynamic DM, the use of this myopic²⁴ strategy leads to an inefficient behaviour, which even may cause instability of the closed decision loop [30].

3.5.5 Extension of Marginal Probability Densities to the Ideal PD

The set \mathbf{I} (3.38) of the ideal pds is determined by linear constraints explicitly expressing the participant's wishes. This set is non-empty as

$$l(a, d) = \tilde{P}(d|a) \prod_{i \in \mathbf{i}} l(a_i) \in \mathbf{I}, \text{ with } l(a_i) = M(a_i | \bar{d}),$$

²⁴ The myopic strategy is mostly looking one-stage-ahead.

where $\tilde{\mathbf{P}}(d|a)$ is an arbitrary pd positive on \mathbf{d} for all conditions in \mathbf{a} . The arbitrariness reflects the fact known from the copula theory [44] that marginal pds do not determine uniquely the joint pd having them. Thus, \mathbf{I} contains more members and an additional condition has to be adopted to make the final well-justified choice. The selection should not introduce extra, participant-independent, DM preferences. The minimum KLD principle, Section 3.2, has this property if it selects the ideal pd $l(a, d)$ from (3.38) as the closest one to the pd $\mathbf{P}(a, d)$ expressing the available knowledge about the closed decision loop with the optimal strategy.

The minimised KLD $\mathbb{D}(\mathbf{I}|\mathbf{P})$ is a strictly convex functional of the pd l from the convex set \mathbf{I} (3.38) due to the assumed positivity of the involved pds $\mathbf{M}(a|d)$ and $\mathbf{P}(a, d)$. Thus, see [18], the constructed l is a unique minimiser of the Lagrangian functional determined by the multipliers $-\ln(\Lambda(a_i))$, $i \in \mathbf{i}$,

$$\begin{aligned} l^e(a, d) &= \arg \min_{l \in \mathbf{I}} \int_{(\mathbf{a}, \mathbf{d})} \left[l(a, d) \ln \left(\frac{l(a, d)}{\mathbf{P}(a, d)} \right) - \sum_{i \in \mathbf{i}} \ln(\Lambda(a_i)) l(a, d) \right] da dd \\ &= \arg \min_{l \in \mathbf{I}} \int_{\mathbf{a}} l(a) \left[\ln \left(\frac{l(a)}{\mathbf{P}(a) \prod_{i \in \mathbf{i}} \Lambda(a_i)} \right) + \int_{\mathbf{d}} l(d|a) \ln \left(\frac{l(d|a)}{\mathbf{P}(d|a)} \right) dd \right] da \\ &= \frac{\mathbf{P}(d|a) \mathbf{P}(a) \prod_{i \in \mathbf{i}} \Lambda(a_i)}{\mathbf{J}(\mathbf{P}, \Lambda)}. \end{aligned} \quad (3.39)$$

The second equality in (3.39) is implied by Fubini theorem [48] and by the chain rule for pds [47]. The result in the third equality follows from the facts that: i) the second term (after the second equality) is a conditional version of the KLD, which reaches its smallest zero value for $l^e(d|a) = \mathbf{P}(d|a)$ and ii) the first term is the KLD minus logarithm of an l -independent normalising constant $\mathbf{J}(\mathbf{P}, \Lambda)$.

The Lagrangian multipliers solve, for $a_{\setminus i} = (a_1, \dots, a_{i-1}, a_{i+1}, \dots, a_{|\mathbf{i}|})$, $\Lambda_{\setminus i} = (\Lambda_1, \dots, \Lambda_{i-1}, \Lambda_{i+1}, \dots, \Lambda_{|\mathbf{i}|})$, $i \in \mathbf{i}$,

$$\begin{aligned} \mathbf{M}(a_i | \vec{d}) &= \frac{\int_{\mathbf{a}_{\setminus i}} \mathbf{P}(a) \prod_{j \in \mathbf{i}} \Lambda_j(a_j) da_{\setminus i}}{\mathbf{J}(\mathbf{P}, \Lambda)} \\ &= \frac{\mathbf{P}_i(a_i) \Lambda_i(a_i)}{\mathbf{J}(\mathbf{P}, \Lambda)} \int_{\mathbf{a}_{\setminus i}} \mathbf{P}(a_{\setminus i} | a_i) \prod_{j \in \mathbf{i} \setminus \{i\}} \Lambda_j(a_j) da_{\setminus i} = \frac{\mathbf{P}_i(a_i) \Lambda_i(a_i)}{\mathbf{J}(\mathbf{P}, \Lambda)} \Phi(a_i, \Lambda_{\setminus i}). \end{aligned} \quad (3.40)$$

By construction, equations (3.40) have a unique solution, which can be found by the successive approximations

$$\frac{{}^k \Lambda_i(a_i)}{\mathbf{J}(\mathbf{P}, {}^k \Lambda)} = \frac{\mathbf{M}(a_i | \vec{d})}{\mathbf{P}_i(a_i) \Phi(a_i, {}^{k-1} \Lambda_{\setminus i})}, \quad (3.41)$$

where an evaluation of the k th approximation ${}^k \Lambda = [{}^k \Lambda_1, \dots, {}^k \Lambda_{|\mathbf{i}|}]$, $k = 1, 2, \dots$ starts from an initial positive guess ${}^0 \Lambda$ and stops after the conjectured stabilisation. The factor $\mathbf{J}(\mathbf{P}, {}^k \Lambda)$ is uniquely determined by the normalisation of the k th approximation of the constructed ideal pd (3.39).

The successive approximations, described by (3.41), do not face problems with division by zero for the assumed $M(a|d), P(a, d) > 0 \Rightarrow M(a_i|\vec{d}), P_i(a_i) > 0, i \in \mathbf{i}$. Numerical experiments strongly support the adopted conjecture about their convergence. If the conjecture is valid then the limit provides the searched solution due to its uniqueness.

3.5.6 Exploring Fully Probabilistic Design

Here, the dynamic exploration aspects of the problem are inspected so that time-dependence is again explicitly expressed. As before, the discussion is preferred against definite results. A practically oriented reader can skip the rest of this section after reading the first section, where the used certainty-equivalent strategy is described.

The repetitive solutions of the same type static DM problems form a single dynamic problem due to the common, recursively-learned, environment and optimal strategy models. The KLD (scaled by $|\mathbf{t}|$) to be minimised over $(S_t)_{t \in \mathbf{t}}$ reads

$$\frac{1}{|\mathbf{t}|} \mathbb{D} \left(\prod_{t \in \mathbf{t}} M_t S_t \middle| \middle| \prod_{t \in \mathbf{t}} l_t \right) = \frac{1}{|\mathbf{t}|} \sum_{t \in \mathbf{t}} \mathbb{E} [\mathbb{D}(M_t S_t \| l_t)].$$

The expectation $\mathbb{E}[\bullet]$ is taken over conditions occurring in the individual summands and

$$\mathbb{D}(M_t S_t \| l_t) = \int_{(a,d)} M_t(a|d) S_t(d) \ln \left(\frac{M_t(a|d) S_t(d)}{l_t(a,d)} \right) da dd. \quad (3.42)$$

Recall, $M_t(a|d) = \int_{\Theta} M(a|d, \theta) F_{t-1}(\theta) d\theta$ and $l_t(a, d)$ is an image of $M_t(a|d)$ and $s_t(d|\phi^o) = \int_{\Theta} s(d|\theta) f_{t-1}(\theta|\phi^o) d\theta$, see Sections 3.5.3, 3.5.5.

The formal solution of this FPD [53] has the following form evaluated backward for $t = |\mathbf{t}|, |\mathbf{t}| - 1, \dots, 1$,

$$\begin{aligned} S_t^o(d) &= \frac{l_t(d) \exp[-\omega_t(d)]}{\int_d l_t(d) \exp[-\omega_t(d)] dd} = \frac{l_t(d) \exp[-\omega_t(d)]}{\gamma(a_1, \dots, a_{t-1}, d_1, \dots, d_{t-1})} \\ \omega_t(d) &= \int_a M_t(a|d) \ln \left(\frac{M_t(a|d)}{\gamma(a_1, \dots, a_{t-1}, a, d_1, \dots, d_{t-1}, d) l_t(a|d)} \right) da \\ &\text{starting from } \gamma(a_1, \dots, a_{|\mathbf{t}|}, d_1, \dots, d_{|\mathbf{t}|}) = 1. \end{aligned} \quad (3.43)$$

The optimal design (3.43) provides a dual strategy [12], which optimally balances exploration and exploitation activities. In the considered static case, the observed data enters the conditions via the learnt models of the environment and the optimal strategy (3.25) determining the used ideal pd $l_t = l_t^e$. The data realisation is influenced by the applied strategy, which can be optimal only if the overall effect of exploitation and exploration is properly balanced. This design is mostly infeasible and some approximate-dynamic-programming technique [52] has to be used.

Certainty Equivalence Decision Strategy

The simplest approximate design, labelled as *certainty equivalence*, replaces the unknown parameter by its current point estimate and assumes that the estimate is uninfluenced by the chosen decisions.

This strategy is asymptotically reasonable as the posterior pds $F_t(\Theta)$, $f_t(\theta)$ form martingales [42] and under the general conditions (met always for Markov chains) they almost surely converge to singular pds [3, 26]. Consequently, the only source of dependence between successive static DM tasks diminishes as the dependence of γ in (3.43) on data disappears. The certainty-equivalent strategy that neglects this dependence breaks the single task (3.42) into the sequence of DM tasks consisting of independently solved one-stage-ahead looking static problems with the solutions of the form (cf. (3.43))

$$S_t^o(d) \propto l_t(d) \exp[-\omega_t(d)], \quad \omega_t(d) = \int_a M_t(a|d) \ln \left(\frac{M_t(a|d)}{l_t(a|d)} \right) da, \quad t \in \mathbf{t}.$$

The transient learning period is, however, critical as – without an exploration – the posterior pds may concentrate on wrong sets, whenever the conditioning data is not informative enough. Here, one of the advantages of the FPD enters the game. The FPD provides a *randomised* optimal strategy and the sampling of decisions from it adds a well-scaled dither (exploring) noise, which diminishes with a proper rate.

Strictly speaking, these statements are conjectures supported by experiments whose samples are in Section 3.6. To get a guaranteed version of sufficiently exploring, almost optimal, strategy, a more sophisticated approximation (3.43) is needed. It requires tailoring of techniques known from approximate dynamic programming [52] to the FPD. It is possible and relatively simple as the strategy optimal in the FPD sense is described explicitly. A widely applicable construction framework is outlined below.

Failing Cautious Decision Strategy

During the whole development, we have been aware (and experiments confirmed it) that exploration is vital for a good overall closed-loop behaviour. Primarily, this made us avoid a direct use of \bar{d} even in the static DM.

For continuous-valued (a, d) , there is an experimental evidence that certainty-equivalent strategy is often exploratory enough even in the standard Bayesian DM. On contrary, in Markov-chain case, there is a proof that deterministic certainty-equivalent strategy may fail with a positive probability [39]. This made us to let one stage-ahead design know that the unknown parameters Θ, θ are in the game. It did not helped as we got cautious strategy [19], which is even less exploring than the certainty equivalent one and worked improperly. The formal mechanism of this failure becomes obvious when noticing that the imprecisions of parameter estimates κ, τ (3.30), (3.31) enter the design only for the design horizon $|\mathbf{t}| > 1$.

Approximate Dynamic Programming in FPD

This section touches of a rich area of approximate dynamic programming [52]. It indicates that within the FPD it is possible to obtain algorithmically feasible approximation of the optimally exploring strategy. The approach is just outlined and thus it suffices to focus on the Markov-chain case to this purpose. The corresponding function $\gamma(\cdot)$, determining the function $\omega(\cdot)$ and the optimal strategy, see (3.43), depends on the sufficient statistics. They consist of the point estimates $\hat{\Theta}_{t-1}, \hat{\theta}_{t-1}$ of the parameters Θ, θ and of the corresponding imprecisions κ_{t-1}, τ_{t-1} , see (3.29), (3.30) and (3.31). Note that the use of this form of sufficient statistics is important as the statistics are expected to converge to finite constants (ideally, $\hat{\Theta}_t, \hat{\theta}_t$ to the correct values of the unknown parameters Θ, θ and imprecisions κ_t, τ_t to zeros) unlike the statistics V_t, v_t .

For this open-ended presentation, it suffices to consider the explicit dependence of $\gamma(\cdot)$ and $\omega(\cdot)$ on $\hat{\Theta}$ and κ only. The first step of the optimal design (3.43) for $t = |t|$ with $\gamma(\cdot) = 1$ gets the form

$$\begin{aligned} \mathbf{S}_t^o(d) &= \mathbf{S}^o(d|V_{t-1}) = \mathbf{S}^o(d|\hat{\Theta}_{t-1}, \kappa_{t-1}) = \frac{l_t(d) \exp[-\omega_t(d, \hat{\Theta}_{t-1}, \kappa_{t-1})]}{\gamma_t(\hat{\Theta}_{t-1}, \kappa_{t-1})} \\ \gamma_t(\hat{\Theta}_{t-1}, \kappa_{t-1}) &= \sum_{d \in \mathbf{d}} l_t(d) \exp[-\omega_t(d, \hat{\Theta}_{t-1}, \kappa_{t-1})] \\ \omega_t(d, \hat{\Theta}_{t-1}, \kappa_{t-1}) &= \sum_{a \in \mathbf{a}} \hat{\Theta}_{t-1}(a|d) \ln \left(\frac{\hat{\Theta}_{t-1}(a|d)}{l_t(a|d)} \right). \end{aligned}$$

Further design steps can be interpreted as value iterations searching for the optimal stationary strategy. Even for the FPD, they converge for large $|t|$ under the rather general conditions [26]. We care about this stationary phase, drop the time subscript at $\gamma_t(\cdot), \omega_t(\cdot)$, and write down a general step of the design in terms of $\gamma(\cdot)$. It holds (cf. Section 3.5.7)

$$\begin{aligned} \gamma(\hat{\Theta}_{t-1}, \kappa_{t-1}) &= \sum_{d \in \mathbf{d}} \left\{ l_t(d) \exp \left[- \sum_{a \in \mathbf{a}} \hat{\Theta}_{t-1}(a|d) \ln \left(\frac{\hat{\Theta}_{t-1}(a|d)}{l_t(a|d)} \right) \right] \right. \\ &\quad \left. \times \exp \left[\sum_{a \in \mathbf{a}} \hat{\Theta}_{t-1}(a|d) \ln (\gamma(\hat{\Theta}_{t-1} + \Delta_{t-1}(\cdot, \cdot, a, d), \kappa_{t-1} \Omega_{t-1}(\cdot, d))) \right] \right\} \\ \Delta_{t-1}(\tilde{a}, \tilde{d}, a, d) &= \frac{\delta(\tilde{d}, d) \kappa_{t-1}(d)}{1 + \delta(\tilde{d}, d) \kappa_{t-1}(d)} (\delta(\tilde{a}, a) - \hat{\Theta}_{t-1}(a|d)), \quad \tilde{a}, a \in \mathbf{a}, \tilde{d}, d \in \mathbf{d} \\ \Omega_{t-1}(\tilde{d}, d) &= \frac{1}{1 + \kappa_{t-1}(\tilde{d}) \delta(\tilde{d}, d)}, \quad \tilde{d}, d \in \mathbf{d}. \end{aligned} \quad (3.44)$$

Let us insert into the right-hand side of $\gamma(\cdot)$ in (3.44) an exponential-type approximation of $\gamma(\cdot) > 0$

$$\gamma(\hat{\Theta}, \kappa) \approx \exp \left[\text{tr}(G\hat{\Theta}) + g\kappa \right]. \quad (3.45)$$

The approximation is parameterised by a matrix G (of the size of the transposed $\hat{\theta}_t$) and by a vector g (of the transposed size of κ_t). It gives a mixture of exponential functions of a similar type on the left-hand side of (3.44). The recovering of the feasible exponential form (3.45) requires an approximation of this mixture by a single exponential function. The key question is what proximity measure should be used. To decide it, it suffices to observe that the optimal strategy $\mathcal{S}^o(d|\hat{\theta}_{t-1}, \kappa_{t-1})$ is the pd of d conditioned on $\hat{\theta}_{t-1}, \kappa_{t-1}$. Thus, $\gamma(\hat{\theta}_{t-1}, \kappa_{t-1})$ can be interpreted as a (non-normalised) marginal pd of $\hat{\theta}_{t-1}, \kappa_{t-1}$, which should be approximated by a feasible pd having the form (3.45). This singles out the approximation in terms of the KLD (3.2). In the considered case, it reduces to a fitting of the moments of $\hat{\theta}_{t-1}, \kappa_{t-1}$ of the left-hands-side mixture by moments of a pd having the exponential form (3.45).

The outlined idea is applicable generally. Algorithmic details for the EF will be elaborated in an independent paper.

3.5.7 Learning of the Environment Model

The environment model serves to the DM task for which the DM preferences are elicited. Thus, it has to be constructed anyway. Within the considered context of repetitive DMs, its standard Bayesian learning is available. It works with an environment model $M(a|d, \theta)$ parameterised by a finite-dimensional unknown parameter $\theta \in \Theta$. The following version of the Bayes rule [4] evolves the pd $F_{t-1}(\theta)$ comprising all information about the unknown parameter θ

$$F_t(\theta) = \frac{M(a_t|d_t, \theta)F_{t-1}(\theta)}{\int_{\Theta} M(a_t|d_t, \theta)F_{t-1}(\theta) d\theta}, \quad (3.46)$$

where the pair (a_t, d_t) is realised in the t th DM task. This pd provides the predictive pd (3.25) used as the environment model $M_{t+1}(a|d)$. It is worth stressing that the applied strategy cancels in the directly applied Bayes rule. Within the considered context, it “naturally” fulfils natural conditions of control (decision making) [47] requiring $S_t(d|\theta) = S_t(d)$ and expressing that the parameter θ is unknown to the used strategy.

For a parametric model from the EF (3.16) and a conjugate self-reproducing pd $F_t(\theta) \propto \exp\langle V_t, \mathbf{C}(\theta) \rangle$, the functional form of the Bayes rule (3.46) reduces to the updating of the sufficient statistic, cf. (3.27),

$$V_t = V_{t-1} + \mathbf{B}(a_t, d_t), \quad V_0 \text{ determines the prior pd and has to make } \mathbf{J}(V_0) < \infty.$$

For the controlled Markov chain used in simulations, V_t is an occurrence table with entries $V_t(a|d)$, $a \in \mathbf{a}$, $d \in \mathbf{d}$. The specific form of the function $\mathbf{B}(a, d)$ (3.17) provides the updating in the form

$$V_t(a|d) = V_{t-1}(a|d) + \delta(a, a_t)\delta(d, d_t), \quad V_0(a|d) > 0 \Leftrightarrow \mathbf{J}(V_0) < \infty, \quad (3.47)$$

for the realised pair (a_t, d_t) and $\delta(a, \tilde{a})$ denoting Kronecker delta. This recursion transforms to the recursion for the point estimates $\hat{\Theta}_t$ (3.29) and the imprecisions κ_t (3.30)

$$\begin{aligned}\hat{\Theta}_t(a|d) &= \hat{\Theta}_{t-1}(a|d) + \delta(d, d_t)\kappa_t(d)(\delta(a, a_t) - \hat{\Theta}_{t-1}(a|d)) \\ \kappa_t(d) &= \frac{\kappa_{t-1}(d)}{1 + \kappa_{t-1}(d)\delta(d, d_t)},\end{aligned}$$

which is used in Section 3.5.6 discussing an approximate dynamic programming.

This learning is conceptually very simple but it is strongly limited by the curse of dimensionality as the involved occurrence tables are mostly too large. Except very short vectors of attributes and decisions with a few possible values, their storing and updating require extremely large memory and, even worse, an extreme number of the observed data. Learning a mixture of low-dimensional approximate environment models relating scalar entries of attributes to scalar entries of the decision, [26, 43], seems to be a systematic viable way around.

Note that if parameter vary either because of physical reasons or due to approximation errors, the pd $F_t(\Theta)$ differs from a correct pd and the situation discussed in connection with changing DM preferences, Section 3.5.2, recurs. Thus, the parameter changes can be respected when complementing the Bayes rule (3.46) by the stabilised forgetting

$$\begin{aligned}F_t(\Theta) &\rightarrow F_t^\phi(\Theta)(F_0(\Theta))^{1-\phi_0} \quad \text{in general case} \\ V_t &\rightarrow \phi V_t + (1 - \phi)V_0 \quad \text{for members of the EF.}\end{aligned}$$

In this context, the forgetting factor $\phi \in [0, 1]$ can be learnt in the usual Bayesian way at least on a discrete grid in $[0, 1]$, see e.g. [36, 40].

3.5.8 Learning of the Optimal Strategy

The construction of the ideal pd $I_t = I_t^e$ strongly depends on availability of the model $P_t(a, d)$ of the closed decision loop with the optimal strategy $P_t(a, d) = M_t(a|d)s_t(d)$, see Section 3.5.5. The Bayes rule is used for learning the environment model, Sections 3.5.7, 3.5.3. This rule could be used for learning the optimal strategy if all past decisions were generated by it. This cannot be expected within the important transient learning period. Thus, we have to decide whether a generated decision comes from the (almost) optimal strategy or not: we have to use a censored data.

If the realised attribute falls in the set of the most desirable attribute values α then we have a strong indicator that the used decision is optimal. When relying only on it, we get an unique form of the learning with the strict data censoring

$$f_t(\theta) = \frac{s(d_t|\theta)^{\chi_{\alpha}(a_t)} f_{t-1}(\theta)}{\int_{\theta} s(d_t|\theta)^{\chi_{\alpha}(a_t)} f_{t-1}(\theta) d\theta}. \quad (3.48)$$

However, the event $a_t \in \boldsymbol{\alpha}$ may be rare or a random consequence of a bad decision within a particular realisation. Thus, an indicator working with “almost optimality” is needed. It has to allow a learning even for $a_t \notin \boldsymbol{\alpha} \Leftrightarrow \chi_{\boldsymbol{\alpha}}(a_t) = 0$. For its design, it suffices to recognise that no censoring can be errorless. Thus, the pd $\mathbf{f}_{t-1}(\theta)$ is an approximate learning result: even if $a_t \in \boldsymbol{\alpha}$, we are uncertain whether the updated pd, labelled $\tilde{\mathbf{f}}_t(\theta) \propto \mathbf{s}(d_t|\theta)\mathbf{f}_{t-1}(\theta)$ coincides with a correct pd $\mathbf{f}_t(\theta)$. In other words, $\tilde{\mathbf{f}}_t$ only approximates the correct pd \mathbf{f}_t . Again, as shown in [5], [28], the KLD is the proper Bayesian expression of their divergence. Thus, $\mathbb{D}(\mathbf{f}_t||\tilde{\mathbf{f}}_t) \leq k_t$ for some $k_t \geq 0$. At the same time, the pd $\mathbf{f}_{t-1}(\theta)$ is the best available guess before processing the realisation (a_t, d_t) . The extension of this knowledge is to be done by the minimum KLD principle, Section 3.2, which provides

$$\mathbf{f}_t(\theta) = \frac{\mathbf{s}^{\phi_t}(d_t|\theta)\mathbf{f}_{t-1}(\theta)}{\int_{\boldsymbol{\theta}} \mathbf{s}^{\phi_t}(d_t|\theta)\mathbf{f}_{t-1}(\theta) d\boldsymbol{\theta}}, \quad \phi_t \in [0, 1]. \quad (3.49)$$

The formula (3.49) resembles (3.48) with the forgetting factor $\phi_t \in [0, 1]$ replacing the value of the indicator function $\chi_{\boldsymbol{\alpha}}(a_t) \in \{0, 1\}$. This resemblance helps us to select the forgetting factor, which is unknown due to the unspecified K_t , as a prediction of the indicator-function value

$$\phi_t = \int_{\boldsymbol{\alpha}} M_{t+1}(a|d_t) da = \int_{\boldsymbol{\alpha}} \int_{\boldsymbol{\theta}} M(a|d_t, \boldsymbol{\theta}) F_t(\boldsymbol{\theta}) d\boldsymbol{\theta} da. \quad (3.50)$$

The use of d_t in the condition complies with checking of the (approximate) optimality of this decision sought before updating. Its use in the condition differentiates the formula (3.50) from (3.24), which cares about an “average” degree of optimality of the past decisions. The current experience with the choice (3.50) is positive but still the solution is of ad-hoc type.

3.6 Illustrative Simulations

All solution steps were tested on simulation experiments, which among others allowed cutting off clear cul-de-sacs of the developed solution. Here, we present a simple illustrative example, which can be confronted with intuitive expectations.

3.6.1 Simulation Set Up

The presentation follows the respective steps of Algorithm 1, see Section 3.4.

1. DM elements

- a. One-dimensional decisions $d \in \mathbf{d} = \{1, 2\}$ and two-dimension observable attributes $(a_1, a_2) \in X_{i \in I} \mathbf{a}_i$, with $\mathbf{a}_1 = \{1, 2\}$, $\mathbf{a}_2 = \{1, 2, 3\}$ coincide with the sets simulated by the environment, see Table 3.3.
- b. Zero-order Markov chains with the general parameterisations (3.17) are used.

- c. The conjugate Dirichlet pds (3.19) are used as priors with $V_0(a|d) = 0.1$ on (\mathbf{a}, \mathbf{d}) , $v_0 = [4, 1]$. The latter choice intentionally prefers a bad strategy. This choice checks learning abilities of the proposed preference elicitation.
2. The sets of the most desirable individual entries are $\alpha_1 = \{1\}$, $\alpha_2 = \{1\}$ giving $\alpha = \{(1, 1)\}$. No change of the DM preferences is assumed.
Further on, the algorithm runs for increasing time $t \in \mathbf{t} = \{1, \dots, 100\}$.
3. The predictive pds serving as the environment model and the model of the optimal strategy are evaluated according to formulae (3.28).
4. The decision \bar{d} (3.12) is evaluated for the uniform weights $w_i = 1/|\mathbf{i}| = 1/2$ reflecting the indifference with respect to the attribute entries.
5. The marginal ideal pds $l_t(a_i) = M_t(a_i|\bar{d})$ are extended to the ideal pd $l_t(a, d) = l_t^e(a, d)$ as described in Section 3.5.5.
6. The FPD is performed in its certainty-equivalent version, see Section 3.5.6. The decision d_t is sampled from the tested strategy and it is fed into the simulated environment described by the transition probabilities in Table 3.3

Table 3.3 The simulated environment with probabilities of the configurations of attributes a responding to decision d being in respective cells. Under the complete knowledge of these probabilities, the optimal strategy selects $d = 2$.

	$a=(1,1)$	$a=(1,2)$	$a=(1,3)$	$a=(2,1)$	$a=(2,2)$	$a=(2,3)$
$d=1$	0.20	0.30	0.10	0.10	0.10	0.20
$d=2$	0.35	0.05	0.05	0.15	0.15	0.25

A fixed seed of random-numbers generator is used in all simulation runs, which makes the results comparable.

7. Bayesian learning of the environment model is performed according to the Bayes rule (3.47) without forgetting.
8. Learning of the optimal strategy runs exactly as proposed in Section 3.5.8.

3.6.2 Simulation Results

The numerical results of experiments show outcomes for: i) the *naive strategy* that exploits directly \bar{d}_t (3.12) as the applied decision; ii) the *optimal strategy* that permanently applies the optimal decision ($d_t = 2$); iii) the *proposed strategy* that samples decisions from the certainty-equivalent result of the FPD.

Table 3.4 summarises relative estimation error of Θ parameterising the environment model

$$(1 - \hat{\Theta}_t(a|d)/\Theta(a|d)) \times 100, a \in \mathbf{a}, d \in \mathbf{d}, \quad (3.51)$$

where $\hat{\Theta}_t(a|d)$ are point estimates (3.29) of simulated transition probabilities $\Theta(a|d)$ listed in Table 3.3. The occurrences of attribute and decision values for the respective tested strategies are in Table 3.5.

Table 3.4 Relative estimation errors [%] (3.51) after simulating $|t| = 100$ samples for the respective tested strategies. The error concerning the most desirable attribute value is in the first column.

<i>Results for the naïve strategy</i>						
$d = 1$	0.9	-69.8	11.9	5.7	5.7	10.4
$d = 2$	52.4	-11.1	-233.3	-11.1	-233.3	33.3
<i>Results for the optimal strategy</i>						
$d = 1$	16.7	-66.7	44.4	-66.7	-66.7	16.7
$d = 2$	-21.3	5.7	5.7	24.5	-32.1	17.0
<i>Results for the proposed strategy</i>						
$d = 1$	23.9	-52.2	13.0	-30.4	-8.7	2.2
$d = 2$	-34.2	9.1	-21.2	49.5	-21.2	21.2

Table 3.5 Occurrences of attribute and decision values among $|t| = 100$ samples for the respective tested strategies. The attribute entries $a_1 = 1$ and $a_2 = 1$ are the most desirable ones.

<i>Results for the naïve strategy</i>			
a_1	value 1: 55 times	value 2: 45 times	
a_2	value 1: 37 times	value 2: 36 times	value 3: 27 times
d	value 1: 100 times	value 2: 0 times	
<i>Results for the optimal strategy</i>			
a_1	value 1: 55 times	value 2: 45 times	
a_2	value 1: 58 times	value 2: 15 times	value 3: 27 times
d	value 1: 0 times	value 2: 100 times	
<i>Results for the proposed strategy</i>			
a_1	value 1: 57 times	value 2: 43 times	
a_2	value 1: 50 times	value 2: 22 times	value 3: 28 times
d	value 1: 40 times	value 2: 60 times	

Numerical outcomes of experiments are complemented by figures characterising the run of the proposed strategy. The left-hand side of Figure 3.2 shows the simulated attributes and the right-hand side shows the used decisions. The left-hand side of Figure 3.3 shows the probability (3.50) of the set $\alpha = \{(1, 1)\}$ of the most desirable attribute values used in the data censoring. The right-hand side displays tuning of the strategy, which gradually swaps from a bad strategy to the optimal one.

A representative fragment of the performed experiments is reported here. It suffices to add that longer runs (1000 samples and more) exhibited the expected convergence of the parameter estimates and of the proposed strategy.

3.6.3 Discussion of the Simulation Results

Estimation results reflected in Table 3.4 show that values of estimation errors are much less significant than their distribution over the parameter entries. They can fairly be evaluated only with respect to the achieved decision quality.

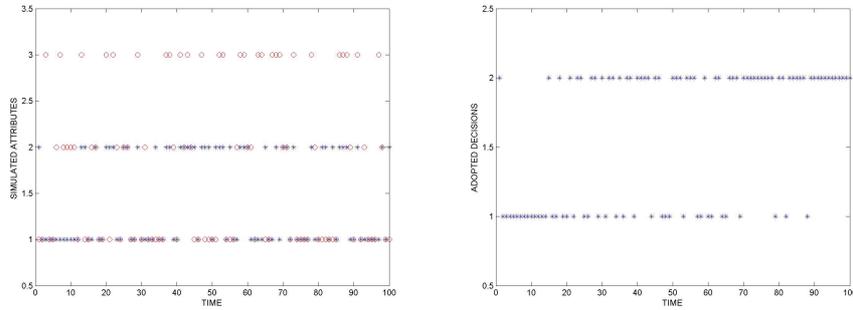


Fig. 3.2 The left figure shows simulated attributes (a_1 stars, a_2 circles). The right figure shows decisions generated by the proposed strategy.

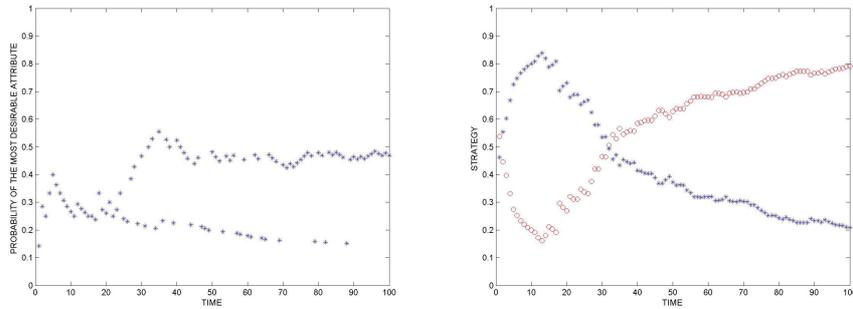


Fig. 3.3 The left figure shows probability $\int_{\alpha} M_{t+1}(a|d_t) da$ (3.50) used in the data censoring, Section 3.5.8. The right figure shows the proposed strategy, i.e. the pd given by the values $S_t(d = 1)$ (stars) and $S_t(d = 2) = 1 - S_t(d = 1)$ (circles).

Table 3.5 confirms that the naive strategy can be very bad indeed. It is also visible that the specific simulated environment, described in Table 3.3, makes only the values of the attribute a_2 sensitive to the decision taken: the simulated example is non-trivial in spite of its apparent simplicity. It is also worth noticing that the proposed strategy generates a relatively large portion of exploring decisions. Figure 3.2 shows that with increasing time the non-optimal decision ($d = 1$) is used less and less often than the optimal one ($d = 2$), as desirable. The similar desirable time dependence is visible in Figure 3.3: smaller probabilities (3.50) of the most desirable attribute pair, used for the data censoring (3.49), occur less and less often. The forgetting value stabilises still well below unity, which conforms with a high uncertainty connected with responses of the simulated environment. The same figure shows that after a relatively short learning period the proposed strategy starts to converge to the optimal one (as mentioned long runs confirmed).

3.7 Concluding Remarks

The chapter proposes a methodology of an automated elicitation of DM preference when the set of the most desirable attribute values is specified. This specification is quite common and can be extended to a preferential hierarchy with tunable weights of the DM preferences.

Many important features of the proposed solution are implied by the fact that the constructed ideal pd respects the relation between attributes and decisions as described by the environment model $M_t(a|d)$. Specifically,

- The proposed ideal pd is *not* fully concentrated on the most desirable attribute value, which reflects that it cannot be reached with certainty.
- The functional form of the ideal pd is determined by the model $M_t(a|d)$: it is not created in an ad-hoc, model independent, manner unlike utilities [33].
- It is always possible to project the constructed ideal pd into a class of feasible pds by using information criterion justified in [5], [28], Section 3.2, if the constructed ideal pd is too complex for numerical treatment or analysis.
- The environment model $M_t(a|d)$ as well as the closed-loop model with the optimal strategy $P_t(a, d) = M_t(a|d)S_t(d)$ are sequentially learnt; consequently, the DM preference description given by the ideal pd $I_t(a, d)$ derived from them, is learned, too.
- The approach can be directly extended to a dynamic DM with a regression-type dependence.
- The involved pds can quantify the joint distributions of discrete-valued and continuous-valued attributes. This simplifies the elicitation of the DM preferences given by categorical as well as numerical attributes.

Suppression or a more firm justification of ad-hoc steps (e.g. choice of forgetting factors in tracking of the changed DM preferences or in data censoring) is the key methodological research direction. At a more practical level, covering other DM preference specifications is the important problem to be addressed. The proposed solution is clearly connected with the DM preference learning presented in [27]. There, an explicit dependence between the environment-model parameter and the parameter of the optimal strategy has been established in a special, practically significant, way. A more detailed and general inspection of this relation is another open problem. The design of specific algorithmic solutions for commonly used environment models is a further topic to be covered.

In spite of the width of the problems hidden behind these open research directions, the proposed methodology appears practically promising.

Acknowledgements. This research and chapter have been strongly influenced by Dr. Tatiana V. Guy who I take as co-author of innovative and interesting parts of this text. I do acknowledge her help very much indeed.

The support of the project GAČR 13-13502S is acknowledged. The support of the project GAČR 102/08/0567 is also acknowledged.

References

1. Abramowitz, M., Stegun, I.: Handbook of Mathematical Functions. Dover Publications, New York (1972)
2. Barndorff-Nielsen, O.: Information and exponential families in statistical theory, New York (1978)
3. Berc, L., Kárný, M.: Identification of reality in Bayesian context. In: Warwick, K., Kárný, M. (eds.) Computer-Intensive Methods in Control and Signal Processing: Curse of Dimensionality, Birkhäuser, Boston, pp. 181–193 (1997)
4. Berger, J.: Statistical Decision Theory and Bayesian Analysis. Springer, New York (1985)
5. Bernardo, J.M.: Expected information as expected utility. *The Annals of Statistics* 7(3), 686–690 (1979)
6. Bohlin, T.: Interactive System Identification: Prospects and Pitfalls. Springer, New York (1991)
7. Boutilier, C.: A POMDP formulation of preference elicitation problems. In: Biundo, S. (ed.) AAAI-2002 Proc. of the Fifth European Conference on Planning, pp. 239–246. AAAI Press, Durham (2002)
8. Boutilier, C., Brafman, R., Geib, C., Poole, D.: A constraint-based approach to preference elicitation and decision making. In: Proceedings of AAAI Spring Symposium on Qualitative Decision Theory, Stanford, CA, pp. 19–28 (1997)
9. Bowong, S., Dimi, J., Kamgang, J., Mbang, J., Tewa, J.: Survey of recent results of multi-compartments intra-host models of malaria and HIV. *Revue ARIMA* 9, 85–107 (2008)
10. Chajewska, U., Koller, D.: Utilities as random variables: Density estimation and structure discovery. In: Proceedings of UAI 2000, pp. 63–71 (2000)
11. Cooke, N.: Varieties of knowledge elicitation techniques. *International Journal of Human-Computer Studies* 41, 801–849 (1994)
12. Feldbaum, A.: Theory of dual control. *Autom. Remote Control* 21(9) (1960)
13. Gajos, K., Weld, D.: Preference elicitation for interface optimization. In: Proceedings of UIST 2005 (2005)
14. Garthwaite, P., Kadane, J., O’Hagan, A.: Statistical methods for eliciting probability distributions. *J. of the American Statistical Association* 100(470), 680–700 (2005)
15. Guy, T.V., Kárný, M., Wolpert, D.H. (eds.): Decision Making with Imperfect Decision Makers. ISRL, vol. 28. Springer, Heidelberg (2012)
16. Viappiani, H.P., Boutilier, S.Z., Learning, C.: complex concepts using crowdsourcing: A Bayesian approach. In: Proceedings of the Second Conference on Algorithmic Decision Theory (ADT 2011), Piscataway, NJ (2011)
17. Haykin, S.: Neural Networks: A Comprehensive Foundation. Macmillan, New York (1994)
18. Horst, R., Tuy, H.: Global Optimization, p. 727. Springer (1996)
19. Jacobs, O., Patchell, J.: Caution and probing in stochastic control. *Int. J. of Control* 16, 189–199 (1972)
20. Jaynes, E.: Information theory and statistical mechanics. *Physical Review Series II* 106(4), 620–630 (1957)
21. Jimison, H., Fagan, L., Shachter, R., Shortliffe, E.: Patient-specific explanation in models of chronic disease. *AI in Medicine* 4, 191–205 (1992)
22. Jorgensen, S.B., Hangos, K.M.: Qualitative models as unifying modelling tool for grey box modelling. *Int. J. of Adaptive Control and Signal Processing* 9(6), 461–562 (1995)
23. Kárný, M.: Towards fully probabilistic control design. *Automatica* 32(12), 1719–1722 (1996)

24. Kárný, M., Andryšek, J., Bodini, A., Guy, T., Kracík, J., Nedoma, P., Ruggeri, F.: Fully probabilistic knowledge expression and incorporation. Tech. Rep. 8-10MI, CNR IMATI, Milan (2008)
25. Kárný, M., Andryšek, J., Bodini, A., Guy, T.V., Kracík, J., Ruggeri, F.: How to exploit external model of data for parameter estimation? *Int. J. of Adaptive Control and Signal Processing* 20(1), 41–50 (2006)
26. Kárný, M., Böhm, J., Guy, T.V., Jirsa, L., Nagy, I., Nedoma, P., Tesař, L.: *Optimized Bayesian Dynamic Advising: Theory and Algorithms*. Springer, London (2006)
27. Kárný, M., Guy, T.: Preference elicitation in fully probabilistic design of decision strategies. In: *Proc. of the 49th IEEE Conference on Decision and Control*. IEEE (2010)
28. Kárný, M., Guy, T.V.: On Support of Imperfect Bayesian Participants. In: Guy, T.V., Kárný, M., Wolpert, D.H. (eds.) *Decision Making with Imperfect Decision Makers*. ISRL, vol. 28, pp. 29–56. Springer, Heidelberg (2012)
29. Kárný, M., Guy, T.V.: Fully probabilistic control design. *Systems & Control Letters* 55(4), 259–265 (2006)
30. Kárný, M., Halousková, A., Böhm, J., Kulhavý, R., Nedoma, P.: Design of linear quadratic adaptive control: Theory and algorithms for practice. *Kybernetika* 21(Supplement to Nos. 3, 4, 5, 6) (1985)
31. Kárný, M., Kroupa, T.: Axiomatisation of fully probabilistic design. *Information Sciences* 186(1), 105–113 (2012)
32. Kárný, M., Kulhavý, R.: Structure determination of regression-type models for adaptive prediction and control. In: Spall, J. (ed.) *Bayesian Analysis of Time Series and Dynamic Models*. ch.12, Marcel Dekker, New York (1988)
33. Keeney, R., Raiffa, H.: *Decisions with Multiple Objectives: Preferences and Value Trade-offs*. John Wiley and Sons Inc. (1976)
34. Koopman, R.: On distributions admitting a sufficient statistic. *Tran. of American Mathematical Society* 39, 399 (1936)
35. Kuhn, H., Tucker, A.: Nonlinear programming. In: *Proc. of the 2nd Berkeley Symposium*, pp. 481–492. University of California Press, Berkeley (1951)
36. Kulhavý, R.: Can we preserve the structure of recursive Bayesian estimation in a limited-dimensional implementation? In: Helmke, U., Mennicken, R., Saurer, J. (eds.) *Systems and Networks: Mathematical Theory and Applications*, vol. I, pp. 251–272. Akademie Verlag, Berlin (1994)
37. Kulhavý, R., Zarrop, M.B.: On a general concept of forgetting. *Int. J. of Control* 58(4), 905–924 (1993)
38. Kullback, S., Leibler, R.: On information and sufficiency. *Annals of Mathematical Statistics* 22, 79–87 (1951)
39. Kumar, P.: A survey on some results in stochastic adaptive control. *SIAM J. Control and Applications* 23, 399–409 (1985)
40. Lainiotis, D.: Partitioned estimation algorithms, i: Nonlinear estimation. *Information Sciences* 7, 203–235 (1974)
41. Linden, G., Hanks, S., Lesh, N.: Interactive assessment of user preference models: The automated travel assistant. In: *Proceedings of User Modelling 1997* (1997)
42. Loeve, M.: *Probability Theory*. van Nostrand, Princeton (1962), Russian translation, Moscow
43. Nagy, I., Suzdaleva, E., Kárný, M., Mlynářová, T.: Bayesian estimation of dynamic finite mixtures. *Int. Journal of Adaptive Control and Signal Processing* 25(9), 765–787 (2011)
44. Nelsen, R.: *An Introduction to Copulas*. Springer, New York (1999)
45. O’Hagan, A., Buck, C.E., Daneshkhah, A., Eiser, J.R., Garthwaite, P.H., Jenkinson, D.J., Oakley, J., Rakow, T.: *Uncertain judgement: eliciting experts’ probabilities*. John Wiley & Sons (2006)

46. Osborne, M., Rubinstein, A.: A course in game theory. MIT Press (1994)
47. Peterka, V.: Bayesian system identification. In: Eykhoff, P. (ed.) Trends and Progress in System Identification, pp. 239–304. Pergamon Press, Oxford (1981)
48. Rao, M.: Measure Theory and Integration. John Wiley, New York (1987)
49. Savage, L.: Foundations of Statistics. Wiley, New York (1954)
50. Shi, R., MacGregor, J.: Modelling of dynamic systems using latent variable and subspace methods. *J. of Chemometrics* 14(5-6), 423–439 (2000)
51. Shore, J., Johnson, R.: Axiomatic derivation of the principle of maximum entropy and the principle of minimum cross-entropy. *IEEE Tran. on Information Theory* 26(1), 26–37 (1980)
52. Si, J., Barto, A., Powell, W., Wunsch, D. (eds.): Handbook of Learning and Approximate Dynamic Programming. Wiley-IEEE Press, Danvers (2004)
53. Šindelář, J., Vajda, I., Kárný, M.: Stochastic control optimal in the Kullback sense. *Kybernetika* 44(1), 53–60 (2008)

