

# Use of Prior Information in Structure Estimation

Miroslav Kárný, Petr Nedoma, Natalia Khaylova, Lenka  
Pavelková

Institute of Information Theory and Automation  
Pod vodárenskou věží 4, 182 08 Prague 8, POB 18, Czech Republic  
school@utia.cas.cz tel. +420-2-66052274,  
fax +420-2-66052068

**Key words:** adaptive control, Bayesian estimation, prior information, ARX model

## 1 Introduction

Adaptive LQG controllers that minimize approximately quadratic criterion while using recursively estimated linear Gaussian models become a standard in academic environment [1, 2, 3]. The version with controlled autoregressive models (ARX) belongs among the most successful ones as confirmed by their full scale applications [4, 5, 6]. At the same time, their potential is far from being adequately exploited. Man-power expensive commission is one of the main reasons of this undesirable state. This fact – that applies also to other adaptive controllers like GPC, MUSMAR etc. [3, 7] – stimulated an extensive project [8, 9, 10]. It aims at creating a complete computer support of the commission. At present, a full solution for single-input single-outputs systems is implemented in software system DESIGNER. It covers (i) data-preprocessing; (ii) quantification of prior information [?]; (iii) selection of the model structure [11]; (iv) off-line estimation [12] that serves for initialization of the on-line estimation as well as an alternative in the stabilized forgetting used for parameter tracking [13]; (v) estimation of the forgetting factor; (vi) tuning of kernels in the optimized quadratic loss so that user's aims and restrictions are met; this also provides off-line prediction of closed-loop behavior [14].

This paper proposes a correct combination of steps (ii) and (iii). The summary of steps (ii) - (iv) serves also to the companion paper [15], that provides a novel solution of the most difficult step (vi). Thus, in addition to particular improvements, the paper pair informs on the "technology used" within the DESIGNER.

Review of necessary results is given in Section 2 leading to the addressed problem formulation and its conceptual solution. Mapping of processed

knowledge on so called fictitious data is contained in Section 3. The results are elaborated to ARX model in Section 4, illustrated in Section 5 and complemented by remarks of Section 6.

The paper is relatively self-containing. The readers less familiar with the adopted Bayesian set up are referred to [12].

## 2 Problem formulation and solution

Here, the essence of the addressed problem and its conceptual solution are described after technical preliminaries. The used ARX model belongs to the exponential family. Its formal handling is very simple when its member-specific details are left aside. This motivates us to present all elements within it.

### 2.1 Basic notation and operations

The following notation is used throughout the text.

Symbol	Meaning
$\equiv$	equality by definition
$x^*$	a set of $x$ -values
$\dot{x}$	the number of elements in the vector or sequence $x$
$f(\cdot \cdot)$	conditional probability (density) functions (p(d)f)
$\mathcal{E}$ , var, cov	expectation, variance and covariance, respectively
$d(t)$	the sequence $(d_1, \dots, d_t)$
$\mathcal{S}$ , $\mathcal{K}$	model structure and knowledge item, respectively
$t$	discrete-time, always the last subscript after ;
$\tau_{\mathcal{K}}$	discrete time of fictitious data expressing knowledge $\mathcal{K}$
$i$	the subscript of the entry $d_{i;t}$ of a data item $d_t$
$I$ , tr, $'$	unit matrix, trace and transposition, respectively
$^{\lfloor\psi}L$	a non-numerical index $\psi$ of a variable $L$
$V, \nu$	statistics describing conjugate pdf to exponential family
$-$	bar distinguishing (flat) pre-prior pdf
$^{\lfloor\Delta}$	mark of increments of statistics obtained from data only

Note that pdfs are distinguished by the identifiers in their arguments. No formal distinction is made between random variable, its realization and argument of pdfs. The correct meaning follows from the context.

The following elementary operations of non-negative pdfs are used [12]

$$\begin{aligned}
\text{Normalization} \quad & \int f(a) da \equiv \int_{a \in a^*} f(a) da = 1, \\
\text{Chain rule} \quad & f(a, b|c) = f(a|b, c)f(b|c), \\
\text{Marginalization} \quad & f(a|c) = \int f(a, b|c) db, \\
\text{Bayes rule} \quad & f(a|b, c) = \frac{f(b|a, c)f(a|c)}{\int f(b|a, c)f(a|c) da} \propto f(b|a, c)f(a|c).
\end{aligned} \tag{1}$$

Bayesian paradigm we exploit operates on the joint pdf of all uncertain variables encountered. It composes this pdf from its elements and derives its particular marginal or conditional versions using (1). It inserts in them the measured realization of any variable that is at disposal.

Here, sequence of multi-variate data  $d(\hat{t}) = (d_1, \dots, d_{\hat{t}})$ , unknown, finite-dimensional parameters  $\Theta_{\mathcal{S}}$  and unknown structures  $\mathcal{S}$  of an adequate model are considered. The joint pdf is composed as follows

$$\begin{aligned}
& \underbrace{\text{joint pdf}}_{f(d(\hat{t}), \Theta_{\mathcal{S}}, \mathcal{S})} = \underbrace{\text{prior pdf} | \mathcal{S}}_{f(\Theta_{\mathcal{S}} | \mathcal{S})} \times \underbrace{\text{prior on } \mathcal{S}}_{f(\mathcal{S})} \times \underbrace{\prod_{t=1}^{\hat{t}} \prod_{i=1}^{\hat{d}} f(d_{i;t} | d_{i+1;t}, \dots, d_{\hat{d};t}, d(t-1), \Theta_{i\mathcal{S}})}_{\text{parameterized model}} \\
& \Theta_{\mathcal{S}} \equiv (\Theta_{1\mathcal{S}}, \dots, \Theta_{\hat{d}\mathcal{S}}).
\end{aligned} \tag{2}$$

## 2.2 Estimation and prediction in exponential family

The parameterized model in (2) of a fixed structure  $\mathcal{S}$  is the central modelling element. Within the control context, when the amount of observed data is permanently increasing, the following models are predominantly used.

**Agreement 2.1 (Exponential family)** *The  $i$ -th parameterized model belongs to the (dynamic) exponential family iff it can be written in the form*

$$\begin{aligned}
f(d_{i;t} | d_{i+1;t}, \dots, d_{\hat{d};t}, d(t-1), \Theta_{i\mathcal{S}}) &= f(d_{i;t} | \psi_{i\mathcal{S};t}, \Theta_{i\mathcal{S}}) \\
&= A(\Theta_{i\mathcal{S}}) \exp [\langle B(\Psi_{i\mathcal{S};t}), C(\Theta_{i\mathcal{S}}) \rangle],
\end{aligned} \tag{3}$$

where  $\Psi'_{i\mathcal{S};t} = [d_{i;t}, \psi'_{i\mathcal{S};t}]$  is data vector, given by a finite dimensional regression vector  $\psi_{i\mathcal{S};t}$ , depending on  $d_{i+1;t}, \dots, d_{\hat{d};t}$  and on  $d(t-1)$ ; it is assumed that the values of all data vectors  $\Psi_{i\mathcal{S};t-1}$ ,  $i = 1, \dots, \hat{d}$ , can be recursively updated using the newest data item  $d_t$  only,  $A(\cdot)$  is a non-negative scalar function defined on  $\Theta_{i\mathcal{S}}^*$ ,

$\langle \cdot, \cdot \rangle$  is a functional that is linear in the first argument,

$B(\cdot)$ ,  $C(\cdot)$  are either vector or matrix functions of compatible, finite dimensions. They are defined on  $\Psi_{iS;t}^*$  and  $\Theta_{iS}^*$ , respectively.

Practical significance of the exponential family becomes obvious when summarizing the corresponding estimation and prediction, [12].

**Proposition 2.1 (Estimation and prediction in exponential family)**

Let the parameterized model have the form (3) and the parameters  $\Theta_S = \Theta(\mathring{S})$  be a priori independent, i.e.  $f(\Theta_S) = \prod_{i=1}^{\mathring{d}} f(\Theta_{iS})$ . Let, moreover, the conjugate prior pdfs  $f(\Theta_{iS})$  [16]

$$f(\Theta_{iS}) \propto A^{\nu_{iS;0}}(\Theta_{iS}) \exp[\langle V_{iS;0}, C(\Theta_{iS}) \rangle] \chi_{\Theta_{iS}^*}(\Theta_{iS}) \equiv \mathcal{G}_{\Theta_{iS}}(V_{iS;0}, \nu_{iS;0}) \quad (4)$$

are used. The conjugate pdfs have the parameterized-model-induced functional form  $\mathcal{G}$ . They are determined by the prior finite-dimensional statistic  $V_{iS;0}$ , by the prior sample counter  $\nu_{iS;0}$  and indicator  $\chi_{\Theta_{iS}^*}(\Theta_{iS})$  of the set  $\Theta_{iS}^*$ .

Then, the parameters  $\Theta_{iS}$  are independent a posteriori and the respective posterior pdfs  $f(\Theta_{iS}|d(t))$  preserve the functional form of the prior pdfs

$$f(\Theta_{iS}|d(t)) = \frac{A^{\nu_{iS;t}}(\Theta_{iS}) \exp[\langle V_{iS;t}, C(\Theta_{iS}) \rangle] \chi_{\Theta_{iS}^*}(\Theta_{iS})}{\mathcal{I}(V_{iS;t}, \nu_{iS;t})} = \frac{\mathcal{G}_{\Theta_{iS}}(V_{iS;t}, \nu_{iS;t})}{\mathcal{I}(V_{iS;t}, \nu_{iS;t})}$$

$$\mathcal{I}(V_{iS;t}, \nu_{iS;t}) = \int A^{\nu_{iS;t}}(\Theta_{iS}) \exp[\langle V_{iS;t}, C(\Theta_{iS}) \rangle] \chi_{\Theta_{iS}^*}(\Theta_{iS}) d\Theta_{iS}. \quad (5)$$

The involved statistics  $V_{iS;t}$ ,  $\nu_{iS;t}$  can be updated recursively

$$V_{iS;t} = V_{iS;t-1} + B(\Psi_{iS;t}), \quad \nu_{iS;t} = \nu_{iS;t-1} + 1 \quad \text{with a priori chosen } V_{iS;0}, \nu_{iS;0}. \quad (6)$$

The predictive pdf, modelling evolution of the  $i$ -th data entry ( $i$ -th channel), is given by the formula

$$f(d_{i;t}|d_{i+1;t}, \dots, d_{\mathring{d};t}, d(t-1), \mathcal{S}) = \frac{\mathcal{I}(V_{iS;t-1} + B(\Psi_{iS;t}), \nu_{iS;t-1} + 1)}{\mathcal{I}(V_{iS;t-1}, \nu_{iS;t-1})}. \quad (7)$$

The overall predictive pdf, given by the structure  $\mathcal{S}$ , is product of pdfs (7) over  $i$ . The joint pdf of data conditioned by the structure  $\mathcal{S}$  is

$$f(d(\mathring{t})|\mathcal{S}) = \prod_{i=1}^{\mathring{d}} \mathcal{L}_i(d(\mathring{t}), \mathcal{S}) \quad \text{with } \mathcal{L}_i(d(\mathring{t}), \mathcal{S}) = \frac{\mathcal{I}(V_{iS;\mathring{t}}, \nu_{iS;\mathring{t}})}{\mathcal{I}(V_{iS;0}, \nu_{iS;0})} \quad (8)$$

called partial likelihood.  $\mathcal{L}_i(d(\mathring{t}), \mathcal{S})$  expresses descriptive abilities of the model having the structure  $\mathcal{S}$  judged with respect to  $i$ -th channel.

Thus, the estimation and prediction can be performed for respective *is*, channel-wise. We can focus the attention on a fixed channel and *drop the index  $i$*  remembering that a scalar variable is predicted.

The estimation and prediction reduce to algebraic operations with the *finite-dimensional statistic*  $V_{\mathcal{S};t}$  and of the *sample counter*  $\nu_{\mathcal{S};t}$ . Moreover, a single type of the integral  $\mathcal{I}(V_{\mathcal{S}}, \nu_{\mathcal{S}})$  has to be evaluated.

The need to get the *complete recursion* explains the requirement for the possibility to update data vector  $\Psi_{\mathcal{S};t}$  recursively. Note that this requirement excludes use of models with unknown moving average noise.

We inspect influence of particular knowledge items  $\mathcal{K} \in \mathcal{K}^* = \{1, \dots, \mathring{\mathcal{K}}\}$  on descriptive abilities of the adopted model. The notation  $\mathcal{L}(d(\mathring{t}), \mathcal{S}, \mathcal{K})$  stresses the use of the prior pdf  $f(\Theta_{\mathcal{S}}|\mathcal{K})$ . Similarly,  $\mathcal{L}(d(\mathring{t}), \mathcal{S}, \mathcal{K}(\mathring{\mathcal{K}}))$  denotes the joint predictive pdf obtained when using the prior pdf  $f(\Theta_{\mathcal{S}}|\mathcal{K}(\mathring{\mathcal{K}}))$  that includes all knowledge items available.

We need the following proposition (the fixed index  $\mathcal{S}$  is dropped).

**Proposition 2.2 (Weighted geometric mean of conjugate pdfs)** *Let  $f(\Theta)$ ,  ${}^{\mathsf{L}}f(\Theta)$  be a pair of pdfs conjugated to the parameterized model in the exponential family, i.e.  $f(\Theta) \propto \mathcal{G}_{\Theta}(V, \nu)$  and  ${}^{\mathsf{L}}f(\Theta) \propto \mathcal{G}_{\Theta}({}^{\mathsf{L}}V, {}^{\mathsf{L}}\nu)$ .*

*Then, their geometric mean  $f_{\lambda} \propto f^{\lambda} {}^{\mathsf{L}}f^{1-\lambda}$ , weighted by the factor  $\lambda \in [0, 1]$ , is the conjugated pdf  $f_{\lambda}(\Theta) \propto \mathcal{G}_{\Theta}(V_{\lambda}, \nu_{\lambda})$  whose statistics are*

$$V_{\lambda} = \lambda V + (1 - \lambda) {}^{\mathsf{L}}V, \quad \nu_{\lambda} = \lambda \nu + (1 - \lambda) {}^{\mathsf{L}}\nu. \quad (9)$$

This proposition serves for tracking of slow parameter changes using the *stabilized forgetting* [13]. There,  $f(\Theta) \equiv f(\Theta|d(t))$  is the posterior pdf of  $\Theta$  based on data  $d(t)$  measured up to the moment  $t$  when the forgetting is applied. The externally supplied *alternative pdf*  ${}^{\mathsf{L}}f(\Theta)$  describes possible changes of estimated parameters before measuring next data. The weight  $\lambda$  called *forgetting factor* is interpreted as the probability that the parameters do not change. The usual *exponential forgetting* is obtained by taking the completely flat alternative  ${}^{\mathsf{L}}f(\Theta) \propto 1$ . Use of the pre-prior pdf  $\bar{f}(\Theta) \propto \mathcal{G}_{\Theta}(\bar{V}, \bar{\nu})$  as the alternative pdf  ${}^{\mathsf{L}}f(\Theta)$  is more wise.

Mostly,  $\bar{f}(\Theta)$  is a flat pdf that respects just finiteness of  $\Theta$  values. It conservatively guarantees that the forgotten posterior pdf does not move on the area of improbable infinite values of  $\Theta$ .

The geometric mean of pdfs serves us also for finding a representant  $f(\Theta|d(\mathring{t}), \mathcal{K}(\mathring{\mathcal{K}}))$  of several pdfs  $f(\Theta|d(\mathring{t}), \mathcal{K})$ ,  $\mathcal{K} \in \mathcal{K}^*$ , each including a piece of prior knowledge  $\mathcal{K}$  about  $\Theta$ . For the measured data  $d(\mathring{t})$  and no prejudice, the degree of belief to each of them coincides with the posterior probability

$f(\mathcal{K}|d(\hat{t})) \propto \mathcal{L}(d(\hat{t}), \mathcal{K})$ . The representant  $f(\Theta|d(\hat{t}), \mathcal{K}(\hat{\mathcal{K}}))$ , called *merger* and motivated in [?], is chosen as the weighted geometric mean

$$\begin{aligned} f(\Theta|d(\hat{t}), \mathcal{K}(\hat{\mathcal{K}})) &\propto \prod_{\mathcal{K}=1}^{\hat{\mathcal{K}}} \left[ f(\Theta|d(\hat{t}), \mathcal{K}) \right]^{f(\mathcal{K}|d(\hat{t}))} \\ &= \mathcal{G}_{\Theta} \left( \sum_{\mathcal{K}=1}^{\hat{\mathcal{K}}} f(\mathcal{K}|d(\hat{t})) V_{\mathcal{K};\hat{t}}, \sum_{\mathcal{K}=1}^{\hat{\mathcal{K}}} f(\mathcal{K}|d(\hat{t})) \nu_{\mathcal{K};\hat{t}} \right). \end{aligned} \quad (10)$$

### 2.3 Structure estimation in nested exponential family

For a fixed functional form of the exponential family, the *model structure*  $\mathcal{S}$  is determined by the allowed entries of the regression vector. By collecting the potential entries into the *richest regression vector*  $\psi_{\mathcal{R};t}$ , the estimation of the model structure can be formulated as a selection of indices in it. They mark those entries that should be used in the proper regression vector  $\psi_{\mathcal{S};t}$ . There are  $2^{\psi_{\mathcal{R}}}$  of such options. This number is usually excessive one and prevents the straightforward Bayesian structure estimation through judging the posterior probabilities of competing structures  $f(\mathcal{S}|d(\hat{t})) \propto \mathcal{L}(d(\hat{t}), \mathcal{S})f(\mathcal{S})$ . These posterior probabilities qualify a posteriori the discrete pointers  $\mathcal{S}$  that have the prior pf  $f(\mathcal{S})$ .

The accumulation of  $V_{\mathcal{S};\hat{t}}$  makes the main computational burden related to the structure estimation. The current implementation in the system DESIGNER [17] relies on nesting of competitive structures  $\mathcal{S}$  in the richest one  $\mathcal{R}$ . The model, given by the data vectors  $\Psi_{\mathcal{S};t}$ , is called nested in the richest one  $\Psi_{\mathcal{R};t}$  if there is a *linear nesting mapping*  $N_{\mathcal{S}}$  such that, cf. (3),

$$B(\Psi_{\mathcal{S};t}) = N_{\mathcal{S}}[B(\Psi_{\mathcal{R};t})]. \quad (11)$$

This notion and Proposition 2.1 imply the following statement.

**Proposition 2.3 (Nesting in exponential family)** *Let the parameterized model with the richest structure belong to the exponential family (3)  $f(d_t|\psi_{\mathcal{R};t}, \Theta_{\mathcal{R}}) = A(\Theta_{\mathcal{R}}) \exp[\langle B(\Psi_{\mathcal{R};t}), C(\Theta_{\mathcal{R}}) \rangle]$ . Let us consider another model  $f(d_t|\psi_{\mathcal{S};t}, \Theta_{\mathcal{S}}) = A(\Theta_{\mathcal{S}}) \exp[\langle B(\Psi_{\mathcal{S};t}), C(\Theta_{\mathcal{S}}) \rangle]$ . Let  $N_{\mathcal{S}}$  be a time-invariant, linear nesting mapping such that  $B(\Psi_{\mathcal{S};t}) = N_{\mathcal{S}}[B(\Psi_{\mathcal{R};t})]$ . Let us assume that  $V_{\mathcal{S};0} = N_{\mathcal{S}}[V_{\mathcal{R};0}]$ .*

*Then, the  $V$ -statistics of the posterior pdfs of both models are related by the nesting mapping  $V_{\mathcal{S};\hat{t}} = N_{\mathcal{S}}[V_{\mathcal{R};\hat{t}}]$  and the posterior probability on the*

structure  $\mathcal{S}$  is given by the formula

$$f(\mathcal{S}|d(\mathring{t})) \propto \frac{\mathcal{I}(N_{\mathcal{S}}[V_{\mathcal{R};\mathring{t}}], \nu_{\mathcal{S};\mathring{t}})}{\mathcal{I}(N_{\mathcal{S}}[V_{\mathcal{R};0}], \nu_{\mathcal{S};0})} f(\mathcal{S}). \quad (12)$$

Thus, for the nested models and *nested prior statistics*, it is sufficient to collect the  $V$ -statistic for the richest structure. It helps but only partially. Full evaluation of the pf values (12) on the complete set of  $2^{\mathring{\psi}_{\mathcal{R}}}$  competitive structures is still impossible. Thus, *maximum a posteriori probability* (MAP) estimate of  $\mathcal{S}$  has to be searched for. The non-normalized values of the pf (12) evaluated during the search provide useful partial information on highly probable structures. The following conceptual search algorithm is used [11].

**Algorithm 2.1 (SEN: MAP structure estimate of nested model)**

Initial phase

- Collect the real-data-dependent increment  ${}^{\mathbb{L}}\Delta V_{\mathcal{R};\mathring{t}}$  of the  $V$ -statistic corresponding to the richest structure of the data vector  $\Psi_{\mathcal{R};t}$

$${}^{\mathbb{L}}\Delta V_{\mathcal{R};\mathring{t}} = \sum_{t=1}^{\mathring{t}} B(\Psi_{\mathcal{R};t}). \quad (13)$$

- Select the prior statistic  $V_{\mathcal{R};0}$  so that  $V_{\mathcal{S};0}$  are nested in it  $\forall \mathcal{S} \in \mathcal{S}^*$ .
- Specify prior pf  $f(\mathcal{S})$  of competitive structures, often as uniform one.

Search phase is run until the pre-specified number of restarts is reached.

1. Initialize the current guess of the structure.

*Empty, richest and user-specified structures of regression vectors are used. These options are complemented by structures selected randomly among a priori possible ones.*

2. Do while the value of the posterior partial likelihood increases.

- (a) Make full search for the best structure within a neighborhood of the current guess of the structure, i.e. maximize within it the posterior partial likelihood

$$\mathcal{L}(d(\mathring{t}), \mathcal{S}) f(\mathcal{S}) = \frac{\mathcal{I}\left(N_{\mathcal{S}}\left[{}^{\mathbb{L}}\Delta V_{\mathcal{R};\mathring{t}} + V_{\mathcal{R};0}\right], \mathring{t} + \nu_{\mathcal{R};0}\right)}{\mathcal{I}(N_{\mathcal{S}}[V_{\mathcal{R};0}], \nu_{\mathcal{R};0})} f(\mathcal{S}).$$

*The neighborhood consists of all structures gained by:*

- *adding a single entry to the current guess of the structure,*
  - *removing a single entry from the current guess of the structure,*
  - *considering structures nested in those named above.*
- (b) *Take the maximizer as a new current guess of the structure.*

## 2.4 Quantification of prior knowledge

Parameter and structure estimation are sensitive to the information content of the learning data. As a rule, the available data are poorly informative. Then, prior knowledge has to be used. In the Bayesian set up, it is feeded through the prior pdf. Here, we outline the way how it can be constructed. The following conditions are specific for technological applications.

- ⊕ Groups of widely accessible knowledge types exist.
- ⊕ Experimental data  $d(\mathring{t})$  measured on the modelled system are available.
- ⊕ Admissible prior pdfs are restricted to conjugate ones, thus the prior knowledge is to be translated into values of the prior statistics  $V, \nu$ .
- ⊖ The person feeding the prior knowledge does not care about the probabilistic tool set exploited.
- ⊖ No supervisor for knowledge elicitation and judgement of the expert competence is at disposal.
- ⊖ Knowledge items processed are expected to be repetitive, not fully consistent and differing in precision and nature. Mutual dependencies of knowledge items are unknown.

The following quantification algorithm respects these conditions [?].

### Algorithm 2.2 (Quantification of prior knowledge)

Initiation phase

- *Select the  $i$ -th parameterized model (2) of a fixed structure  $\mathcal{S}$  in the exponential family you deal with.*
- *Collect the real-data-dependent increment  ${}^{\text{L}}\Delta V_{\mathcal{S};i}$  of the  $V$ -statistics according to (13) for  $\mathcal{R} = \mathcal{S}$ .*
- *Split the existing knowledge into internally consistent knowledge items.*
- *Select the pre-prior pdf  $\bar{f}(\Theta_{\mathcal{S}}) \propto \mathcal{G}_{\Theta_{\mathcal{S}}}(\bar{V}_{\mathcal{S}}, \bar{\nu}_{\mathcal{S}})$  on unknown parameters  $\Theta_{\mathcal{S}}$  that expresses the common knowledge available.*
- *Initialize the auxiliary scalar normalization factor  $s = 0$ .*

Quantification phase *runs for the internally consistent knowledge items  $\mathcal{K} \in \mathcal{K}^*$ .*



- Translate the knowledge item  $\mathcal{K}$  into the fictitious-data  $d(\hat{t}_{\mathcal{K}})$  dependent increments  ${}^{\mathcal{L}}\Delta V_{\mathcal{S};\hat{t}_{\mathcal{K}}}$ ,  ${}^{\mathcal{L}}\Delta \nu_{\mathcal{S};\hat{t}_{\mathcal{K}}}$  of the pre-prior statistics  $\bar{V}_{\mathcal{S}}$ ,  $\bar{\nu}_{\mathcal{S}}$  so that  $V_{\mathcal{SK};0} = {}^{\mathcal{L}}\Delta V_{\mathcal{S};\hat{t}_{\mathcal{K}}} + \bar{V}_{\mathcal{S}}$  and  $\nu_{\mathcal{SK};0} = {}^{\mathcal{L}}\Delta \nu_{\mathcal{S};\hat{t}_{\mathcal{K}}} + \bar{\nu}_{\mathcal{S}}$  reflect the processed knowledge item, i.e.  $f(\Theta_{\mathcal{S}}|\mathcal{K}) \propto \mathcal{G}_{\Theta_{\mathcal{S}}}(V_{\mathcal{SK};0}, \nu_{\mathcal{SK};0})$ .
- Evaluate the descriptive abilities gained by exploiting this knowledge on real data  $d(\hat{t})$  and update the normalization factor  $s$

$$\mathcal{L}(d(\hat{t}), \mathcal{S}, \mathcal{K}) = \frac{\mathcal{I}({}^{\mathcal{L}}\Delta V_{\mathcal{SK};\hat{t}} + V_{\mathcal{SK};0}, \hat{t} + \nu_{\mathcal{SK};0})}{\mathcal{I}(V_{\mathcal{SK};0}, \nu_{\mathcal{SK};0})}, \quad s = s + \mathcal{L}(d(\hat{t}), \mathcal{S}, \mathcal{K}). \quad (14)$$

Merging phase combines particular knowledge items into the merger (10)

$$f(\Theta_{\mathcal{S}}|d(\hat{t}), \mathcal{K}(\hat{\mathcal{K}})) \propto f(d(\hat{t})|\Theta_{\mathcal{S}}) \prod_{\mathcal{K}=1}^{\hat{\mathcal{K}}} [f(\Theta_{\mathcal{S}}|\mathcal{K})]^{f(\mathcal{K}|d(\hat{t}), \mathcal{S})} \quad (15)$$

$$f(\mathcal{K}|d(\hat{t}), \mathcal{S}) = \frac{\mathcal{L}(d(\hat{t}), \mathcal{S}, \mathcal{K})}{s}.$$

It gives  $f(\Theta_{\mathcal{S}}|d(\hat{t}), \mathcal{K}(\hat{\mathcal{K}})) \propto$

$$\mathcal{G}_{\Theta_{\mathcal{S}}} \left( \underbrace{{}^{\mathcal{L}}\Delta V_{\mathcal{S};\hat{t}} + \underbrace{\sum_{\mathcal{K}=1}^{\hat{\mathcal{K}}} f(\mathcal{K}|d(\hat{t}), \mathcal{S}) V_{\mathcal{SK};0}}_{V_{\mathcal{SK}(\hat{\mathcal{K}});0}}}_{V_{\mathcal{SK}(\hat{\mathcal{K}});\hat{t}}} + \underbrace{\sum_{\mathcal{K}=1}^{\hat{\mathcal{K}}} f(\mathcal{K}|d(\hat{t}), \mathcal{S}) \nu_{\mathcal{SK};0}}_{\nu_{\mathcal{SK}(\hat{\mathcal{K}});0}} \right).$$

It remains to specify the meaning of internally consistent knowledge item and to show how to construct the increments of statistics on fictitious data. These aspects are covered in Section 3.

## 2.5 Addressed problem and its conceptual solution

The problem, we are facing, stems from the fact that usually the *prior sufficient statistic*  $V_{\mathcal{SK};0}$  expressing the piece of knowledge  $\mathcal{K}$  within the structure  $\mathcal{S}$  is *not nested* in the statistics corresponding to that with the richest data vector. In other words, the *efficient nested structure estimation* Algorithm 2.1 *cannot be directly used whenever a non-trivial prior knowledge is to be exploited*. This fact was overlooked in our former implementations and publications and caused worse estimation results than expected.

Knowing the problem, the remedy is rather straightforward. The following conceptual algorithm is used.

### Algorithm 2.3 (Structure estimation with prior knowledge)

Initial phase

- Select the parameterized model in the exponential family and the richest structure of the underlying regression vector  $\psi_{\mathcal{R}}$ .
- Select the pre-prior pdf  $\bar{f}(\Theta_{\mathcal{R}}) \propto \mathcal{G}_{\Theta_{\mathcal{R}}}(\bar{V}_{\mathcal{R}}, \bar{\nu}_{\mathcal{R}})$  on unknown parameters  $\Theta_{\mathcal{R}}$  that express the common knowledge available while requiring that the same knowledge is expressed for all nested structures of interest.

Structure estimation with nested prior knowledge

- Apply the algorithm SEN 2.1 in order to find a pre-selected number  $\hat{\mathcal{S}}$  of structures  $\mathcal{S} \in \mathcal{S}^* = \{1, 2, \dots, \hat{\mathcal{S}}\}$  having the highest posterior probabilities when using the restricted prior knowledge described by the pdf  $\bar{f}(\Theta_{\mathcal{R}})$ .

Inclusion of prior knowledge for promising structures  $\mathcal{S} \in \mathcal{S}^*$

1. Apply Algorithm 2.2 for the fixed structure  $\mathcal{S}$  to get the statistics of the best merger  $V_{\mathcal{SK}(\hat{\mathcal{K}}); \tau}, \nu_{\mathcal{SK}(\hat{\mathcal{K}}); \tau}, \tau \in \{0, \hat{t}\}$ , cf. (15).
2. Evaluate descriptive abilities of the best merger, within the given structure  $\mathcal{S}$ , cf. (15),

$$\mathcal{L}(d(\hat{t}), \mathcal{S}) = \frac{\mathcal{I}(V_{\mathcal{SK}(\hat{\mathcal{K}}); \hat{t}}, \nu_{\mathcal{SK}(\hat{\mathcal{K}}); \hat{t}})}{\mathcal{I}(V_{\mathcal{SK}(\hat{\mathcal{K}}); 0}, \nu_{\mathcal{SK}(\hat{\mathcal{K}}); 0})}.$$

3. Provide  $\hat{\mathcal{S}} \in \text{Arg max}_{\mathcal{S}^*} \mathcal{L}(d(\hat{t}), \mathcal{S}) f(\mathcal{S})$  as the structure estimate and its statistics  $V_{\hat{\mathcal{S}}\mathcal{K}(\hat{\mathcal{K}}); \hat{t}}, \nu_{\hat{\mathcal{S}}\mathcal{K}(\hat{\mathcal{K}}); \hat{t}}$  as initial conditions for the subsequent on-line estimation and as the alternative pdf needed in the stabilized forgetting.

## 3 Fictitious data

Here, the construction of the common information basis, i.e. fictitious data, is recalled and refined [?]. It allows us to cope with knowledge items of a different nature in a unified way.

### 3.1 Internally consistent fictitious data blocks

Some information sources provide knowledge pieces  $\mathcal{K}$  in form of data blocks  $d(\hat{\tau}_{\mathcal{K}})$ . They include obsolete or interpolated data measured on the system in question or data measured on a similar system, data from identification experiments violating usual working conditions – like measurement of step response – and data gained from a realistic simulation.

The data block  $d(\hat{\tau}_{\mathcal{K}})$  expressing the knowledge piece  $\mathcal{K}$  is called internally consistent iff  $f(\Theta_S|\mathcal{K})$  coincides with a flattened version of the posterior pdf  $f(\Theta_S|d(\hat{\tau}_{\mathcal{K}})) \propto \mathcal{G}_{\Theta_S}(V_{S;\hat{\tau}_{\mathcal{K}}}, \nu_{S;\hat{\tau}_{\mathcal{K}}})$ . Let us describe it in a detail.

The posterior pdf is obtained by the Bayes rule applied to the pre-prior pdf  $\bar{f}(\Theta_S) \propto \mathcal{G}_{\Theta_S}(\bar{V}_S, \bar{\nu}_S)$  with the stabilized forgetting. The *forgetting* is used in order to counteract mismodelling. The pre-prior pdf is used as the alternative pdf. Thus, the adequate forgetting factor  $\lambda$  is to be chosen only. A comparison of descriptive abilities of the posterior pdfs obtained for various forgetting factors serves well to this purpose. It is done anyway during merging of individual knowledge pieces. Thus, it suffices to take pdfs  $f(\Theta_S|d(\hat{\tau}_{\mathcal{K}}), \lambda) \propto \mathcal{G}_{\Theta_S}(V_{S\lambda;\hat{\tau}_{\mathcal{K}}}, \nu_{S\lambda;\hat{\tau}_{\mathcal{K}}})$  gained for different  $\lambda$  as different knowledge pieces. It is done from now on and the reference to  $\lambda$  is suppressed.

Mismodelling or obsolete nature of the fictitious data blocks imply that the pdfs

$$f(\Theta_S|d(\hat{\tau}_{\mathcal{K}})) = \mathcal{G}_{\Theta_S}(V_{S;\hat{\tau}_{\mathcal{K}}}, \nu_{S;\hat{\tau}_{\mathcal{K}}}) \equiv \mathcal{G}_{\Theta_S}({}^{\Delta}V_{S;\hat{\tau}_{\mathcal{K}}} + \bar{V}_S, {}^{\Delta}\nu_{S;\hat{\tau}_{\mathcal{K}}} + \bar{\nu}_S) \quad (16)$$

reflect system properties approximately only. Consequently, these pdfs have to be adequately flattened before merging such a piece of knowledge. This is of an extreme importance as, for instance, simulators may provide a huge amount of data that may over-fit the posterior pdfs at wrong positions. The geometric mean serves well for flattening of these pdfs, cf. (16),

$$\begin{aligned} f(\Theta_S|\mathcal{K}) &\propto [f(\Theta_S|d(\hat{\tau}_{\mathcal{K}}))]_{S\mathcal{K}}^{\Lambda} [\bar{f}(\Theta_S)]^{1-\Lambda_{S\mathcal{K}}} \\ &\propto \mathcal{G}_{\Theta_S}(\Lambda_{S\mathcal{K}}V_{S;\hat{\tau}_{\mathcal{K}}} + (1-\Lambda_{S\mathcal{K}})\bar{V}_S, \Lambda_{S\mathcal{K}}\nu_{S;\hat{\tau}_{\mathcal{K}}} + (1-\Lambda_{S\mathcal{K}})\bar{\nu}_S) \\ &= \mathcal{G}_{\Theta_S}(\Lambda_{S\mathcal{K}}{}^{\Delta}V_{S;\hat{\tau}_{\mathcal{K}}} + \bar{V}_S, \Lambda_{S\mathcal{K}}{}^{\Delta}\nu_{S;\hat{\tau}_{\mathcal{K}}} + \bar{\nu}_S). \end{aligned} \quad (17)$$

The scalar  $\Lambda_{S\mathcal{K}} \in [0, 1]$ , called *flattening factor*, is selected to maximize the descriptive ability, i.e.

$$\Lambda_{S\mathcal{K}} = \arg \max_{\Lambda \in [0,1]} \frac{\mathcal{I}({}^{\Delta}V_{S;\hat{\tau}_{\mathcal{K}}} + \Lambda({}^{\Delta}V_{S;\hat{\tau}_{\mathcal{K}}} + \bar{V}_S; \hat{\tau}_{\mathcal{K}} + \Lambda({}^{\Delta}\nu_{S;\hat{\tau}_{\mathcal{K}}} + \bar{\nu}_S))}{\mathcal{I}(\Lambda({}^{\Delta}V_{S;\hat{\tau}_{\mathcal{K}}} + \bar{V}_S; \hat{\tau}_{\mathcal{K}}) + \Lambda({}^{\Delta}\nu_{S;\hat{\tau}_{\mathcal{K}}} + \bar{\nu}_S))}. \quad (18)$$

The values of the flattening factor  $\Lambda_{\mathcal{S}\mathcal{K}}$  are structure-dependent. Thus, the optimized flattening spoils possible nesting of  ${}^{\Delta}V_{\mathcal{S};\tau_{\mathcal{K}}}$ .

### 3.2 Construction of fictitious data

Here, we focus on those prior knowledge items that *do not have directly the form of data blocks* but can be interpreted as the *expected system response on a pre-specified stimulus*. Static gain, a point on frequency response, time-constants serve as their prominent examples. Such a knowledge item  $\mathcal{K}$  is *always uncertain* to some degree. It can be interpreted as a collection of partial characterizations of the several predictors. Each of them is expressed in terms of its prior pdf  $f(\Theta_{\mathcal{S}}|\tau_{\mathcal{K}})$

$$f(d_{\tau_{\mathcal{K}}}|\psi_{\mathcal{S};\tau_{\mathcal{K}}}) = \int f(d_{\tau_{\mathcal{K}}}|\psi_{\mathcal{S};\tau_{\mathcal{K}}}, \Theta_{\mathcal{S}}) f(\Theta_{\mathcal{S}}|\tau_{\mathcal{K}}) d\Theta_{\mathcal{S}} \quad (19)$$

for respective regression vectors  $\psi_{\mathcal{S};\tau_{\mathcal{K}}}$ ,  $\tau_{\mathcal{K}} = 1, \dots, \tau_{\mathcal{K}}^{\circ}$ . Mostly, the  $\tau_{\mathcal{K}}$ -th part of the knowledge piece  $\mathcal{K}$  can be expressed as initial moments of the pdf (19). Formally,

$$h(\psi_{\mathcal{S};\tau_{\mathcal{K}}}) = \int H(d_{\tau_{\mathcal{K}}}, \psi_{\mathcal{S};\tau_{\mathcal{K}}}) f(d_{\tau_{\mathcal{K}}}|\psi_{\mathcal{S};\tau_{\mathcal{K}}}) dd_{\tau_{\mathcal{K}}}. \quad (20)$$

$h(\psi_{\mathcal{S};\tau_{\mathcal{K}}})$  and  $H(\Psi_{\mathcal{S};\tau_{\mathcal{K}}}) = H(d_{i;\tau_{\mathcal{K}}}, \psi_{\mathcal{S};\tau_{\mathcal{K}}})$  are known vector functions of the indicated arguments. When there is no pdf  $f(\Theta_{\mathcal{S}}|\tau_{\mathcal{K}})$  fulfilling (19), (20), then this information source is inconsistent and has to be split into several, internally consistent, knowledge sources. Then, the restrictions (19), (20) do not determine fully the constructed pdf  $f(\Theta_{\mathcal{S}}|\tau_{\mathcal{K}})$ . Pragmatic reasons make us to search within the class of conjugate pdfs. Moreover, it is reasonable to construct such a pdf  $f(\Theta_{\mathcal{S}}|\tau_{\mathcal{K}})$  that reflects *just* the knowledge item expressed by (19), (20). Thus, it makes sense to choose such a pdf  $f(\Theta_{\mathcal{S}}|\tau_{\mathcal{K}})$  that is the nearest one to the flat pre-prior pdf  $\bar{f}(\Theta_{\mathcal{S}})$ . The choice is made among those meeting (19), (20). The Kullback-Leibler distance [18]  $\mathcal{D}(f||\bar{f}) = \int f(\Theta_{\mathcal{S}}|\tau_{\mathcal{K}}) \ln [f(\Theta_{\mathcal{S}}|\tau_{\mathcal{K}})/\bar{f}(\Theta_{\mathcal{S}})] d\Theta_{\mathcal{S}}$  is used as an adequate proximity measure.

Both the found  $f(\Theta_{\mathcal{S}}|\tau_{\mathcal{K}})$  and pre-prior pdf  $\bar{f}(\Theta_{\mathcal{S}})$  belong to the same conjugate class. Thus, their ratio can be interpreted as a product of the parameterized model at some fictitious data vectors, leading to the statistics increments  ${}^{\Delta}V_{\mathcal{S};\tau_{\mathcal{K}}}$ ,  ${}^{\Delta}\nu_{\mathcal{S};\tau_{\mathcal{K}}}$ . The knowledge item  $\mathcal{K}$  is supposed to be internally consistent. Consequently, fictitious data vectors obtained for various  $\tau_{\mathcal{K}}$  should be processed by using the Bayes rule. The result is then scaled

in the same way as the consistent data block, Section 3.1. Thus,  $f(\Theta_S|\mathcal{K})$  has the form (17) with statistics

$$\mathbb{L}^{\Delta}V_{\mathcal{S};\hat{\tau}_{\mathcal{K}}} = \sum_{\tau_{\mathcal{K}}=1}^{\hat{\tau}_{\mathcal{K}}} \mathbb{L}^{\Delta}V_{\mathcal{S};\tau_{\mathcal{K}}}, \quad \mathbb{L}^{\Delta}\nu_{\mathcal{S};\hat{\tau}_{\mathcal{K}}} = \sum_{\tau_{\mathcal{K}}=1}^{\hat{\tau}_{\mathcal{K}}} \mathbb{L}^{\Delta}\nu_{\mathcal{S};\tau_{\mathcal{K}}}, \quad \text{where} \quad (21)$$

$\mathcal{G}_{\Theta_S} = \left( \mathbb{L}^{\Delta}V_{\mathcal{S};\tau_{\mathcal{K}}} + \bar{V}_S, \mathbb{L}^{\Delta}\nu_{\mathcal{S};\tau_{\mathcal{K}}} + \bar{\nu}_S \right)$  are pdfs minimizing the Kullback-Leibler distance to  $\mathcal{G}_{\Theta_S}(\bar{V}_S, \bar{\nu}_S)$  under restrictions (19), (20). The factor  $\Lambda_{\mathcal{SK}}$  results from the maximization (18).

Specific optimization steps are elaborated for the ARX model in the next section.

## 4 Application to normal ARX model

The system DESIGNER, in which the presented results are predominantly used, relies on the normal ARX model and its variants. This makes us to specialize the general solution to this case.

### 4.1 Estimation and prediction with normal ARX model

The normal ARX model belongs to the exponential family with the following correspondence to its general form (3)

$$f(d|\psi, \Theta) = \mathcal{N}_d(\theta'\psi, r) = A(\Theta) \exp[\langle B(\Psi), C(\Theta) \rangle] \quad \text{with} \quad (22)$$

$$\Theta = [\theta, r] = [\text{regression coefficients, noise variance}], \quad A(\Theta) = (2\pi r)^{-0.5}$$

$$\langle B, C \rangle = \text{tr}[B'C], \quad B(\Psi) = \Psi\Psi', \quad C(\Theta) = (2r)^{-1}[-1, \theta']'[-1, \theta'].$$

This correspondence determines the conjugate prior (4) in the form known as the Gauss-inverse-Wishart (*GiW*) pdf [19]

$$\mathcal{G}_{\Theta}(V, \nu) = r^{-0.5(\nu + \hat{\psi} + 2)} \exp \left\{ -\frac{1}{2r} \text{tr} \left( V [-1, \theta']' [-1, \theta'] \right) \right\}. \quad (23)$$

The  $(\hat{\Psi}, \hat{\Psi})$ -dimensional *extended information matrix*  $V$  can be chosen symmetric and must be *positive definite*. Otherwise, the function  $\mathcal{G}_{\Theta}(V, \nu)$  cannot be normalized to a pdf. Thus, we can use the numerically advantageous  $L'DL$  decomposition of this matrix [20]

$$V = L'DL, \quad L = \text{lower triangular with unit diagonal}, \quad D = \text{diagonal}$$

$$L = \begin{bmatrix} 1 & 0 \\ \mathbb{L}^d\psi_L & \mathbb{L}^{\psi}L \end{bmatrix}, \quad D = \text{diag} \left[ \mathbb{L}^dD, \mathbb{L}^{\psi}D \right], \quad \mathbb{L}^dD = \text{scalar}. \quad (24)$$

**Proposition 4.1 (Basic properties and moments of the GiW pdf)**

The conjugate GiW pdf has the following alternative expression

$$\begin{aligned}\mathcal{G}_\Theta(L, D, \nu) &= \frac{r^{-0.5(\nu+\mathring{\psi}+2)}}{\mathcal{I}(L, D, \nu)} \exp \left\{ -\frac{1}{2r} \left[ (\theta - \hat{\theta})' C^{-1} (\theta - \hat{\theta}) + \mathbb{I}^d D \right] \right\} \\ \hat{\theta} &= \mathbb{I}^\psi L^{-1} \mathbb{I}^{d\psi} L = \text{least-squares (LS) estimate of } \theta \\ C &= \mathbb{I}^\psi L^{-1} \mathbb{I}^\psi D^{-1} \left( \mathbb{I}^\psi L' \right)^{-1} = \text{LS covariance factor} \\ \mathbb{I}^d D &= \text{LS remainder.}\end{aligned}\tag{25}$$

The normalization integral is

$$\begin{aligned}\mathcal{I}(L, D, \nu) &= \Gamma(0.5\nu) \mathbb{I}^d D^{-0.5\nu} \prod_{j=1}^{\mathring{\psi}} \mathbb{I}^\psi D_{jj}^{-0.5} 2^{0.5\nu} (2\pi)^{0.5\mathring{\psi}} \quad \text{with} \\ \Gamma(x) &= \int_0^\infty z^{x-1} \exp(-z) dz < \infty \quad \text{for } x > 0.\end{aligned}\tag{26}$$

Thus, it is finite iff  $\nu > 0$  and  $V$  is positive definite  $\Leftrightarrow D_{jj} > 0$ ,  $j = 1, \dots, \mathring{\Psi}$ . Under this condition, the normalization integral  $\mathcal{I}$  is positive.

The GiW pdf has the following moments

$$\begin{aligned}\mathcal{E}[r|L, D, \nu] &= \frac{\mathbb{I}^d D}{\nu - 2} = \hat{r}, \quad \text{var}[r|L, D, \nu] = \frac{2\hat{r}^2}{\nu - 4} \\ \mathcal{E}[\theta|L, D, \nu] &= \mathbb{I}^\psi L^{-1} \mathbb{I}^{d\psi} L = \hat{\theta} \\ \text{cov}[\theta|L, D, \nu] &= \frac{\mathbb{I}^d D}{\nu - 2} \mathbb{I}^\psi L^{-1} \mathbb{I}^\psi D^{-1} \left( \mathbb{I}^\psi L' \right)^{-1} = \hat{r} C.\end{aligned}\tag{27}$$

Proposition 2.1 specializes to the following normal variant.

**Proposition 4.2 (Estimation & prediction with ARX model)** Let the normal ARX model (22) be considered, together with the conjugate GiW prior pdf  $\mathcal{G}_\Theta(L_0, D_0, \nu_0)$  (25) and the alternative pdf  $\mathcal{G}_\Theta(\mathbb{I}^a L, \mathbb{I}^a D, \mathbb{I}^a \nu)$  in the stabilized forgetting with the forgetting factor  $\lambda \in [0, 1]$ .

Then, the posterior pdf is the GiW pdf  $\mathcal{G}_\Theta(L_t, D_t, \nu_t)$  and its sufficient statistics evolve according to the recursions

$$\begin{aligned}L_t' D_t L_t &= \underbrace{\begin{bmatrix} L_{t-1} \\ \Psi_t' \\ \mathbb{I}^a L \end{bmatrix}'}_{\mathbf{L}'} \underbrace{\text{diag} \begin{bmatrix} \lambda D_{t-1} \\ \lambda \\ (1 - \lambda) \mathbb{I}^a D \end{bmatrix}}_{\mathbf{D}} \underbrace{\begin{bmatrix} L_{t-1} \\ \Psi_t' \\ \mathbb{I}^a L \end{bmatrix}}_{\mathbf{L}} \\ \nu_t &= \lambda(\nu_{t-1} + 1) + (1 - \lambda) \mathbb{I}^a \nu, \quad L_0, D_0, \nu_0 \text{ given a priori.}\end{aligned}\tag{28}$$

The rectangular matrix  $\mathbf{L}'$  is mapped on  $[L'_t, 0]$  by the regular matrix  $\mathcal{T}$  that diagonalizes the  $(\ddot{\Psi}, \ddot{\Psi})$ -left upper corner in  $\mathcal{T}^{-1} \text{diag}[D, 1](\mathcal{T}')^{-1}$ .

The predictive pdf is known as the Student pdf. For any data vector  $\Psi = [d, \psi']'$ , its values can be found numerically as the ratio (7).

## 4.2 Internally consistent fictitious data blocks

Processing of internally consistent data blocks coincides with Bayesian estimation of the ARX normal model. The one-to-one mapping of the extended information matrix  $V$  on the LS quantities, Proposition 4.1, implies that its updating is equivalent to recursive least squares [12]. It is equipped with tracking ability through the stabilized forgetting. Numerically, its  $L'DL$  decomposition is evaluated by using an efficient, rank-one updating [20] that creates the mapping  $\mathcal{T}$ , see Proposition 4.2.

For the assumed positive definite  $\bar{V}_S$  and  $\bar{\nu}_S > 0$ , Proposition 4.1 implies that the partial likelihood is continuous and bounded function for  $\Lambda \in [0, 1]$ . Thus, the maximizer (18) exist. It can be simply found by using, for instance, golden-section based maximization.

## 4.3 Construction of fictitious data

The most common case of prior knowledge type is solved here as the practically used example of its quantification. Specifically, initial moments of the predictive pdfs  $f(d_{\tau_K} | \psi_{\tau_K})$  are assumed to be given for a fixed regression vector  $\psi_{\tau_K}$  (index of the fixed structure  $\mathcal{S}$  is suppressed)

$$\begin{aligned} \hat{d}_{\tau_K} &= \int d_{\tau_K} f(d_{\tau_K} | \psi_{\tau_K}) dd_{\tau_K} \\ {}^{\text{L}}d_{r_{\tau_K}} &= \int (d_{\tau_K} - \hat{d}_{\tau_K})^2 f(d_{\tau_K} | \psi_{\tau_K}) dd_{\tau_K}. \end{aligned} \quad (29)$$

It corresponds to the knowledge  $h(\psi_{\tau_K}) = [\hat{d}_{\tau_K}, {}^{\text{L}}d_{r_{\tau_K}}]$ ,  $H(d_{\tau_K}, \psi_{\tau_K}) = [d_{\tau_K}, (d_{\tau_K} - \hat{d}_{\tau_K})^2]$  in (20). For the ARX model, the restriction on the constructed prior pdf (29) becomes, Proposition 4.1,

$$\hat{d}_{\tau_K} = \hat{\theta}'_{\tau_K} \psi_{\tau_K}, \quad {}^{\text{L}}d_{r_{\tau_K}} = \hat{r}_{\tau_K} (1 + \zeta_{\tau_K}), \quad \zeta_{\tau_K} = \psi'_{\tau_K} C_{\tau_K} \psi_{\tau_K}, \quad (30)$$

where  $\hat{\theta}'_{\tau_K}, \hat{r}_{\tau_K}, C_{\tau_K}$  are LS equivalents of the statistics  $V_{\tau_K}$  resulting from the minimization of the KL distance to the pre-prior pdf.

Typically, the expert provides the range  $[\underline{d}_{\tau_K}, \bar{d}_{\tau_K}]$  of the response  $d_{\tau_K}$  on the stimulus coded in  $\psi_{\mathcal{S}; \tau_K}$ . Then, neglecting a small asymmetry of the Student pdf, we choose  $\hat{d}_{\tau_K} = 0.5 (\underline{d}_{\tau_K} + \bar{d}_{\tau_K})$  and  ${}^{\text{L}}d_{r_{\tau_K}} = [0.5 (\bar{d}_{\tau_K} - \underline{d}_{\tau_K})]^2$ .

The pre-prior pdf used in the minimization is assumed in the form

$$\bar{f}(\Theta) = \mathcal{G}_\Theta(I, \text{diag}[\underbrace{\text{l}^d_\varepsilon, \varepsilon [1, \dots, 1]}_{\hat{\psi}}], \bar{\nu}). \quad (31)$$

It is given by positive scalars  $\text{l}^d_\varepsilon, \varepsilon, \bar{\nu}$ . Such a pdf expresses the common knowledge that the parameters are finite and, importantly, *this knowledge is preserved for all nested structures*.

Using Proposition 4.1 for the pre-prior pdf (31), the optimized Kullback-Leibler distance is (the subscript  $\tau_K$  is also temporarily suppressed whenever possible)

$$\begin{aligned} 2\mathcal{D}(f||\bar{f}) &= \omega(\nu) + \varepsilon \text{tr}[C] - \ln|C| - \bar{\nu} \ln(1 + \zeta) + \bar{\nu} \ln(\hat{r}) \\ &+ \frac{\nu}{(\nu - 2)\hat{r}} \left[ \varepsilon \hat{\theta}' \hat{\theta} + \text{l}^d_\varepsilon \right], \end{aligned} \quad (32)$$

where  $\omega(\nu)$  includes all terms depending on  $\nu$  only. The optimization of this function with respect to  $\nu$  is rather involved and, importantly, its results do not have an intuitive support. This makes us to minimize the function (32) with respect to the remaining arguments only and interpret the results as updating of  $\bar{V}$  by a rank-one  $\text{l}^\Delta V$  matrix defined by a *fictitious data vector*  $\Psi$ . Then, the restrictions (30) determine even  $\text{l}^\Delta \nu$  uniquely.

**Proposition 4.3 (Optimal fictitious data vector)** *The LS quantities  $\hat{\theta}$ ,  $C$ ,  $\hat{r}$  minimizing the function (32) under the restrictions (30) are obtained by the least-squares-type updating of the pre-prior statistics (31) using the fictitious data vector*

$$\Psi'_{\tau_K} = \left[ \hat{d} \left( \frac{\rho}{\sqrt{x}} + \sqrt{x} \right), \sqrt{x} \psi' \right], \quad \rho = \frac{\varepsilon}{\psi' \psi}. \quad (33)$$

The weight  $x > 0$  is determined by the following formulas

$$\begin{aligned} q &= \frac{\nu_{\tau_K}}{\nu_{\tau_K} - 2} (\alpha \rho + \beta), \quad \alpha = \frac{\hat{d}^2}{\text{l}_{d_r}}, \quad \beta = \frac{\text{l}^d_\varepsilon}{\text{l}_{d_r \tau_K}} > 0 \\ b &= \rho + 1 - q + \bar{\nu}, \quad c = \rho(-\bar{\nu} + q) + q, \quad x = 0.5 \left( -b + \sqrt{b^2 + 4c} \right). \end{aligned} \quad (34)$$

The corresponding  $\nu_{\tau_K}$  is specified by the implicit formula

$$\nu_{\tau_K} = 2 + (1 + \rho + x) \left( \frac{\beta}{\rho + x} + \frac{\alpha \rho}{x} \right). \quad (35)$$

The coupled equations (34) (35) have a unique solution in the meaningful domain  $x > 0$ ,  $\nu_{\tau_K} > 2$ .



*Proof:* The minimization of the function (32) with respect to  $\hat{\theta}$  gives directly

$$\hat{\theta} = \frac{\hat{d}\psi}{\psi'\psi} \quad (36)$$

irrespective of other variables. Inserting this  $\hat{\theta}$  into the optimized function (32) and using the second restriction in (30) for expressing  $\hat{r} = \text{ld}_r/(1 + \zeta)$ , we get the following function minimized with respect to  $C$

$$\begin{aligned} 2\mathcal{D}(f||\bar{f}) &= \omega(\nu) + \varepsilon \text{tr}[C] - \ln|C| - \bar{\nu} \ln(1 + \zeta) + (1 + \zeta)q \\ q &= \frac{\nu}{\nu - 2} \left[ \underbrace{\rho \frac{\hat{d}^2}{\text{ld}_r}}_{\alpha} + \underbrace{\frac{\text{ld}_\varepsilon}{\text{ld}_r}}_{\beta} \right], \quad \rho = \frac{\varepsilon}{\psi'\psi}. \end{aligned}$$

Taking its derivatives with respect to  $C$  and using the identity  $\frac{\partial \ln|C|}{\partial C} = C^{-1}$ , we get the necessary condition for minimum

$$C^{-1} = \varepsilon I + x\psi\psi' \text{ with } x = -\frac{\bar{\nu}}{1 + \zeta} + q. \quad (37)$$

This implicit definition,  $\zeta = \psi' C \psi$ , is resolved using the formula  $(\varepsilon I + x\psi\psi')^{-1} = \varepsilon^{-1} [I - \frac{x\psi\psi'}{\varepsilon + x\psi'\psi}]$ . It gives the equation  $x = \frac{-\bar{\nu}(\rho+x)}{\rho+x+1} + q$ , that converts into the quadratic one for  $x$ ,  $x^2 + \underbrace{[\rho + 1 - q + \bar{\nu}]}_b x - \underbrace{[\rho(-\bar{\nu} + q) + q]}_c = 0$ . For  $\nu$  sufficiently close to 2 (from right), we get  $c > 0$  and the equation has the unique positive (meaningful) solution  $x = 0.5 \left( -b + \sqrt{b^2 + 4c} \right)$ .

The found form of the updating of the pre-prior covariance factor  $\bar{C} = \varepsilon^{-1} I$  (37) implies that the fictitious regression vector corresponding to the  $\tau_K$ -th part of the knowledge item is simply  $\psi_{\tau_K} = \sqrt{x}\psi$ . The derived formula for  $\hat{\theta}$  (36) is obtained if we take fictitious output  $d_{\tau_K} = \frac{\hat{d}}{\sqrt{x}}(\rho + x)$ .

The least-squares remainder, Proposition 4.1, that corresponds to this updating has the value  $\text{ld}D = \text{ld}_\varepsilon + \frac{\hat{d}^2 \rho(\rho+x)}{x}$ . At the same time, the estimate of the noise variance meeting the given restrictions has the form  $\hat{r} = \frac{\text{ld}D}{\nu-2} = \frac{\text{ld}_r(\rho+x)}{1+\rho+x}$ . Thus, the found results can be interpreted as updating by the fictitious data vector  $\Psi_{\tau_K} = [d_{\tau_K}, \psi'_{\tau_K}]'$  iff we take

$$\nu_{\tau_K} - 2 = (1 + \rho + x) \left( \frac{\beta}{\rho + x} + \frac{\alpha\rho}{x} \right).$$

Inserting the relationship between  $x$  and  $q$  from (37), we can express the searched  $\nu_{\tau_K} - 2$  occurring in the above equation as a function of  $x$ . It gives a

third order algebraic equation for  $x$ . It has always real solution. Standard but lengthy analysis shows uniqueness of the solution in the meaningful domain.  $\square$

#### 4.4 Practical examples of prior knowledge

Here, we list common prior pieces knowledge available and ways of the data vectors  $\Psi'_i = [\hat{d}_i, \psi'_i]$  construction (fixed subscripts  $\tau_K, \mathcal{S}$  are still suppressed). The multi-variate data items  $d_t$  and the common case of the state in the phase form  $[d'_{t-1}, \dots, d'_{t-\delta}, 1]$  of the order  $\delta$  are considered. The structure of the data vector is described by the index  $i$ , pointing to  $i$ -th output channel, and by the list  $l_i$  of indices  $(j, \delta)$  with  $j \in \{1, \dots, d\}$  and  $\delta \in \{0, \dots, \delta\}$ . The indices say that the data on the  $j$ -th channel  $d_{j;t-\delta}$  are in the constructed regression vector  $\psi_i$ .

In all cases discussed below, the *entries of  $\psi_i$  that are not explicitly mentioned are set to zero.*

The knowledge of the *static gain*  $\hat{d}_i = \hat{g}$  of the  $i$ -th channel on stimulus from the  $j$ -th channel is expressed by setting  $d_{i;t-\delta} = \hat{g}$  and  $d_{j;t-\delta} = 1$  for all delays  $\delta$  in the list  $l_i$ .

It is shown in [21] that the knowledge of a point of the *frequency response* on the  $j$ -th channel – given by the module  $\hat{a}(\omega)$  and phase  $\phi(\omega)$  shift at frequency  $\omega$  – is determined by a pair of data vectors with  $d_{i;\delta} = \hat{a}(\omega) \sin(\phi(\omega) + \delta\omega)$ ,  $d_{j;\delta} = \sin(\delta\omega)$  and  $d_{i;\delta} = \hat{a}(\omega) \cos(\phi(\omega) + \delta\omega)$ ,  $d_{j;\delta} = \sin(\delta\omega)$ . The range of  $a(\omega) = [\underline{a}(\omega), \bar{a}(\omega)]$  can often be well specified. The uncertainty in the phase  $\phi(\omega)$  is simply reflected by considering a relatively coarse grid within the uncertainty range and processing each case as an individual data item. The subsequent merging (15) cares about the proper weighting.

Knowledge of *cut-off frequency* is a special case of frequency knowledge with practically zero amplitude, i.e. the amplitude range is given by the point estimate of the standard deviation of the noise. Introduction of this knowledge for several frequencies higher than the cut-off one excludes isolated pass through the zero level. Again, it generates several competitive knowledge items balanced by the merging.

The knowledge of a range of the *dominant time constant* is implemented by modelling lower and upper envelope on the impulse response generated by the first order models with time constants equal to the specified bounds on the time constant. Data are filled from the average trajectory into  $\Psi_i$  while distance of envelopes determine the variance  $\mathbb{L}^{d_r}$ .

*Envelope of measured data*, obtained from a periodic measurement, is handled in the same way.

*Smoothness of the step response* [22] can be respected by enforcing its second order difference to be close to zero.

Note, that the lengths of the samples in "simulated" responses have to be limited so that stationary values are not repeated too much. Otherwise, the assumption on internal consistency, i.e. applicability of the Bayes rule, would be violated.

#### 4.5 Overall algorithm for normal ARX model

Here, we put the algorithmic element together for the normal ARX model. The recommended options correspond with the pre-processed data  $d(\hat{t})$  without outliers, suppressed measurement noise and *data scaled* so that their means are approximately zero and variances are about one. The evaluation is organized so that computation resources are preserved as much as possible.

The explicit reference is made here to the channel treated (subscript  $i$ ).

##### Algorithm 4.1 (Structure estimation with prior pdf)

Initial phase

- *Select the grids on:*  
*forgetting factors  $\{\lambda\}$ , used for processing of internally consistent data blocks,*  
*phases  $\{\phi(\omega)\}$  that complete definitions of the frequency response,*  
*frequencies  $\{\omega_c\}$  that guarantee that frequency response is close to zero above the cut-off frequency.*
- *Select the number of repetitive starts in the nested structure estimation algorithm SEN 2.1.*
- *Select the order  $\hat{\delta}_{\mathcal{R}}$  of the richest data vector  $\Psi_{\mathcal{R};t} = [d'_t, \dots, d'_{t-\hat{\delta}_{\mathcal{R}}}, 1]'$  that includes all potential entries when predicting all modelled entries  $d_{i;t}$ ,  $i = 1, \dots, \hat{d}$ , of the data item  $d_t$ .*
- *Specify statistics  $\bar{L}_{\mathcal{R}} = I$ ,  $\bar{D}_{\mathcal{R}} = \text{diag}[\text{ }^{\text{d}}\varepsilon, \varepsilon \overbrace{[1, \dots, 1]}^{\hat{\psi}_{\mathcal{R}}}]$  and  $\bar{\nu}$  determining the flat pre-prior pdf on the richest possible parameterization.*

*Requirements on finiteness of the a priori assigned expectation of  $r$  and need for a flat pdf  $\bar{f}(r)$  hint to use  $\bar{\nu} = 3$ , see (27). For this choice,  $\text{ }^{\text{d}}\varepsilon$  is the variance of the unpredictable part of the modelled data. It is sufficient to consider a few categories of the noise-to-signal ratio.*

For instance, the values (0.1%, 1%, 5%, 10%, 50%) correspond with  $\mathbb{L}^d_\varepsilon = (1e-6, 1e-4, 2.5e-3, 1e-2, 0.25)$ .

For stable systems, that are predominantly treated, the auto-regression coefficients do not cross the value  $\gamma = \begin{pmatrix} \delta_{\mathcal{R}} \\ 0.5\delta_{\mathcal{R}} \end{pmatrix}$ . The regression coefficients are, as a rule, much smaller. Properties of the GiW pdf imply that  $\varepsilon \approx 25 \mathbb{L}^d_\varepsilon / \gamma^2$  is an appropriate conservative option.

- Select the number  $\hat{S}$ , say several tens, of competitive structures to be refined by using prior information.
- Use the available real data  $d_t, t = 1, \dots, \hat{t}$  to update  $L'DL$  decomposition of the increment of the extended information matrix corresponding to the richest data vectors  $\Psi_{\mathcal{R};t}$ , see (13),

$$\mathbb{L}^{\Delta} L'_{\mathcal{R};t} \mathbb{L}^{\Delta} D_{\mathcal{R};t} \mathbb{L}^{\Delta} L_{\mathcal{R};t} = \begin{bmatrix} \mathbb{L}^{\Delta} L_{\mathcal{R};t-1} \\ \Psi'_{\mathcal{R};t} \end{bmatrix}' \text{diag} \begin{bmatrix} \mathbb{L}^{\Delta} D_{\mathcal{R};t-1} \\ 1 \end{bmatrix} \begin{bmatrix} \mathbb{L}^{\Delta} L_{\mathcal{R};t-1} \\ \Psi_{\mathcal{R};t} \end{bmatrix}$$

$$\mathbb{L}^{\Delta} \nu_{\mathcal{R};t} = \mathbb{L}^{\Delta} \nu_{\mathcal{R};t-1} + 1, \text{ with } \mathbb{L}^{\Delta} L_{\mathcal{R};0} = I, \mathbb{L}^{\Delta} D_{\mathcal{R};0} = 0, \mathbb{L}^{\Delta} \nu_{\mathcal{R};0} = 0.$$

- Evaluate the  $L'DL$  decomposition of the extended information matrix corresponding to the richest data vectors  $\Psi_{\mathcal{R};t}$ , i.e. add the statistics of the pre-prior pdf

$$L'_{\mathcal{R};\hat{t}} D_{\mathcal{R};\hat{t}} L_{\mathcal{R};\hat{t}} = \begin{bmatrix} \mathbb{L}^{\Delta} L'_{\mathcal{R};\hat{t}} \\ I \end{bmatrix}' \text{diag} \begin{bmatrix} \mathbb{L}^{\Delta} D_{\mathcal{R};\hat{t}} \\ \bar{D}_{\mathcal{R}} \end{bmatrix} \begin{bmatrix} \mathbb{L}^{\Delta} L_{\mathcal{R};\hat{t}} \\ I \end{bmatrix}.$$

- Set the channel index  $i = 1$ .

Cycle over indices  $i$  of the modelled entries in data records

- Set the auxiliary description of the structure  $\hat{S} = \emptyset$ ,  $\hat{\mathcal{L}}_i = -\infty$  needed for the MAP estimation.

Structure estimation with nested prior knowledge

- Select the factors of the pre-prior and posterior extended information matrices,  $\bar{L}_{i\mathcal{R}}$ ,  $\bar{D}_{i\mathcal{R}}$ ,  $L_{i\mathcal{R};\hat{t}}$ ,  $D_{i\mathcal{R};\hat{t}}$ , as well as of the increment  $\mathbb{L}^{\Delta} L_{i\mathcal{R};\hat{t}}$ ,  $\mathbb{L}^{\Delta} D_{i\mathcal{R};\hat{t}}$ , corresponding to the  $i$ -th predicted data entry  $d_{i;\hat{t}}$  and the richest regression vector  $\psi_{i\mathcal{R};\hat{t}}$ .

They are nested in  $L_{\mathcal{R};\hat{t}}$ ,  $D_{\mathcal{R};\hat{t}}$ ,  $\bar{L}$ ,  $\bar{D}$ , and  $\mathbb{L}^{\Delta} L_{\mathcal{R};\hat{t}}$ ,  $\mathbb{L}^{\Delta} D_{\mathcal{R};\hat{t}}$ . The  $L'DL$  decompositions spoiled by this selection have to be recovered using the rank-one corrections [20].

- Apply the algorithm *SEN* (2.1) giving *L'**DL* factors of the pre-prior extended information matrices  $\bar{L}_{iS}$ ,  $\bar{D}_{iS}$  and their data-dependent increments  ${}^{\Delta}L_{iS;\hat{i}}$ ,  ${}^{\Delta}D_{iS;\hat{i}}$ . They correspond to the most probable structures  $S \in \mathcal{S}^*$  found when just the nested pre-prior knowledge is used.
- Select  $S \in \mathcal{S}^*$ .

Inclusion of prior knowledge for promising structures

- Select a knowledge item  $\mathcal{K}$  in the list  $\mathcal{K}_i^* = \{1, \dots, \mathring{\mathcal{K}}_i\}$  available for the  $i$ -th channel.
- Set the normalization factor needed in merging  $s = 0$ .

Processing of knowledge items

- Do if the individual knowledge item  $\mathcal{K}$  has to be converted into fictitious data vectors

- Set  $L_{iS;0} = \bar{L}_{iS}$ ,  $D_{iS;0} = \bar{D}_{iS}$ ,  $\nu_{iS;0} = \bar{\nu}$ .
- For  $\tau_{\mathcal{K}} = 1, \dots, \mathring{\tau}_{\mathcal{K}}$ 
  - \* Generate data reflecting  $\tau_{\mathcal{K}}$ -th part of the knowledge item  $\mathcal{K}$  given by  $\hat{d}_{iS;\tau_{\mathcal{K}}}$ ,  ${}^{\Delta}r_{iS;\tau_{\mathcal{K}}}$  and  $\psi_{iS;\tau_{\mathcal{K}}}$ , cf. (29).
  - \* Evaluate fictitious data vectors  $\Psi_{iS;\tau_{\mathcal{K}}}$  and its  $\nu_{iS;\tau_{\mathcal{K}}}$ , cf. (33),

$$\Psi'_{iS;\tau_{\mathcal{K}}} = \left[ \hat{d}_{iS;\tau_{\mathcal{K}}} \left( \frac{\rho_{iS;\tau_{\mathcal{K}}}}{\sqrt{x_{iS;\tau_{\mathcal{K}}}}} + \sqrt{x_{iS;\tau_{\mathcal{K}}}} \right), \sqrt{x_{iS;\tau_{\mathcal{K}}}} \psi'_{iS;\tau_{\mathcal{K}}} \right].$$

\* Update

$$\begin{aligned} L'_{iS;\tau_{\mathcal{K}}} D_{iS;\tau_{\mathcal{K}}} L_{iS;\tau_{\mathcal{K}}} &= \begin{bmatrix} L_{iS;\tau_{\mathcal{K}}-1} \\ \Psi'_{iS;\tau_{\mathcal{K}}} \end{bmatrix}' \text{diag} \begin{bmatrix} D_{iS;\tau_{\mathcal{K}}-1} \\ 1 \end{bmatrix} \begin{bmatrix} L_{iS;\tau_{\mathcal{K}}-1} \\ \Psi'_{iS;\tau_{\mathcal{K}}} \end{bmatrix} \\ \nu_{iS;\tau_{\mathcal{K}}} &= \nu_{iS;\tau_{\mathcal{K}}-1} + {}^{\Delta}\nu_{iS;\tau_{\mathcal{K}}}. \end{aligned}$$

- Run the estimation with the stabilized forgetting for the selected forgetting factors  $\lambda$  and with the alternative pdf given by  $\bar{L}_{iS}$ ,  $\bar{D}_{iS}$ ,  $\bar{\nu}$  if the knowledge item  $\mathcal{K}$  is formed by an internally consistent data block.

Store the results into  ${}^{\Delta}L_{iS;\mathring{\tau}_{\mathcal{K}}}$ ,  ${}^{\Delta}D_{iS;\mathring{\tau}_{\mathcal{K}}}$  and  ${}^{\Delta}\nu_{iS;\mathring{\tau}_{\mathcal{K}}}$

- Select the flattening factor  $\Lambda_{i\mathcal{SK}}$  maximizing the partial likelihood

$$\mathcal{L}_i(d(\hat{t}), \mathcal{S}, \mathcal{K}) = \frac{\mathcal{I}\left(V_{i\mathcal{SK};\hat{t}}, \hat{t} + \Lambda_{i\mathcal{SK}} \text{}^{\Delta}\nu_{i\mathcal{S};\hat{\tau}_{\mathcal{K}}} + \bar{\nu}_{i\mathcal{S}}\right)}{\mathcal{I}\left(\Lambda_{i\mathcal{SK}} \text{}^{\Delta}L'_{i\mathcal{S};\hat{\tau}_{\mathcal{K}}} \text{}^{\Delta}D_{i\mathcal{S};\hat{\tau}_{\mathcal{K}}} \text{}^{\Delta}L_{i\mathcal{S};\hat{\tau}_{\mathcal{K}}} + \bar{L}'_{i\mathcal{S}} \bar{D}_{i\mathcal{S}} \bar{L}_{i\mathcal{S}}, \Lambda_{i\mathcal{SK}} \text{}^{\Delta}\nu_{i\mathcal{S};\hat{\tau}_{\mathcal{K}}} + \bar{\nu}_{i\mathcal{S}}\right)}$$

$$V_{i\mathcal{SK};\hat{t}} = \text{}^{\Delta}L'_{i\mathcal{S};\hat{t}} \text{}^{\Delta}D_{i\mathcal{S};\hat{t}} \text{}^{\Delta}L_{i\mathcal{S};\hat{t}} + \Lambda_{i\mathcal{SK}} \text{}^{\Delta}L'_{i\mathcal{S};\hat{\tau}_{\mathcal{K}}} \text{}^{\Delta}D_{i\mathcal{S};\hat{\tau}_{\mathcal{K}}} \text{}^{\Delta}L_{i\mathcal{S};\hat{\tau}_{\mathcal{K}}} + \bar{L}'_{i\mathcal{S}} \bar{D}_{i\mathcal{S}} \bar{L}_{i\mathcal{S}}$$

$$s = s + \mathcal{L}_i(d(\hat{t}), \mathcal{S}, \mathcal{K})$$

Notice that  $\text{}^{\Delta}L'_{i\mathcal{S};\hat{t}}$ ,  $\text{}^{\Delta}D_{i\mathcal{S};\hat{t}}$  and  $\bar{L}'_{i\mathcal{S}}$ ,  $\bar{D}_{i\mathcal{S}}$ ,  $\bar{\nu}_{i\mathcal{S}}$  depend only on the structural indices  $i, \mathcal{S}$  and not on  $\mathcal{K}$ .

- Take a new knowledge item  $\mathcal{K}$ , if the list  $SK_i$  is not exhausted, and go to Processing of knowledge items. Otherwise continue.
- Set  $L_{i\mathcal{SK}(\hat{\mathcal{K}})} = I$ ,  $D_{i\mathcal{SK}(\hat{\mathcal{K}})} = \bar{D}_{i\mathcal{S}}$ ,  $\nu_{i\mathcal{SK}(\hat{\mathcal{K}})} = \bar{\nu}_{i\mathcal{S}}$ .
- Select  $\mathcal{K} \in \mathcal{K}_i^*$ . Evaluation of the merger within the structure  $\mathcal{S}$
- Update

$$f_i(\mathcal{K}|d(\hat{t}), \mathcal{S}) = \frac{\mathcal{L}_i(d(\hat{t}), \mathcal{S}, \mathcal{K})}{s}, \quad L'_{i\mathcal{SK}(\hat{\mathcal{K}})} D_{i\mathcal{SK}(\hat{\mathcal{K}})} L_{i\mathcal{SK}(\hat{\mathcal{K}})}$$

$$= L'_{i\mathcal{SK}(\hat{\mathcal{K}})} D_{i\mathcal{SK}(\hat{\mathcal{K}})} L_{i\mathcal{SK}(\hat{\mathcal{K}})} + f_i(\mathcal{K}|d(\hat{t}), \mathcal{S}) \text{}^{\Delta}L'_{i\mathcal{S};\hat{\tau}_{\mathcal{K}}} \text{}^{\Delta}D_{i\mathcal{S};\hat{\tau}_{\mathcal{K}}} \text{}^{\Delta}L_{i\mathcal{S};\hat{\tau}_{\mathcal{K}}}$$

$$\nu_{i\mathcal{SK}(\hat{\mathcal{K}})} = \nu_{i\mathcal{SK}(\hat{\mathcal{K}})} + f_i(\mathcal{K}|d(\hat{t}), \mathcal{S}) \text{}^{\Delta}\nu_{i\mathcal{S};\hat{\tau}_{\mathcal{K}}}.$$

- Select a new  $\mathcal{K}$  if their list  $\{1, \dots, \hat{\mathcal{K}}_i\}$  is not exhausted and go to Evaluation of the merger within the structure  $\mathcal{S}$ . Otherwise continue.
- Evaluate partial likelihood expressing assigned to the structure  $\mathcal{S}$

$$\mathcal{L}_i(d(\hat{t}), \mathcal{S}) = \frac{\mathcal{I}\left(\text{}^{\Delta}L'_{i\mathcal{S};\hat{t}} \text{}^{\Delta}D_{i\mathcal{S};\hat{t}} \text{}^{\Delta}L_{i\mathcal{S};\hat{t}} + L'_{i\mathcal{SK}(\hat{\mathcal{K}})} D_{i\mathcal{SK}(\hat{\mathcal{K}})} L_{i\mathcal{SK}(\hat{\mathcal{K}})}, \hat{t} + \nu_{i\mathcal{SK}(\hat{\mathcal{K}})}\right)}{\mathcal{I}\left(L'_{i\mathcal{SK}(\hat{\mathcal{K}})} D_{i\mathcal{SK}(\hat{\mathcal{K}})} L_{i\mathcal{SK}(\hat{\mathcal{K}})}, \nu_{i\mathcal{SK}(\hat{\mathcal{K}})}\right)}.$$

- Set  $\hat{\mathcal{S}} = \mathcal{S}$ ,  $\hat{\mathcal{L}}_i = \mathcal{L}_i(d(\hat{t}), \mathcal{S})$  and store the statistics corresponding to the posterior pdf  $\text{}^{\Delta}L'_{i\hat{\mathcal{S}};\hat{t}} \text{}^{\Delta}D_{i\hat{\mathcal{S}};\hat{t}} \text{}^{\Delta}L_{i\hat{\mathcal{S}};\hat{t}} + L'_{i\hat{\mathcal{S}}\mathcal{K}(\hat{\mathcal{K}})} D_{i\hat{\mathcal{S}}\mathcal{K}(\hat{\mathcal{K}})} L_{i\hat{\mathcal{S}}\mathcal{K}(\hat{\mathcal{K}})}$  and  $\hat{t} + \nu_{i\hat{\mathcal{S}}}$  if  $\mathcal{L}_i(d(\hat{t}), \mathcal{S}) > \hat{\mathcal{L}}_i$ .

- *Select a new structure  $\mathcal{S}$  if the list  $\mathcal{S}^*$  of the most probable ones is not exhausted and go to the Inclusion of prior knowledge for promising structures. Otherwise continue.*
- *Offer the structure  $\hat{\mathcal{S}}$  as the recommended one for  $i$ -th channel with the corresponding stored posterior statistics.*
- *Increase  $i$  and go to Cycle over indices  $i$  of modelled data entries if  $i \leq \mathring{d}$ . Otherwise stop.*

The algorithm provides also, till unavailable, estimate of the best order in factorization (2). It must, however, be complemented by a check against incorrect dependence loops.

It is important that the algorithm can cope with problems having  $\mathring{d}$  of the order several tens.

## 5 Illustrative example

Contribution of prior knowledge to structure estimation results is illustrated on single-input single-output simulated system. It corresponds to two-dimensional  $d_t = [y_t, u_t]'$ . The system input  $u_t$  is generated as white normal noise with variance 0.3. The modelled system output  $y_t$  is *simulated* by the ARX model determined by the “objective parameter”  ${}^{\circ}\Theta = [{}^{\circ}\theta, {}^{\circ}r]$ . It is usually written in form of the difference equation driven by the normalized white normal noise  $e_t \sim \mathcal{N}_{e_t}(0, 1)$

$$\begin{aligned}
y_t &= {}^{\circ}a_1 y_{t-1} + {}^{\circ}a_2 y_{t-2} + {}^{\circ}b_0 u_t + {}^{\circ}b_1 u_{t-1} + \sqrt{{}^{\circ}r} e_t \\
&= \underbrace{[1.81, 0.8187, 0.0438, 0.00468]}_{{}^{\circ}\theta'} \underbrace{[y_{t-1}, y_{t-2}, u_t, u_{t-1}]}_{\psi_t} + \underbrace{\sqrt{0.001}}_{\sqrt{{}^{\circ}r}} e_t \\
&\Leftrightarrow f(y_t | u_t, d(t-1), {}^{\circ}\Theta) = f(y_t | \psi_t, {}^{\circ}\Theta) = \mathcal{N}_{y_t}({}^{\circ}\theta' \psi_t, {}^{\circ}r).
\end{aligned}$$

The “real” data  $d(\mathring{t}) = d(300)$  were “measured” on this system (MATLAB simulation with the seed of normal generator equal to 1).

The SEN algorithm 2.1 was applied with the richest structure of the phase form given by the order  $\mathring{\delta}_{\mathcal{R}}$  and the recommended nested pre-prior pdf were used. The number of restarts was 10 and  $\mathring{\mathcal{S}} = 10$  of the best structures were stored giving the significant entries marked by \* and the related posterior pfs

S	$a_1$	$a_2$	$a_3$	$b_0$	$b_1$	$b_2$	$f(\mathcal{S} d(300))$
1	*	*					0.6935
2	*	*			*		0.2572
3	*	*	*				0.0185
4	*	*		*			0.0111
5	*	*		*	*		0.0077
$\vdots$			$\vdots$				$\vdots$

The correct structure was assigned the 5th largest probability. Inclusion of the prior knowledge  $\mathcal{K} = 1 \equiv$  static gain is in range  $[1.0310 \quad 1.0518]$  made the largest pf  $f(\mathcal{S} = 2|d(300), \mathcal{K} = 1)$  corresponding to the second order model with non-zero  $b_1$ .

... ..

The following analysis is based on data sample only and on the prior knowledge about the system:

1. the system static gain is in the range
2. the dominant time constants in ranges:  
 $\lambda_1 = [0.8254 \quad 0.8421]$ ;  $\lambda_2 = [1.1237 \quad 1.3734]$

The conclusion are:

- the SEN algorithm of nested structure estimation does not discover the dependency of output on input;
- the "true" structure (appears among "best" structures with a low probability (the last in the table above).

In this case, dependency of output on input is discovered. The prior knowledge of static gain carries the information about dependency but not about dynamics - it can explain the absence of the  $b_0$  in the result.

Now, the knowledge of the time constants is applied. The prior knowledge carries the information about dependency and system dynamics. This time, the "true" structure  $a_1, a_2, b_0, b_1$  is discovered.

A range of prior knowledge and their combinations were tested in the task of the structure estimation. The experiments show that use of prior pdf is not a magic tool that can solve any problem – it only slightly enlarges the range of successful solutions.

CCC>> MK: information on some improvement to former version, some information on weights <<CCC



## 6 Concluding remarks

Significance of the inclusion of prior knowledge into black-box models is still underestimated. The theory and algorithms presented in this paper solve this problem to a significant extent. The available practical experience confirms that the use of even vague knowledge may decide on the success or failure of structure estimation and consequently on the success or failure of the controller design.

A wider and more precise use of the prior knowledge is especially important in the context of the prior design of the advanced controller with incomplete knowledge [15].

The theory and algorithms have been elaborated for the LQG-type-design set up. The same problems are met out of this class and the adopted methodology may serve to them, too.

The problem of the joint processing of the prior knowledge and structure estimation is solved here for the exponential family but elaborated for the ARX models only. It determines a direction of a further development: the controlled-Markov-chain case should be elaborated, too. It would strengthen dynamic modelling of systems with dynamic discrete data.

### Acknowledgements

This research has been partially supported by grants GA AVČR S1075102, S1075351 and GAČR 102/03/0049.

## References

- [1] K.J. Astrom and B. Wittenmark, *Adaptive Control*, Addison-Wesley Publishing Company, Reading, Massachusetts, 1989.
- [2] P.E. Wellstead and M.B. Zarrop, *Self-tuning Systems*, John Wiley & Sons, Chichester, 1991.
- [3] E. Mosca, *Optimal, Predictive, and Adaptive Control*, Prentice Hall, 1994.
- [4] M. Kárný, A. Halousková, J. Böhm, R. Kulhavý, and P. Nedoma, “Design of linear quadratic adaptive control: Theory and algorithms for practice”, *Kybernetika*, vol. 21, 1985, Supplement to No. 3, 4, 5, 6.
- [5] A. Johnson, “LQG applications in the process industries”, *Chemical Engineering Practice*, vol. 48, no. 16, pp. 2829–2838, 1993.

- [6] M. Elbelkacemi, A. Lachhab, M. Limouri, D. Dahhou, and A. Essaid, “Adaptive control of a water supply system”, *Control Engineering Practice*, vol. 9, no. 3, pp. 343–349, 2001.
- [7] D.W. Clarke, *Advances in Model-Based Predictive Control*, Oxford University Press, Oxford, 1994.
- [8] M. Kárný and A. Halousková, “Implementing LQG adaptive control: a CAD approach”, in *Preprints of the 9th IFAC/IFORS Symposium on Identification and System Parameter Estimation*, pp. 1585–1590. AKA PRINT Nyomdaipari, Budapest, 1991.
- [9] M. Kárný and A. Halousková, “Pre-tuning of self-tuners”, in *Advances in Model-Based Predictive Control*, D. Clarke, Ed., pp. 333–343. Oxford University Press, Oxford, 1994.
- [10] J. Bůcha, M. Kárný, P. Nedoma, J. Böhm, and J. Rojíček, “Designer 2000 project”, in *International Conference on Control '98*, London, September 1998, pp. 1450–1455, IEE.
- [11] M. Kárný and R. Kulhavý, “Structure determination of regression-type models for adaptive prediction and control”, in *Bayesian Analysis of Time Series and Dynamic Models*, J.C. Spall, Ed. Marcel Dekker, New York, 1988, chapter 12.
- [12] V. Peterka, “Bayesian system identification”, in *Trends and Progress in System Identification*, P. Eykhoff, Ed., pp. 239–304. Pergamon Press, Oxford, 1981.
- [13] R. Kulhavý and M. B. Zarrop, “On general concept of forgetting”, *International Journal of Control*, vol. 58, no. 4, pp. 905–924, 1993.
- [14] M. Kárný, T. Jeníček, and W. Ottenheimer, “Contribution to prior tuning of LQG selftuners”, *Kybernetika*, vol. 26, no. 2, pp. 107–121, 1990.
- [15] M. Novák and J. Böhm, “Automated multivariable adaptive controller design”, *IEEE Trans. on Control theory and its applications*, 2003, This issue.
- [16] J.M. Bernardo and A.F.M. Smith, *Bayesian theory*, John Wiley & Sons, Chichester, New York, Brisbane, Toronto, Singapore, 1997, 2nd edition.

- [17] P. Nedoma, M. Kárný, and J. Böhm, “Software tools for use of prior knowledge in design of LQG adaptive controllers”, in *Preprints of the IFAC Workshop on Adaptive Systems in Control and Signal Processing*, Glasgow, August 1998, pp. 425–429, IFAC.
- [18] S. Kullback and R. Leibler, “On information and sufficiency”, *Annals of Mathematical Statistics*, vol. 22, pp. 79–87, 1951.
- [19] A. Zellner, *An Introduction to Bayesian Inference in Econometrics*, J. Wiley, New York, 1976.
- [20] G.J. Bierman, *Factorization Methods for Discrete Sequential Estimation*, Academic Press, New York, 1977.
- [21] N. Khaylova, “Exploitation of prior knowledge in adaptive control design”, Tech. Rep., FAV ZČU, University of West Bohemia, Pilsen, Czech Republic, 2001, PhD Thesis.
- [22] M. Kárný, “Quantification of prior knowledge about global characteristics of linear normal model”, *Kybernetika*, vol. 20, no. 5, pp. 376–385, 1984.