

## ONE-SIDED APPROXIMATION OF BAYES RULE AND ITS APPLICATION TO REGRESSION MODEL WITH CAUCHY NOISE

MIROSLAV KÁRNÝ, KATALIN M. HANGOS

The paper presents results which aim to broaden applicability of Bayesian identification to non-standard problems by using a systematic approach to the design of a feasible approximation of the optimal, but unfeasible solution.

The essence of the theory, admitting to generate global approximants in real time, is outlined. The approximation task formulated as a properly chosen problem of mathematical programming is applied to the practically important case of the regression with heavy tailed noise.

### 1. INTRODUCTION

Bayesian statistics has proved to be a suitable basis for understanding of a broad variety of identification tasks, for instance, parameter estimation [11], and tracking [8] in closed control loop or model comparison and system structure determination [4], [7]. Practical utility of the approach is, however, restricted to a relatively narrow range of numerically feasible problems. The limitations are especially urgent whenever real-time implementation is required (as for adaptive predictors and/or controllers).

The paper presents current results of the author's systematic attempts [5], [6] to provide Bayesian identification with, hopefully, well-founded method(s) to produce a feasible, real-time implementable approximation of the theoretical solution.

In the paper, the theory presented in [6] is recalled and its application to the regression model with Cauchy-distributed noise given.

The essence of the theory can be summarized as follows: A simple upper bound on the Kullback-Leibler distance between the unfeasible conditional probability density function  $p(w | t)$  (c.p.d.f.), describing an unknown parameter  $w$ , and an admissible approximant  $\hat{p}(w | t)$ , is found. The bound is based on a one-sided, recursively feasible, (lower) bound on  $p(w | t)$ . Optimization of the bound leads to the rational choice of  $\hat{p}(w | t)$  which becomes a problem of mathematical programming.

The chosen application case, regression model with Cauchy-distributed noise,

models the important estimation problem with outliers. The outlier-robust recursive estimation of regression coefficients is gained by a straightforward but nontrivial application of the present theory.

## 2. PRELIMINARIES

Necessary notions, facts and relation are summarized in the section.

### Bayesian identification

Let data items  $d(t) = (y(t), u(t))$  composed of the system output  $y(t)$  and the system input  $u(t)$  be measured on the system at discrete time instants labeled by  $t$ . Let, moreover,  $p(y(t) | t - 1; u(t), w) = m(t | w)$  be a known model of the observed system, namely, conditional probability density function (c.p.d.f.) of the output  $y(t)$  conditioned on the data measured up to and including  $t - 1$ :  $d(1 \dots t - 1) = (d(1), d(2), \dots, d(t - 1))$ , on the current input  $u(t)$  and a finite-dimensional constant unknown parameter  $w$ . To keep notational simplicity we have confined ourselves to c.p.d.f.'s with respect to Lebesgue measure.

The objective of the Bayesian identification is the determination of the c.p.d.f.  $p(w | t) = p(w | d(1 \dots t))$  under the so-called natural conditions of control [11], i.e. under the following restriction, on the class of admissible control strategies

$$(1) \quad p(u(t) | t - 1; w) = p(u(t) | t - 1), \quad t = 1, 2, 3 \dots$$

Under the conditions (1), the required c.p.d.f.  $p(w | t)$  evolves according to the following form of the Bayes rule

$$(2) \quad p(w | t) = m(t | w) p(w | t - 1) / c(t), \quad t = 1, 2, 3 \dots$$

$$(3) \quad c(t) = \int m(t | w) p(w | t - 1) dw$$

where  $p(w | 0) = p_0(w)$  is a user-specified prior p.d.f.

### Unnormalized version of c.p.d.f. $p(w | t)$

Let us join with any c.p.d.f.  $p(w | t)$  the class of its unnormalized versions,  $\langle p(w | t) \rangle$ . A nonnegative integrable function  $l(w | t)$  belongs to  $\langle p(w | t) \rangle$  iff

$$(4) \quad 0 < \int l(w | t) dw = g(t) < \infty$$

and

$$(5) \quad p(w | t) = l(w | t) / g(t).$$

For simplicity and because of close connections, we shall call members of  $\langle p(w | t) \rangle$  the likelihood functions.

Clearly,  $\langle p(w | t) \rangle$  can be viewed as an equivalence class: two likelihood functions,

say  $l, \bar{l}$ , are equivalent if there exists a positive finite constant, say  $k$ , such that  $l = k\bar{l}$ . The following lemma states that the equivalence class  $\langle p(w | t - 1) \rangle$  is mapped by the Bayes rule on the equivalence class  $\langle p(w | t) \rangle$ , irrespectively of the used parametrization.

**Lemma 1.** Let  $l(w | t - 1) \in \langle p(w | t - 1) \rangle$ , then for any positive finite  $k$  the function

$$(6) \quad l(w | t) = k m(t | w) l(w | t - 1) \in \langle p(w | t) \rangle.$$

The feature is invariant with respect to any regular transformation of the  $w$ -space.

**Proof.** The invariancy (useful in applications) is implied by the observation that the rules for recomputation of  $p(w | t)$ , when transforming the  $w$ -space, extend to their likelihoods. The rest of the simple proof is omitted.  $\square$

#### Class of approximating c.p.d.f.'s

Let us introduce the class of approximating c.p.d.f.'s, say  $\hat{P}(t)$ , with a generic member  $\hat{p}(w | t)$  or  $\bar{p}(w | t)$ . The class  $\hat{P}(t)$  is characterized by the two requirements:

- $\hat{P}(t) \subset P(t) =$  c.p.d.f.'s on the  $w$ -space, an approximant is required to be also the c.p.d.f.,
- the functions  $\hat{p}(w | t) \in \hat{P}(t)$  are "feasible" from the computational viewpoint.

Each approximant (being a p.d.f.) generates an equivalence class of its likelihoods. The collection of all (approximating) likelihoods will be denoted by  $\langle P(t) \rangle$ ,  $\langle \hat{P}(t) \rangle$ , and its generic members by  $l(w | t)$ ,  $(\bar{l}(w | t), \tilde{l}(w | t))$ , respectively.

#### The Kullback-Leibler distance

Let us choose a distance measure the need for which arises in the approximation task. The (dis)similarity measure of probability density functions

$$(7) \quad I(\hat{p}, p) = \int \hat{p}(w | t) \ln (\hat{p}(w | t)/p(w | t)) dw$$

called the Kullback-Leibler distance (or relative Shannon entropy) is known to have a lot of nice properties. For instance, there are estimates of an increase of the risk in any decision task due to an approximation involved [10]. They can be used to estimate the impact of the approximation on the quality of prediction and/or control.

For our treatment, the following elementary facts are needed [9]:

- the distance  $I$  is always nonnegative being zero only if the compared functions differ at most on a zero-measure set,
- the distance  $I = \infty$  if there is a measurable set  $A$  such that  $\int_A \hat{p}(w | t) dw > 0$  and  $p(w | t)$  is zero on this set.

Consequently, it holds

$$(8) \quad I(\hat{p}, p) < \infty \Rightarrow (\text{if } p(w | t) = 0 \text{ then } \hat{p}(w | t) = 0)$$

with the exception of a set of zero measure. (As usual, the abbreviation a.e. for almost everywhere will be used to denote that some feature holds up to such a set.)

Under general conditions [12]  $p(w | t)$  converges to the Dirac delta function concentrated on the “true”  $w$ . Thus  $\hat{p}(w | t)$  guaranteeing

$$\overline{\lim}_{t \rightarrow \infty} I(\hat{p}(w | t), p(w | t)) < +\infty$$

gives a strongly consistent estimate of  $w$  under the conditions referred.

#### Use of $\arg \inf(\cdot)$ symbol

Throughout the paper, the short-hand, appealing but imprecise notation

$$l = \arg \inf_{l \in L} J(l)$$

will be used ( $J$  is a functional on functions  $l$  from a set  $L$ ). The symbol  $\arg \inf(\cdot)$  will be interpreted as follows:

- the existence of a minimizing argument within  $L$  is generically supposed (if it is not the case, an  $\varepsilon$ -approximation,  $\varepsilon$  a small positive number,  $l_\varepsilon \in L$  for which  $J(l_\varepsilon) \leq \inf_L J(l) + \varepsilon$  is used instead),
- a specific representant is uniquely chosen if there are more minimizing functions (minimizing arguments are equivalent with respect to the task solved).

With the above agreement unnecessarily cumbersome notation is avoided.

### 3. ONE-SIDED APPROXIMATION

The section presents a one-shot approximation of the c.p.d.f.  $p(w | t)$ , the idea of which will be easily extended to the real-time-feasible approximation treated later.

The approximation would generate a  $\hat{p}(w | t) \in \hat{P}(t)$  which is close to  $p(w | t)$  in terms of the Kullback-Leibler distance  $I$ .

The need for some approximation arises whenever the “exact” c.p.d.f.  $p(w | t)$  is too complex for a particular treatment. Consequently, the approximation itself must not require the evaluation of the c.p.d.f.  $\hat{p}(w | t)$ . Instead of the distance  $I$ , an upper bound on it is searched for, which can be determined without full knowledge of  $p(w | t)$ . A solution in this vein is described by the following simple but fundamental lemma.

**Lemma 2.** Let  $\langle \hat{P}_p(t) \rangle$  be a nonempty subset of  $\langle \hat{P}(t) \rangle$  with members having

a likelihood fulfilling the inequality

$$(9) \quad \bar{l}(w | t) \leq p(w | t) \quad \text{a.e.}$$

Then

$$(10) \quad \inf_{\hat{P}(t)} I(\hat{p}, p) \leq \inf_{\langle \hat{P}_p(t) \rangle} (-\ln \int \bar{l}(w | t) dw).$$

Proof. Let  $\bar{l}(w | t)$  be any member of  $\langle \hat{P}_p(t) \rangle$ ; then it holds for any  $\hat{p}(w | t) \in \hat{P}(t)$

$$(11) \quad \hat{p}(w | t) \ln (\hat{p}(w | t) / p(w | t)) \leq \hat{p}(w | t) \ln (\hat{p}(w | t) / (\bar{l}(w | t) / \int \bar{l}(w | t) dw)) - \hat{p}(w | t) \ln (\int \bar{l}(w | t) dw) \quad \text{a.e.}$$

By integrating the inequality (11) we have

$$(12) \quad I(\hat{p}, p) \leq I(\hat{p}, \bar{l} / \int \bar{l}) - \ln (\int \bar{l}(w | t) dw).$$

The first term on the right-hand side of the inequality (12) achieves its minimum (zero) for

$$(13) \quad \hat{p}(w | t) = \bar{p}(w | t) = \bar{l}(w | t) / \int \bar{l}(w | t) dw \in \hat{P}_p(t) \subset \hat{P}(t).$$

The tightest bound (10) is obtained for

$$(14) \quad \bar{l}(w | t) = \operatorname{arginf}_{\langle \hat{P}_p(t) \rangle} (-\ln (\int \bar{l}(w | t) dw)). \quad \square$$

**Corollary 1.** If  $\hat{P}_p(t) = \hat{P}(t)$ , i.e. for any  $\hat{p}(w | t)$  there is a positive and finite  $k(t)$  such that  $k(t) \hat{p}(w | t) \leq p(w | t)$ , then

$$(15) \quad \inf_{\hat{P}(t)} I(\hat{p}, p) \leq -\ln (\sup \int \bar{l}(w | t) dw),$$

where sup is taken for  $\bar{l}(w | t) \leq p(w | t)$  a.e.,  $\bar{l}(w | t) \in \langle \hat{P}(t) \rangle$ .

Lemma 2 and Corollary 1 support the following construction of the approximant:  
Choose

$$\hat{p}(w | t) = \bar{l}(w | t) / \int \bar{l}(w | t) dw$$

where

$$(16) \quad \bar{l}(w | t) = \operatorname{arg sup} \int \bar{l}(w | t) dw,$$

where sup is taken for  $\bar{l}(w | t) \leq p(w | t)$  a.e.,  $\bar{l}(w | t) \in \langle \hat{P}(t) \rangle$ .

#### Remarks.

(i) The proposed way of constructing the upper bound on  $I(\hat{p}, p)$  using a point-wise bound on  $p(w | t)$  can be supported by some heuristics which takes into account relations between approximation theory and statistics [2]. The decisive reason for making this choice is, however, rather pragmatic one: this is the only estimate we have found up to now which has a global character and restricts the related multivariate integration (usually unfeasible) to  $\int \bar{l}(w | t) dw$ , which can be performed for carefully chosen class  $\hat{P}(t)$ , at least to the level needed for the optimization task (16).

(ii) The task (16) requires almost complete knowledge of  $p(w | t)$ , but it gives guideline how to construct a feasible (one-sided) approximation of  $p(w | t)$  in real time. The advantage of optimization in the class  $\langle \hat{P}_p(t) \rangle$  will be more transparent in the recursive case.

#### 4. ONE-SIDED RECURSIVE APPROXIMATION: APPROXIMATION OF THE BAYES RULE

A recursive variant of the foregoing approximation is presented here and its alternative formulation is given.

The starting point is a simple extension of Lemma 2.

**Lemma 3.** Let

$$\begin{aligned} \hat{l}(w | t - 1) \in \langle \hat{P}_{\langle p \rangle}(t - 1) \rangle &= \{ \hat{l}(w | t - 1) \in \langle \hat{P}(t - 1) \rangle : \hat{l}(w | t - 1) \leq \\ &\leq l(w | t - 1) \text{ a.e. for some } l(w | t - 1) \in \langle p(w | t - 1) \rangle \}. \end{aligned}$$

A. If the set

$$(18) \quad \langle \hat{P}_{\langle m \rangle}(t) \rangle = \{ \hat{l}(w | t) \in \hat{P}(t) : \hat{l}(w | t) \leq m(t | w) \hat{l}(w | t - 1) \text{ a.e.} \}$$

is nonempty then

$$(19) \quad \inf_{\hat{P}(t)} I(\hat{\beta}, p) \leq \varkappa(t) + \inf_{\hat{l}(w | t) \in \langle \hat{P}_{\langle m \rangle}(t) \rangle} (-\ln (\int \hat{l}(w | t) dw))$$

where  $\varkappa(t)$  is a constant independent of the approximant used.

B. The likelihood minimizing upper bound (19) can be generated recursively in the time course while  $\langle \hat{P}_{\langle m \rangle}(t) \rangle$  is nonempty, using just knowledge of  $p(w | 0) = p(w)$  and of  $\hat{l}(w | t - 1)$ ,  $m(t | w)$  at each time instant.

**Proof.**

Case A. By definition,  $\hat{l}(w | t - 1) \in \langle \hat{P}_{\langle p \rangle}(t - 1) \rangle$  means

$$(20) \quad \hat{l}(w | t - 1) \leq l(w | t - 1) = k p(w | t - 1)$$

for some positive finite  $k$ . Consequently, for any  $\hat{\beta}(w | t) \in \hat{P}(t)$  and any  $\hat{l}(w | t) \leq m(t | w) \hat{l}(w | t - 1)$ , using the Bayes rule (2), (3), we have

$$(21) \quad \begin{aligned} &\hat{\beta}(w | t) \ln (\hat{\beta}(w | t) / p(w | t)) = \\ &= \hat{\beta}(w | t) \ln ((\hat{\beta}(w | t) c(t) k) / (m(t | w) p(w | t - 1) k)) \leq \\ &\leq \hat{\beta}(w | t) \ln ((\hat{\beta}(w | t) c(t) k) / (\hat{l}(w | t) / \int \hat{l}(w | t) dw)) - \hat{\beta}(w | t) \ln (\int \hat{l}(w | t) dw) \text{ a.e.} \end{aligned}$$

Integrating (21) over the entire  $w$ -space and proceeding as in the proof of Lemma 2 we arrive at (19) with

$$(22) \quad \varkappa(t) = \ln (c(t) k).$$

Case B. It is sufficient to notice that nonnegativity of  $m(t | w)$  justifies the im-

plication:

$$(23) \quad \begin{aligned} \hat{l}(w | t) \leq m(t | w) \hat{l}(w | t - 1) &\Rightarrow \hat{l}(w | t) \leq m(t | w) \hat{l}(w | t - 1) = \\ &= p(w | t) c(t) k \end{aligned}$$

or equivalently

$$(24) \quad \langle \hat{P}_{\langle m \rangle}(t) \rangle \subset \langle \hat{P}_{\langle p \rangle}(t) \rangle.$$

The recursion has only to start with

$$(25) \quad \hat{l}(w | 0) \in \langle \hat{P}_{\langle p \rangle}(0) \rangle. \quad \square$$

The above lemma gives support for our constructive procedure which generates approximants from the class  $\hat{P}(t)$  in real time:

Choose

$$\hat{p}(w | t) = \hat{l}(w | t) / \int \hat{l}(w | t) dw$$

where

$$(26) \quad \hat{l}(w | t) = \arg \sup \int \hat{l}(w | t) dw,$$

where sup is taken for  $\hat{l}(w | t) \leq m(t | w) \hat{l}(w | t - 1)$  a.e.,  $\hat{l}(w | t) \in \langle \hat{P}(t) \rangle$ .

#### Remarks.

(iii) The structure of  $\hat{P}(t)$  guaranteeing  $\langle \hat{P}_{\langle m \rangle}(t) \rangle$  to be nonempty must be carefully chosen when selecting classes of approximants for  $t = 1, 2, 3, \dots$ . Clearly, the form of the system model  $m(t | w)$  has to be taken into account. The parts of the  $w$ -space where  $m(t | w)$  falls to (is) zero must be respected. This remark is a hint how to choose  $\hat{P}(t)$ .

(iv) The methodological consequence of Lemmas 1 and 3, which should be noticed, is the independence of the solution of any positive renormalization of the function  $m(t | w)$ . This fact makes the lack of knowledge of  $c(t)$  (3) unimportant, makes possible to choose an appropriate numerical level of evaluation and, above of all, gives an intuitive support for conjecture that the procedure (26) can converge although no forgetting is used, cf. [5].

The next lemma interprets alternatively the task (26) yielding an insight into the procedure adopted.

**Lemma 4.** The task (26) is equivalent to the "maximin" problem

$$(27) \quad \hat{p}(w | t) = \arg \sup_{\bar{p}(w|t) \in \hat{P}(t)} \operatorname{ess\,inf}_w \bar{p}(w | t) / \bar{p}(w | t)$$

where  $\bar{p}(w | t)$  denotes " $\hat{p}(w | t - 1)$  updated by the Bayes rule"

$$(28) \quad \bar{p}(w | t) = m(t | w) \hat{p}(w | t - 1) / \int m(t | w) \hat{p}(w | t - 1) dw.$$

Proof. Any  $\hat{l}(w | t) \in \langle \hat{P}(t) \rangle$  can be written in the form

$$(29) \quad \hat{l}(w | t) = \gamma(t) \bar{l}(w | t)$$

where  $\gamma(t)$  is a positive finite constant and  $\tilde{l}(w | t) \in \langle \hat{P}(t) \rangle$ . Consequently, the additional parameter  $\gamma(t)$  can be optimized when solving (26). An admissible pair  $\gamma(t), \tilde{l}(w | t)$  has to fulfil the inequality

$$(30) \quad \gamma(t) \tilde{l}(w | t) \leq m(t | w) \tilde{l}(w | t - 1) \quad \text{a.e.}$$

which is equivalent to

$$(31) \quad \gamma(t) \leq \operatorname{ess\,inf}_w m(t | w) \tilde{l}(w | t - 1) / \tilde{l}(w | t) = \Gamma(\tilde{l}).$$

The proof requires to demonstrate the equality of  $\gamma(t)$  and  $\Gamma(\tilde{l})$  for the optimal approximant. By contradiction, let for the optimal  $\tilde{l} = \gamma^* l^*$  be  $\gamma^*(t) < \Gamma(l^*)$ . Then the likelihood  $l^0(w | t) = \Gamma(l^*) l^*(w | t) \in \langle \hat{P}(t) \rangle$  and the inequality  $l^0(w | t) \leq m(t | w) \tilde{l}(w | t - 1)$  is fulfilled. For its integral the inequality

$$(32) \quad \int l^0(w | t) \, dw = \Gamma(l^*) \int l^*(w | t) \, dw = \\ = \Gamma(l^*) / \gamma^*(t) \int \tilde{l}(w | t) \, dw > \int \tilde{l}(w | t) \, dw$$

holds which contradicts optimality of the pair  $\gamma^*(t), l^*(w | t)$ .

Due to the optimality with respect to (26), the pair  $\gamma^*(t), l^*(w | t)$  fulfils, for any admissible  $\tilde{l}(w | t)$ , the inequality

$$\int \tilde{l}(w | t) \, dw \leq \gamma^*(t) \int l^*(w | t) \, dw = \Gamma(l^*) \int l^*(w | t) \, dw = \\ = \operatorname{ess\,inf}_w \bar{p}(w | t) / p^*(w | t) (\int m(t | w) \tilde{l}(w | t - 1) \, dw).$$

The above bound is reached for  $\tilde{l} = l^*$  and it is the highest one for the solution of the task (27).

Reversing the way of reasoning it can be found that the solution of the task (27) solves the task (26).  $\square$

#### Remarks.

(v) The task (27) benefits substantially from the structure of the space  $\hat{P}(t)$ ; the function  $\hat{p}(w | t)$  has to be p.d.f., i.e. nonnegative one integrating to unity. Just this structure implies that two functions  $\bar{p}, \hat{p}$  with common support have to "cross" each other.

(vi) The optimization task (26) is generally a difficult problem with functional restrictions. The class of practically feasible problems is, however, enlarged as test cases indicate.

### 5. APPLICATION TO LINEAR REGRESSION MODEL WITH CAUCHY NOISE

The section demonstrates how the proposed theory can be used.

At first, the treated system model is introduced and necessity to approximate



$p(w | t)$  shown. Then the class of approximating c.p.d.f.'s is chosen and the results of optimization (26) ((27)) are summarized.

### System Model

The linear regression model (for simplicity with single output)

$$(33) \quad y(t) = \Theta' z(t) + e(t)$$

with unknown coefficients  $\Theta$ , with the regressor vector  $z(t)$  being a known function of  $d(1 \dots t-1)$ ,  $u(t)$  (which can be recursively updated), and with the white noise term  $e(t)$ , is often used as a system model when designing adaptive predictors and/or controllers. Recursive least squares are used as a standard estimation procedure. The only but substantial drawback of this procedure is its sensitivity to outliers. A lot of modifications has been proposed recently, see [3]. The Bayesian approach to this problem is based on a proper modelling of the noise term and on the consistent application of the Bayes rule.

Recursive least squares are known to arise in Bayesian identification of the regression model (33) when  $e(t)$  is assumed to be normally distributed with constant (possibly unknown) dispersion  $\omega^{-1}$ . Normal distribution, however, assigns non-negligible probability within the interval of the width  $(2 \div 3) \omega^{-1/2}$ . If the presence of outliers is expected a more proper modelling is needed: the p.d.f. describing the noise has to fall to zero much slowly than the normal p.d.f. The Cauchy distribution

$$(34) \quad p(e(t) | \omega) = (2/\pi) \omega^{1/2} / (1 + \omega e^2(t))$$

is an extreme example to which we shall apply the developed theory hoping to achieve a robust, in the discussed sense, estimation procedure of the regression coefficients.

To demonstrate the need for some approximation and to make a preliminary step in approximation, we write explicitly the exact form of the likelihood  $l(w | t)$ . Defining

$$(35) \quad \begin{aligned} \varepsilon(t, w) &= \omega^{1/2}(y(t) - \Theta' z(t)) = \\ &= -[y(t), z'(t)] \begin{bmatrix} -\omega^{1/2} \\ \Theta \omega^{1/2} \end{bmatrix} = -h'(t) \begin{bmatrix} -\omega^{1/2} \\ \Theta \omega^{1/2} \end{bmatrix} \end{aligned}$$

where

$$(36) \quad w = (\Theta, \omega) \quad h'(t) = [y(t), z'(t)],$$

the system model can be written in the form (omitting the  $w$ -independent constant, cf. Remark (iv))

$$(37) \quad m(t | w) = \omega^{1/2} (1 + \varepsilon^2(t, w))^{-1}.$$

The likelihood functions for this model take the form

$$(38) \quad l(w | t) = \prod_{\tau=1}^t m(\tau | w) p(w | 0) = \omega^{t/2} \prod_{\tau=1}^t (1 + \varepsilon^2(\tau, w))^{-1} p(w | 0).$$

The exact evaluation of  $l(w | t)$  and consequently of the c.p.d.f.  $p(w | t)$  requires storing the full process history; the need for approximation arises.

### Class of Approximants

The assumed form of approximating likelihoods is

$$(39) \quad \tilde{l}(w | t) = \gamma(t) \omega^{v(t)/2} \left( 1 + \omega \begin{bmatrix} -1 \\ \Theta \end{bmatrix}' V(t) \begin{bmatrix} -1 \\ \Theta \end{bmatrix} \right)^{-v(t)}$$

which should be optimally "shaped" by a positive scalar factor  $v(t)$  and by a positive definite matrix  $V(t) > 0$ . The positive scalar  $v(t)$  has to evolve according to the relation

$$(40) \quad v(t) = v(t-1) + 1.$$

The reasons for this choice can be shortly summarized by discussing particular components in (39):

- (a)  $\gamma(t)$ : is a scaling factor the role of which can be understood in deep by inspecting the proof of Lemma 4.
- (b)  $\omega^{v(t)/2}$ :  $l(\cdot)$ ,  $\tilde{l}(\cdot)$  should be close each other, the part of  $l(\cdot)$ , which can be updated recursively, should form a part of  $\tilde{l}(\cdot)$ .

The use of the power  $v(t)$  instead of  $t$  just admits to incorporate some prior value, the recursion (40) corresponds simply to that of counting the number of samples ( $t$ ).

- (c)  $\left( 1 + \omega \begin{bmatrix} -1 \\ \Theta \end{bmatrix}' V(t) \begin{bmatrix} -1 \\ \Theta \end{bmatrix} \right)^{-v(t)}$ : see Lemma 3 and Remark (iii) stating that care

has to be taken to guarantee nonemptiness of  $\langle \hat{P}_{\langle m \rangle}(t) \rangle$  for all  $t$ . We have to be able to find approximate likelihoods fulfilling the inequalities

$$\tilde{l}(w | t) \leq m(t | w) \tilde{l}(w | t-1) \quad \text{for } t = 1, 2, 3, \dots$$

Since support of  $l(t | w)$  forms the entire "natural"  $w$ -space =  $\{\omega > 0, \Theta \text{ being real vector}\}$ , the behaviour of  $\tilde{l}(w | t)$ , when some entry of  $w$  approaches to infinity, i.e. when  $l(w | t)$  falls to zero, is decisive. Clearly,  $\tilde{l}(w | t)$  have to decay with a higher rate than  $l(w | t)$  for  $w$  tending to infinity. On the other hand tightness of the approximation requires the same rate of decrease in this area. Inspecting the form of  $l(w | t)$  (38) as a function of  $w$  (for  $w$  approaching infinity) we find that  $l$  behaves as the term discussed.

The following, theorem summarizes the results of the optimization for the regression model with a heavy-tailed noise.

**Theorem.** The c.p.d.f.  $p(w | t)$  of the parameter  $w = (\Theta, \omega)$ ,  $\dim(\Theta) = \varrho$ , of the "Cauchy-tailed" regression model

$$(41) \quad m(t | w) = \omega^{1/2} (1 + \omega(y(t) - \Theta' z(t))^2)^{-1}$$

can be optimally approximated (in the sense (27)) within the class (39) by the c.p.d.f.

$$(42) \quad \hat{p}(w | t) = \gamma(t) \omega^{v(t)/2} \left( 1 + \omega \begin{bmatrix} -1 \\ \Theta \end{bmatrix} V(t) \begin{bmatrix} -1 \\ \Theta \end{bmatrix} \right)^{-v(t)}$$

where the statistics of this distribution are recursively updated as follows

$$(43) \quad v(t) = v(t-1) + 1, \quad v(0) > 0 \text{ given,}$$

$$(44) \quad V(t) = (1 - \beta(t)) V(t-1) + \beta(t) \begin{bmatrix} y(t) \\ z(t) \end{bmatrix} \begin{bmatrix} y(t) \\ z(t) \end{bmatrix}^T, \quad V(0) > 0 \text{ given.}$$

The weight  $\beta(t)$  is a real root of the polynomial

$$(45) \quad \mu(\beta) = \sum_{i=0}^3 \mu_i \beta^i$$

with the coefficients (time-index omitted)

$$(46) \quad \begin{aligned} \mu_0 &= 2, \quad \mu_1 = -v + (v - \varrho - 1) \xi - (v - \varrho + 2) \xi_z \\ \mu_2 &= \xi \xi_z + (\varrho - 1) \xi - (2v - \varrho) \xi_z, \quad \mu_3 = -(v - 1) \xi \xi_z \end{aligned}$$

where

$$(47) \quad \xi(t) = 1 - \begin{bmatrix} y(t) \\ z(t) \end{bmatrix}^T V^{-1}(t-1) \begin{bmatrix} y(t) \\ z(t) \end{bmatrix}$$

$$(48) \quad \xi_z(t) = 1 - z(t)^T V_z^{-1}(t-1) z(t)$$

$$(49) \quad V(t) = \begin{bmatrix} V_y(t) & V_{zy}(t) \\ V_{zy}(t) & V_z(t) \end{bmatrix}^{-1} \varrho$$

For  $v \geq \varrho + 2$ , this root lies always in the interval

$$(50) \quad (\beta_0, 2/(v-1))$$

with

$$(51) \quad \beta_0 = 1/(v - (v-1)\xi^-)$$

$$(52) \quad \xi^- = \min(0, \xi).$$

Proof. The proof is outlined in the Appendix. It is based on a chain of technical lemmas which specify the solution of the "maximin" task (27).  $\square$

#### Remarks.

(vi) Using the decomposition (49) the following well-known "least-square" quantities are defined:

Point estimate of regression coefficients

$$(53) \quad \hat{\Theta} = V_z^{-1} V_{zy}.$$

Point estimate of the noise dispersion

$$(54) \quad \hat{\omega}^{-1} = V_y - V_{zy} V_z^{-1} V_{zy}.$$

(vii) The quantities  $\hat{\Theta}$ ,  $\hat{\omega}$  are well-known in connection with (recursive) least squares. The essence of the approach compresses, in the given case, into the choice of the weighting factor  $\beta(t)$ . Ordinary least squares (normed version) use

$$(55) \quad \beta_{LS}(t) = 1/v(t).$$

A simple analysis as well as numerical experience have shown that  $\beta(t)$  can be greater than (55) especially when starting the identification. The rate of convergence can be increased in this way. The claimed robustness is achieved because  $\beta(t)$  can be much smaller than (55) whenever the prediction error is too high.

(viii) The computational similarity with the recursive least squares makes possible to use, almost without changes, an effective LD-factorized numerical implementation [1].

## 6. CONCLUSIONS

The reported theory is by no means complete. The results gained up to now are, however, rather promising. The proposed way of approximating the Bayes identification

- is able to produce well-motivated approximative identification algorithms working in real time,
- has a global character taking into account both deterministic and stochastic ingredients of the model,
- can be, in a rather straightforward way, extended to the problems with time-dependent parameters, i.e. to nonlinear filtering.

The presented application example supports our belief that a proper direction of the research has been chosen. By using a unified theoretical and algorithmic approach we have arrived at reasonable, outlier robust estimation of regression coefficients. The resulting algorithm can be viewed as recursive weighted least squares with well-grounded choice of the weights.

## 7. APPENDIX: PROOF OF THEOREM

In order to solve the task (26) or its equivalent (27) the dependence of the integral  $\int \tilde{l}(w | t) dw$  on the parameters determining  $\tilde{l}(w | t)$  (i.e. on  $v(t)$ ,  $V(t)$ ) has to be found.

**Lemma 5.** For likelihoods in the class (39), it holds

$$(56) \quad \int \tilde{l}(w | t) dw = \gamma(t) |V_z|^{(n-\varrho)/2} |V|^{-(n-\varrho+1)/2} \tilde{\varphi}(t)$$

where  $\tilde{\varphi}(\cdot)$  is a term independent of  $v(t)$ ,  $V(t)$

$$(57) \quad V = \begin{bmatrix} \cdot & \cdot \\ \cdot & V_z \end{bmatrix} \begin{matrix} 1 \\ \varrho \end{matrix}, \quad \varrho = \dim(\Theta).$$

Proof. The proof is based essentially on the following substitution in the computed integral

$$(58) \quad x = V^{1/2}(t-1) \begin{bmatrix} -1 \\ \Theta \end{bmatrix} \omega^{1/2}$$

where  $V^{1/2}(t-1)$  denotes any square root of the positive definite matrix  $V(t-1)$ . Further details are omitted.  $\square$

The following lemma determines the form of the optimal  $\hat{l}(w | t)$  in the class (39), i.e. the optimal form of the matrix  $V(t)$ . Moreover, it makes the optimization problem more transparent by an appropriate transformation of the  $w$ -space. For this reason the subsequent notations will be used:

$$(59) \quad J = (V^{-1/2}(t-1))' V(t) V^{-1/2}(t-1)$$

$$(60) \quad f = (V^{-1/2}(t-1))' h(t)$$

$$(61) \quad \xi = 1 - f'f$$

$$(62) \quad \mathbf{1}' = [1, 1]$$

$$(63) \quad v = v(t).$$

**Lemma 6.** The optimal solution of the approximation task (27), say  $\hat{l}(w | t)$ , from the class (39) can be searched for in the subclass of (39) determined by

$$(64) \quad V(t) = (V^{1/2}(t-1))' J V^{1/2}(t-1),$$

where the positive definite matrix  $J$  takes the form ( $I$  denotes the unit matrix of an appropriate dimension)

$$(65) \quad J = (1 - \beta)I + (\beta - \alpha)ff' / f'f.$$

The matrices  $J$  are parametrized by the two-dimensional vector

$$(66) \quad v = \begin{bmatrix} \alpha \\ \beta \end{bmatrix}$$

lying in the *open* half-plane

$$(67) \quad v < \mathbf{1}.$$

The computation of the "min" part of the task (27)

$$(68) \quad \gamma(\hat{l}) = \operatorname{ess\,inf}_w \bar{p}(w | t) / \bar{\beta}(w | t)$$

reduces to the two-dimensional minimization (over some  $x' = [x_1, x_2]$ )

$$(69) \quad \gamma(\hat{l}) / \int \hat{l}(w | t) dw = \operatorname{ess\,inf}_{0 \leq x, x' \leq 1} (1 - v'x)^v / (1 - [\xi, 1]x).$$

Proof. The proof is based on a sequence of regular transformations of the  $w$ -space using the invariance of the task (27).

Starting with the substitution according to Eq. (58) the "inf" part of the optimiza-

tion task reads

$$(70) \quad \gamma(\bar{l}) = \operatorname{ess\,inf}_x (1 + x' J x)^{\nu} / ((1 + x' x)^{\nu-1} (1 + x' f f' x)) \int \bar{l}(w | t) dw$$

where unimportant (with respect to (27)) multiplicative factor  $\gamma(t-1)/\gamma(t)/\int m(t|w) \cdot \bar{l}(w|t-1) dw$  is omitted (cf. proof of Lemma 4) and the relations (35), (36), (39), (40), (44)–(48) are used

Before applying further transformations we shall prove that the optimal  $J$  can be searched for in the form (65). The use of Cholesky (lower triangular) square root  $V^{1/2}(t-1)$  in the definition (59) implies that (cf. (56))

$$(71) \quad \int \bar{l}(w | t) dw = \gamma(t) |J_z|^{(v-\varrho)/2} |J|^{-(v-\varrho+1)/2} \varphi(t, V(t-1))$$

where  $J_z$  corresponds to  $V_z$  in (42) and  $\varphi(\cdot, \cdot)$  to  $\tilde{\varphi}(\cdot)$  in (56). By inspecting the formulae (70), (71) it is can be seen that the value of  $\gamma(\bar{l})$  is *not influenced by any orthogonal transformation*

$$(72) \quad \mathcal{F}: J \rightarrow \mathcal{F} J \mathcal{F}'$$

*restricted to the orthogonal complement of the subspace spanned on the direction  $f$*  (including those describing permutation of entries). Consequently, the optimal  $J$  must have at most two different eigenvalues: one related to the span  $f$  and another one to its complement. A symmetric matrix of such a type takes the form (65) and the admissible area (67) expresses the required definiteness of  $J(V(t))$ .

Using the formulae proved it remains to map the real  $(\varrho + 1)$  dimensional space on the unit sphere by substituting

$$(73) \quad \bar{x} = x/\sqrt{(1 + x' x)}$$

and observing that

$$(74) \quad 1/(1 + x' x) = 1 - \bar{x}' \bar{x}.$$

The remaining substitution exemplifies splitting of the  $\bar{x}$ -space into the span of  $f$  and its complement. Because of invariance of the task with respect to the orthogonal transformation (72), it is sufficient to take  $\bar{x}$  as a linear combination of  $f$  and any orthogonal vector, say  $f^c$ , reducing the minimization just to the two-dimensional space of weights, say  $\pm\sqrt{\kappa_1}$ ,  $\pm\sqrt{\kappa_2}$ , i.e.

$$(75) \quad \bar{x} = \pm\sqrt{\kappa_1} f / \sqrt{f' f} \pm \sqrt{\kappa_2} f^c / \sqrt{f^c f^c}, \quad f' f^c = 0.$$

The use of relations (65), (66), (73)–(75) in the task (70) proves the formula (69).  $\square$

By computing  $\int \bar{l}(w | t) dw$  (with the help of (41)) for the class of approximants (39) with the kernel  $V(t)$  determined according to the formula (64) and with the matrix  $J$  given by (65) the following corollary can be verified.

**Corollary 2.** The approximation task (27) for the system model (37) is solved by  $\hat{p}(w | t)$  in the class (39) with  $V(t)$  (64), (65) when the parameter  $\hat{v}$  (66) is found as

$$(76) \quad \hat{v} = \arg \sup_{v < 1} \inf_{0 \leq \kappa, \kappa' \leq 1} H(v, \kappa)$$

where

$$(77) \quad H(v, \kappa) = (1 - v\kappa)^v (1 - [\xi, 1] \kappa)^{-1} (1 - \beta)^{-v/2} (1 - \alpha)^{-(v-\epsilon+1)/2} \cdot \left(1 - v' \begin{bmatrix} \delta \\ 1 - \delta \end{bmatrix}\right)^{(v-\epsilon)/2}$$

with

$$(78) \quad f' = \begin{bmatrix} \cdot & f'_z \\ \cdot & \cdot \end{bmatrix}$$

and

$$(79) \quad \delta = f'_z f_z / f' f \in (0, 1).$$

Proof. It only suffices to recall that the square matrix  $I - gg'$  where  $g$  is a nonzero  $n$ -vector has eigenvalues equal to 1 except one equal to  $1 - g'g$ .  $\square$

Now, we shall exclude some approximants as nonoptimal ones using necessary conditions of optimality. The following lemma states that the infimum in (69) is attained at the boundary of the admissible area in the *generic* case.

**Lemma 7.** If

$$(80) \quad \alpha - \beta\xi \neq 0$$

then the infimum in (69) is attained at the boundary of the region  $0 \leq \kappa, \kappa' \mathbf{1} \leq 1$ .

Proof. Differentiating the logarithm of the minimized function and setting the gradient equal to zero it can be found with some algebra that the stationary point  $\hat{\kappa}$  has to fulfil the pair of linear equations

$$(81) \quad A\hat{\kappa} = B; \quad A = \begin{bmatrix} -(v-1)\beta & \alpha - v\beta\xi \\ \xi\beta - v\alpha & -(v-1)\xi\alpha \end{bmatrix}, \quad B = \begin{bmatrix} 1 - \beta v \\ \xi - v\alpha \end{bmatrix}.$$

The square on the left-hand side of (80) is (positively) proportional to the determinant of  $A$ , so that under (80) we can find the unique stationary point by solving (81)

$$(82) \quad \hat{\kappa}' = [\alpha - \xi, 1 - \beta] / (\alpha - \beta\xi).$$

The restriction (67) and nonnegativity of admissible stationary points  $\hat{\kappa}$  imply

$$(83) \quad 1 > \beta \Rightarrow \alpha - \beta\xi > 0 \Rightarrow 1 > \alpha \geq \xi.$$

Taking into account the requirement  $\kappa' \mathbf{1} \leq 1$  it can be concluded that

$$(84) \quad 1 \geq \kappa' \mathbf{1} = (\alpha - \xi + 1 - \beta) / (\alpha - \beta\xi) \Rightarrow (1 - \beta)(1 - \xi) \leq 0$$

which is in contradiction to (83). Continuity of the minimized function and compactness of the admissible area of  $\kappa$  guarantee the rest.  $\square$

The following lemma restricts the domain of the optimal approximants to the "degeneracy" line

$$(85) \quad \alpha = \beta\xi.$$

**Lemma 8.** In the class (39) there is no optimal approximant of  $\beta(w | t)$  having the likelihood (38) while the parameter (66) is restricted by the condition (80).

*Proof.* Let the pair  $(\hat{v}, \hat{z})$  determine the optimal solution of the task (76) (equivalent to (27)). Because of the "sup"-part of (76), it must hold

$$(86) \quad H(v, z) \leq H(\hat{v}, \hat{z}) \quad \text{for any } v < 1.$$

The admissible range of  $v$  is an open set, so that the vanishing gradient of  $\ln(H(\hat{v}, \hat{z}))$  with respect to  $v$  gives the necessary optimality condition of the pair  $(\hat{v}, \hat{z})$ . This condition can be rearranged into the form

$$(87) \quad v\hat{z} = (1 - \hat{v}'\hat{z}) \left[ \begin{array}{l} -\frac{(v - \varrho)\delta}{2\left(1 - \hat{v}'\begin{bmatrix} \delta \\ 1 - \delta \end{bmatrix}\right)} + \frac{(v - \varrho + 1)}{2\left(1 - \hat{v}'\begin{bmatrix} 1 \\ 0 \end{bmatrix}\right)} \\ -\frac{(v - \varrho)(1 - \delta)}{2\left(1 - \hat{v}'\begin{bmatrix} \delta \\ 1 - \delta \end{bmatrix}\right)} + \frac{v}{2\left(1 - \hat{v}'\begin{bmatrix} 0 \\ 1 \end{bmatrix}\right)} \end{array} \right] = (1 - \hat{v}'\hat{z}) F.$$

Lemma 6 implies that the optimal  $\hat{z}'$  has one of the forms

$$(88) \quad (\alpha_1, 0), \quad (0, \alpha_2), \quad (\alpha_1, 1 - \alpha_1); \quad 0 \leq \alpha_1, \alpha_2 \leq 1.$$

By checking the possibilities (73), we shall find that no such a pair exists. For instance, let  $\hat{z}' = (0, \alpha_2)$ . For the optimal  $\hat{v}, \hat{z}$ , it implies (see (87))

$$(89) \quad F_1 = 0.$$

This condition can be rearranged as

$$(90) \quad 1' \left[ \begin{array}{l} \delta \\ (v - \varrho + 1)(1 - \delta) \end{array} \right] = \hat{v}' \left[ \begin{array}{l} \delta \\ (v - \varrho + 1)(1 - \delta) \end{array} \right]$$

which is contradictory for  $\hat{v} < 1$  and  $\delta \in (0, 1)$  (cf. (79)). Similar reasoning can be performed for the remaining cases and it is therefore omitted.  $\square$

By combining the results of Lemmas 6, 7, 8 and using formulae (61), (65), (67), (69) and (70) we have:

**Corollary 3.** The matrix  $V(t)$  determining the optimal approximant is given by the formula (64) with

$$(91) \quad J = (1 - \hat{\beta})I + \hat{\beta}ff'$$

and

$$(92) \quad \hat{\beta} = \arg \sup_{\beta < 1} \inf_{\xi \leq \psi \leq 1} \tilde{H}(\beta, \psi)$$

where

$$(93) \quad \tilde{H}(\beta, \psi) = (1 - \beta\psi)^v (1 - \psi)^{-1} (1 - \beta)^{-v/2} (1 - \beta\xi)^{-(v-\varrho+1)/2} \cdot (1 - \beta\xi_z)^{-(v-\varrho)/2}$$



with (cf. (61), (78))

$$(94) \quad \xi = 1 - f'f, \quad \xi_z = 1 - f'_z f_z,$$

$$(95) \quad \xi^- = \min(0, \xi).$$

*Proof.* The proof employs the mentioned formulae and the definition of the new independent, one-dimensional variable

$$(96) \quad \psi = (\xi, 1) \times \quad \text{for which} \quad \xi^- \leq \psi \leq 1$$

instead of  $\times$ .  $\square$

The lemma below summarizes results of the analytical minimization of  $\tilde{H}(\beta, \psi)$  (93) with respect to  $\psi$ .

**Lemma 9.** It holds

$$(97) \quad \inf_{\xi^- \leq \psi \leq 1} \tilde{H}(\beta, \psi) = \begin{cases} \tilde{H}(\beta, \xi^-) & \text{for } \beta \leq \beta_0 \\ v^{\varrho}/(v-1)^{(v-1)} \beta(1-\beta)^{(v-2)/2} (1-\beta\xi)^{-(v-\varrho+1)/2} & \text{for } \beta \geq \beta_0 \\ (1-\beta\xi_z)^{(v-\varrho)/2} & \text{for } \beta \geq \beta_0 \end{cases}$$

where

$$(98) \quad \beta_0 = 1/(v - (v-1)\xi^-).$$

*Proof.* This simple exercise in optimization is omitted.  $\square$

The next lemma helps to avoid the branch  $\beta \leq \beta_0$  of (97) when computing the "sup"-part of (93) to determine  $\hat{\beta}$ .

**Lemma 10.** The optimal weight  $\hat{\beta}$  in (91) is given by

$$(99) \quad \hat{\beta} = \arg \sup_{\beta_0 \leq \beta < 1} \beta(1-\beta)^{(v-2)/2} (1-\beta\xi)^{-v-\varrho+1)/2} (1-\beta\xi_z)^{(v-\varrho)/2} = \\ = \arg \sup_{\beta_0 \leq \beta < 1} H_0(\beta).$$

*Proof.* The function (97) is a continuous function of  $\beta$  even at the point  $\beta_0$ . Thus, if we prove  $\tilde{H}(\beta, \xi^-)$  to be an increasing function of  $\beta$  for  $\beta \leq \beta_0$ , we have

$$(100) \quad \hat{\beta} = \arg \sup_{\beta} \inf_{\psi} \tilde{H}(\beta, \psi) \geq \beta_0.$$

To this end we shall prove the positivity of the derivative of  $\ln(\tilde{H}(\beta, \xi^-))$  with respect to  $\beta$  (for  $\beta \leq \beta_0$ ), i.e.

$$(101) \quad \bar{H}(\beta) = 2 \partial \tilde{H}(\beta, \xi^-) / \partial \beta / \tilde{H}(\beta, \xi^-) = \\ = -2v\xi^- / (1 - \beta\xi^-) + v/(1 - \beta) + (v - \varrho + 1) \xi / (1 - \beta\xi) - (v - \varrho) \xi_z : \\ : (1 - \beta\xi_z).$$

The last term in (101) is a decreasing function of  $\xi_z \leq 1$ , thus on replacing  $\xi_z$  by 1 we find

$$(102) \quad \bar{H}(\beta) \geq -2v\xi^- / (1 - \beta\xi^-) + \varrho / (1 - \beta) + (v - \varrho + 1) \xi / (1 - \beta\xi).$$

The next estimate, proving the lemma, is obtained by observing that  $\xi^- \leq \xi$  and noting that the last term in (102) is an increasing function of  $\xi$ , consequently,

$$(103) \quad \bar{H}(\beta) \geq -2\xi^-(v + \varrho - 1)/(1 - \beta\xi^-) + \varrho/(1 - \beta) > 0$$

because  $\xi^- \leq 0$ . □

The concluding "optimization" lemma states the necessary optimality conditions and specifies the interval in which the optimizer can be searched for.

**Lemma 11.** The maximizer of the function in (99)

$$(104) \quad H_0(\beta) = \beta(1 - \beta)^{(v-2)/2} (1 - \beta\xi)^{(v-\varrho+1)/2} (1 - \beta\xi_z)^{(v-\varrho)/2}$$

can be found as a real root of the polynomial

$$(105) \quad \mu(\beta) = \sum_{i=0}^3 \mu_i \beta^i$$

with the coefficients

$$(106) \quad \begin{aligned} \mu_0 &= 2, \quad \mu_1 = -v + (v - \varrho - 1)\xi - (v - \varrho + 2)\xi_z, \\ \mu_2 &= \xi\xi_z + (\varrho - 1)\xi - (2v - \varrho)\xi_z, \quad \mu_3 = -(v - 1)\xi\xi_z. \end{aligned}$$

This root has to be in the interval

$$(107) \quad \langle \beta_0, 2/(v - 1) \rangle.$$

*Proof.* It holds

$$(108) \quad \begin{aligned} \bar{H}_0(\beta) &= 2 \partial H_0(\beta) / \partial \beta / H_0(\beta) = \\ &= 2/\beta - (v - 2)/(1 - \beta) + (v - \varrho + 1)\xi/(1 - \beta\xi) - (v - \varrho)\xi_z/(1 - \beta\xi_z). \end{aligned}$$

The numerator of (108) can be verified to have the form (105), with the coefficients (106). The last term in (108) is an increasing function of  $\xi_z \leq \xi$ . Replacing  $\xi_z$  by  $\xi$  we get the estimate

$$(109) \quad \bar{H}_0(\beta) \leq 2/\beta - (v - 2)/(1 - \beta) + \xi/(1 - \beta\xi).$$

The last term in (109) is an increasing function of  $\xi \leq 1$  which gives the following upper bound

$$(110) \quad \bar{H}_0(\beta) \leq 2/\beta - (v - 3)/(1 - \beta) = (2 - (v - 1)\beta)/(\beta(1 - \beta)) < 0$$

for  $\beta > 2/(v - 1)$ , which determines the interval (107). □

**Remark.**

(ix) The upper bound of the interval (107) proved to be a reasonable starting point for an iterative gradient-type search of the optimal  $\hat{\beta}$  because of the guaranteed sign of the gradient.

(Received February 2, 1988.)

## REFERENCES

- [1] G. J. Biermann: Factorization Methods for Discrete Sequential Estimation. Academic Press, New York 1977.
- [2] B. Harris and G. Heinell: The relation between statistical decision theory and approximation theory. In: Optimizing Method in Statistics — Proceedings of an International Conference (J. S. Rustagi, ed.). Academic Press, New York—San Francisco—London 1979.
- [3] P. J. Huber: Robust Statistics. J. Wiley, New York 1981.
- [4] M. Kárný: Algorithms for determining the model structure of a controlled system. *Kybernetika* 19 (1983), 164—178.
- [5] M. Kárný nad K. M. Hangos: Approximation of Bayes rule. Preprints of the 7th IFAC/IFORS Symposium on Identification and System parameter Estimation, Vol. 1, pp. 1755—1760, York 1985.
- [6] M. Kárný and K. M. Hangos: One-sided approximation of Bayes rule: theoretical background. Preprints of the 10th IFAC World Congress, Vol. 10, pp. 312—317, Munich 1987.
- [7] M. Kárný and R. Kulhavy: 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.
- [8] R. Kulhavy: Restricted exponential forgetting in real-time identification. *Automatica* 23 (1987), 589—600.
- [9] S. Kullback: Information Theory and Statistics. J. Wiley, New York 1959.
- [10] A. Perez: Information,  $\varepsilon$ -sufficiency and data reduction problems. *Kybernetika* 1 (1965), 297—323.
- [11] V. Peterka: Bayesian approach to system identification. In: Trends and Progress in System Identification (P. Eykhoff, ed.), Chap. 8, pp. 239—304. (IFAC Series for Graduates, Research Workers and Practicing Engineers 1.) Pergamon Press, Oxford 1981.
- [12] Y. Yashin: Martingal approach to identification of stochastic systems. Preprints of the 7th IFAC/IFORS Symposium on Identification and System parameter Estimation, Vol. 2, pp. 1755—1760, York 1985.

*Ing. Miroslav Kárný, CSc., Ústav teorie informace a automatizace ČSAV (Institute of Information Theory and Automation — Czechoslovak Academy of Science), Pod vodárenskou věží 4, 182 08 Praha 8, Czechoslovakia.*  
*Dr. Katalin M. Hangos, Computer and Automation Institute — Hungarian Academy of Sciences, Kende utca 13—17, H 1502 Budapest, Hungary.*