

NOT FOR QUOTATION
WITHOUT PERMISSION
OF THE AUTHOR

**SOME ASPECTS OF MODEL TUNING PROCEDURE:
INFORMATION-THEORETIC ANALYSIS**

A.I.Yashin

April 1985
WP-85-24

Working Papers are interim reports on work of the International Institute for Applied Systems Analysis and have received only limited review. Views or opinions expressed herein do not necessarily represent those of the Institute or of its National Member Organizations.

INTERNATIONAL INSTITUTE FOR APPLIED SYSTEMS ANALYSIS
2361 Laxenburg, Austria

ACKNOWLEDGMENTS

I would like to thank Susanne Stock for her help in preparing the manuscript.

SOME ASPECTS OF MODEL TUNING PROCEDURE: INFORMATION-THEORETIC ANALYSIS

A.I.Yashin

1. INTRODUCTION

Computer or mathematical models are not exact representation of reality: lack of knowledge, technical restrictions and particular modeling goals make it necessary to approximate the real system in various ways. Nevertheless, the procedures by which the models are adjusted to observed data are often based on the assumptions that the real system has the same structure as the model and differs only in the values of certain parameters. These particular values usually should be included in the feasible set of the parameter values, and this fact, together with some additional conditions, usually provides the convergence property for many individual algorithms [1].

However, in reality all of these assumptions are generally false. Even if the structure of the system corresponds to the structure of the model, the real parameters values often do not belong to the presupposed feasible set. Moreover, mathematicians often consciously diminish this set in order to simplify the estimation algorithms. For instance they approximate the bounded compact set of parameter values by a set consisting of a finite number of points, thus increasing the chances that the real parameter values will be excluded.

It is therefore both remarkable and surprising to find that despite these false assumptions and approximations, the parameter estimation algorithms often still converge! The model resulting from this tuning procedure will of course not coincide with the real system, and this rises the natural question: how far is this computer model from reality?

When considering this question it is necessary to have some way of measuring the distance between individual models. One of measure of divergence was introduced by Bhattacharya [2]; Kullback [3] also formulated some measure of information distance. However these measures were not proper metrics. Baram and Sandell [4] later introduced a modified version of Kullback measure, which have been shown to be a proper distance metric. They applied this approach to linear Gaussian systems and models; in this paper it is generalized to a wider class of systems.

2. NOTATIONS AND DEFINITION

Assume that the variety of models of the real system may be characterized by a parameter β , which takes values from the parameter set B . In view of Bayesian formulation of the problem, we will assume β to be a random variable defined on some probabilistic space (Ω, H, P) . Let $\xi_n(\omega), n \geq 0$ be some random process (observation) adapted to some nondecreasing family of σ -algebras $\mathbf{H} = (H_n)_{n \geq 0}$, $H_\infty = H$ in Ω . We shall denote by $\bar{\mathbf{H}} = (\bar{H}_n)_{n \geq 0}$, $\bar{H}_\infty = \bar{H}$ the family of σ -algebras generated by the process $\xi_n, n \geq 0$, where

$$\bar{H}_n = \sigma\{\xi_m, m \leq n\}$$

is σ -algebra in Ω generated by the process ξ_t up to time t .

In the case of continuous time observation process $\xi_t, t \geq 0$ we assume the non-decreasing right continuous family of σ -algebras $\mathbf{H} = (H_t)_{t \geq 0}$ to be given, where $H_\infty = H$ and H_0 is completed by P -zero sets from H . We also introduce the family of σ -algebras $\bar{\mathbf{H}} = (\bar{H}_t)_{t \geq 0}$, where

$$\bar{H}_t = \sigma\{\xi_u, u \leq t\}$$

If the set of the parameter values is finite or denumerable we will denote by $\pi_j(n)$, (or $\pi_j(t)$) the *a posteriori* probabilities of events $\{\beta = \beta_j\}$, $j \in B$ given observations ξ_k , $k \leq n$, (ξ_u , $u \leq t$).

For any $A \in \bar{H}$, $x \in B$ we denote by $\mathbf{P}^x(A)$, the family of probability measures

$$\mathbf{P}^x(A) = P(A | \beta = x).$$

Let $\bar{\mathbf{P}}_n^x(A)$, $\bar{\mathbf{P}}^x(A)$, $x \in B$, $n \geq 0$ be the restrictions of the $\mathbf{P}^x(A)$ on σ -algebras \bar{H} , \bar{H}_n , respectively. Assume also that for any $x, y \in B$ we have $\bar{\mathbf{P}}_n^x \sim \bar{\mathbf{P}}_n^y$. Define $Z_n^{x,y}$ as a Radon-Nicodim derivative

$$Z_n^{x,y} = \frac{d\bar{\mathbf{P}}_n^x}{d\bar{\mathbf{P}}_n^y}$$

and let

$$\alpha_n^{x,y} = Z_n^{x,y} (Z_n^{y,x})^{-1}.$$

It is easy to see that if the \bar{H}_{n-1} -conditional distributions of ξ_n , $n \geq 0$ have densities $f^x(z | \bar{H}_{n-1})$, $x \in B$ then

$$\alpha_n^{x,y} = \frac{f^x(\xi_n | \bar{H}_{n-1})}{f^y(\xi_n | \bar{H}_{n-1})}$$

3. SOME BAYESIAN PARAMETER ESTIMATION ALGORITHM

Before deriving our main results, we will first consider some Bayesian parameter estimation algorithms for different observation schemes.

a) Assume that $\xi_n, n \geq 0$ is given by the formula

$$\xi_n = \vartheta_n + C_t \varepsilon_{2n}.$$

where ϑ_n satisfies the recursive stochastic equation

$$\vartheta_{n+1} = \beta \vartheta_n + b_t \varepsilon_{1n+1}$$

Here $\varepsilon_{1n}, \varepsilon_{2n}, n \geq 0$ are the sequences of independent Gaussian random variables with zero mean and variance equal to one, and β is an unknown parameter. Assuming that β takes its values from some finite set $B_k = \{\beta_1, \beta_2, \dots, \beta_k\}$ the *a posteriori* probabilities are

$$\pi_j(n+1) = \pi_j(n) \left[\frac{\frac{1}{D_j(n+1)} e^{-\frac{(\xi_{n+1} - \beta_j m_n^j)^2}{2D_j^2(n)}}}{\sum_{i=1}^k \frac{\pi_i(n)}{D_i(n+1)} e^{-\frac{(\xi_{n+1} - \beta_i m_n^i)^2}{2D_i^2(n)}}} - 1 \right]$$

where m_n^j are Kalman estimates of ϑ_n given $\{\beta = \beta_j\}$ and $D_j(n)$ are functions of the conditional variance $\gamma_j(n)$ [5]

$$D_j(n) = (b^2 + C^2 + \beta_j^2 \gamma_j(n))^{\frac{1}{2}}$$

b) Consider the continuous (in time) observation process ξ_t given by the stochastic differential equation

$$d\xi_t = A(\beta, \xi_t)dt + C_t dW_t$$

where $W_t, t \geq 0$ is the H -adapted Wiener process, β is an unknown parameter and C_t is H -adapted positive function. Assuming again that the number of parameter values is finite, we have for $\pi_j(t) = P(\beta = \beta_j | \bar{H}_t)$ [6].

$$d\pi_j(t) = \pi_j(t) \frac{(A(\beta_j, \xi_t) - \bar{A}(\xi_0^t))}{C_t^2} (d\xi_t - \bar{A}(\xi_0^t)dt), \quad \pi_j(0) = p_j, \quad j = 1, 2, \dots, k,$$

where

$$\bar{A}(\xi_0^t) = \sum_{i=1}^k \pi_i(t) A(i, \xi_t).$$

c) Consider an observation made by a continuous-state jumping process with unknown transition intensities $\lambda_{i,j}^\beta$. Once again assuming a finite number of values for β we have the following equations for *a posteriori* probability $\pi_j(t)$ [7]

$$\pi_j(t) = \pi_j(0) + \sum_{u \leq t} \pi_j(u) \left(\frac{\lambda_{\xi_u - \xi_u}^j}{\bar{\lambda}_{\xi_u - \xi_u}} - 1 \right) - \int_0^t \pi_j(u) (\lambda_{\xi_u - \xi_u}^j - \bar{\lambda}_{\xi_u - \xi_u}) du$$

where

$$\bar{\lambda}_{\xi_u - \xi_u} = \sum_{i=1}^k \pi_i(u) \lambda_{\xi_u - \xi_u}^i$$

The necessary and sufficient conditions of convergence with probability one for *a posteriori* probabilities to respective indicators were given in the papers [1, 8, 9] in terms of absolute continuity and singularity of some special families of probability distributions. Papers demonstrated the applications of the general theory to various particular forms of the random processes.

4. AUXILIARY RESULTS

One of the central places in the proof of the main convergence result in [5, 9, 1] was the relation between *a posteriori* probabilities and likelihood ratio in the case of denumerable or finite number of the parameter values. More exactly the following lemma is true:

Lemma 1. *Let for any $i = j$ and $n \geq 0$ measure \bar{P}_n^j is equivalent to the measure \bar{P}_n^i . Then \bar{P}^j -a. s. the next equality is true:*

$$\frac{\pi_i(n)}{\pi_j(n)} = \frac{p_i}{p_j} Z_n^{i,j} \quad (1)$$

The proof of this lemma follows from the definition of the likelihood ratio $Z_n^{i,j}$. The equality (1) yields that

$$\{Z_\infty^{i,j} = 0\} = \{\pi_j(\infty) = 1\} = \{\beta = \beta_j\}$$

According to the papers [1, 8, 9] this property guarantees the following result of convergence: (remind that we still deal with the case when the parameter value corresponding to the real system belongs to the feasible set of the parameter values B).

Theorem 1. *Let for any $i = j$, $n \geq 0$, $\bar{P}_n^i \sim \bar{P}_n^j$. Then the condition $\bar{P}^i \perp \bar{P}^j$ is equivalent to the condition*

$$\lim_{n \rightarrow \infty} \pi_j(n) = I(\beta = \beta_j), \text{ P-a.s.}$$

The proof of this theorem is based on the property that singularity set for the measures \bar{P}^i and \bar{P}^j coincide \mathbf{P}^j -a. s. with the set $\{\beta = \beta_j\}$.

If the real parameter value β_k does not belong to the feasible set variables $\pi_i(n), i \in B$ calculated in section 3 are already not the *a posteriori probabilities*, but some functionals of the observable process ξ_n .

Taking them as *a posteriori* probabilities, the observer expects to get the convergence one of $\pi_i(n), i \in B$ (say $\pi_{i_0}(n)$) to 1 and interpret this result as if the real parameter value is equal to i_0 . However this is actually a false conclusion. The questions which arise in this relation are: When does the convergency fact for some of the $\pi_i(n), i \in B$ really take place? What does it mean when $\pi_{i_0}(n)$ tends to 1 for some $i_0 \in B$? In order to answer these questions we need some auxiliary results.

Assume that the real system corresponds to a parameter value k such that $k \in B$. Introduce the function $I_n^k(x, y) = E_k \ln \alpha_n^{x, y}$ [4] and define the measure of distance

$$d_n(x, y) = |I_n^k(x, y)|.$$

Lemma 2. *Function $d_n(i, j)$ is pseudo-metric. That is, the following equalities hold:*

$$d_n(x, y) = d_n(y, x)$$

$$d_n(x, x) = 0$$

$$d_n(x, k) + d_n(k, y) \geq d_n(x, y)$$

The proof of this lemma is done in [4].

Lemma 3. For any $x, y \in B$, $n \geq 0$ we have

$$I_n^z(x, y) \geq 0.$$

Proof. From the definition of the $I_n^z(x, y)$

$$I_n^z(x, y) = E_x(\ln \alpha_n^{z, y} | \bar{H}_{n-1}) = E_x(E_x(\ln \alpha_n^{z, y} | \bar{H}_{n-1})) = E_x(E_y(\psi(\alpha_n^{z, y}) | \bar{H}_{n-1}))$$

where $\psi(t) = t \ln t$. According to the theorem of the mean, $\psi(\alpha_n^{z, y})$ can be represented as follows:

$$\psi(\alpha_n^{z, y}) = \psi(1) + (\alpha_n^{z, y} - 1)\psi'(1) + \frac{1}{2}(\alpha_n^{z, y} - 1)^2\psi''(\delta_n^{z, y}),$$

where $\delta_n^{z, y}$ varies between $\alpha_n^{z, y}$ and 1. It is not difficult to see that

$$E_y(\psi(\alpha_n^{z, y}) | \bar{H}_{n-1}) = \frac{1}{2}E_y\left(\frac{(\alpha_n^{z, y} - 1)^2}{\delta_n^{z, y}} | \bar{H}_{n-1}\right) \geq 0$$

Lemma 4. Let $d_n(k, x) \leq d_n(k, y)$. Then

$$I_n^k(x, y) \geq 0$$

Proof. From the definition of the $I_n^k(x, y)$, we can write

$$\begin{aligned} I_n^k(x, y) &= E_k \ln \alpha_n^{z, y} = E_k \ln f^z(\xi_n | \bar{H}_{n-1}) - E_k \ln f^y(\xi_n | \bar{H}_{n-1}) \\ &= -E_k[\ln f^k(\xi_n | \bar{H}_{n-1}) - \ln f^z(\xi_n | \bar{H}_{n-1})] + E_k[\ln f^k(\xi_n | \bar{H}_{n-1}) - \ln f^y(\xi_n | \bar{H}_{n-1})] \\ &= -E_k \ln \alpha_n^{k, z} + E_k \ln \alpha_n^{k, y} \end{aligned}$$

From Lemma 3 for any $z \in B$

$$E_k \ln \alpha_n^{k, z} \geq 0$$

and thus

$$I_n^k(x,y) \geq 0 \text{ if } d_n(k,y) \geq d_n(k,x)$$

5. RESULTS

Assume that the process $\ln \alpha_n^{x,y}$ is ergodic, i.e.,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{m=1}^n \ln \alpha_m^{x,y} = E_k \ln \alpha^{x,y} = I^k(x,y)$$

Theorem 2. *If $d(k,x) > d(k,y)$ then*

$$Z_n^{x,y} \rightarrow 0, \text{ P-a.s.}$$

If it is known that $Z_n^{x,y} \rightarrow 0$ P-a.s., then

$$d(x,y) \geq d(k,y)$$

Proof. Note that from Lemma 4, the inequality $d(k,x) > d(k,y)$ yields $I^k(x,y) < 0$ and consequently

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{m=1}^n \ln \alpha_m^{x,y} < 0 \text{ P-a.s.}$$

This means that

$$\ln Z_n^{x,y} = \sum_{m=1}^n \ln \alpha_m^{x,y} = -\infty \text{ P}^k\text{-a.s.}$$

and consequently

$$Z_n^{x,y} \rightarrow 0 \text{ P}^k\text{-a.s.}$$

thus proving the first part of the theorem.

In order to prove the second part of the theorem we assume that $Z_n^{x,y} \rightarrow 0$ but that $d(k,x) < d(k,y)$. This yields

$$\lim_{n \rightarrow \infty} \sum_{m=1}^n \ln \alpha_m^{x,y} = I^k(x,y) > 0,$$

from which

$$\ln Z_{\infty}^{z,y} = \sum_{m=1}^{\infty} \ln \alpha_m^{z,y} = \infty$$

and the theorem is proved by contradiction.

Example. Assume that the sequence ξ_n is a finite state ergodic Markov chain on any of the probability spaces $(\Omega, \bar{H}, \bar{P}^i)$, $i \in B$, where B is a finite set. Let $p_{l,m}^i$, $l, m = \overline{1, k}$ be the transition probabilities for one step. It is not difficult to find (see also [8]) that $\alpha_n^{i,j}$ is given by the formula

$$\alpha_n^{i,j} = \frac{P_{\xi_{n-1}, \xi_n}^i}{P_{\xi_{n-1}, \xi_n}^j}$$

Well known results from the Markov chain theory (see [10] for instance) show that the process $\ln \alpha_n^{i,j}$ is ergodic. Thus if the Bayesian algorithm for $\pi_j(n)$ converges to 1 for some particular j_0 it means that this j_0 is the point from B that is the nearest (in the sense of information distance $d(k, x)$) to the real parameter value k .

REFERENCES

1. A.I. Yashin, *Bayesian Approach To Parameter Estimation: Convergence Analysis*, WP-83-67, International Institute For Applied Systems Analysis, Laxenburg, Austria (July 1983).
2. A. Bhattacharya, "On Measure Of Divergence Between Two Statistical Populations Defined By Probability Distributions," *Bulletin Calcutta Mathematical Society* 35, pp.99-104 (1943).
3. S. Kullback, *Information Theory And Statistics*, Wiley, New York (1959).
4. Y. Baram and N.R. Sandell, "An Information Theoretic Approach To Dynamical System Modeling And Identification," *IEEE Transactions Automatic Control* AC-23(1), pp.61-66 (1978).

5. N.M. Kuznetsov, A.V. Lubkov, and A.I. Yashin, "About Consistency Of Bayesian Estimates In Adaptive Kalman Filtration Scheme ," *Automatic and Remote Control (translated from Russian)*(4), pp.47-56 (1981).
6. R.S. Liptzer and A.N. Shiryaev, *Statistics of Random Processes*, Springer-Verlag, Berlin and New York (1978).
7. A.I. Yashin, "Filtering of Jumping Processes," *Automatic and Remote Control* 5, pp.52-58 (1970).
8. A.I. Yashin, "Sostoyatel'nost' Bayesovskikh Otsenok Parametrov (Consistency of Bayesian Parameter Estimates)," *Problemi Peredachi Informacii (in Russian)* (1), pp.62-72 (1981).
9. N.M. Kuznetsov and A.I. Yashin, "On the Conditions of the Identifiability of Partially Observed Systems," *Doklady Akademii Nauk SSSR (in Russian)* 259(4), pp.790-793 (1981).
10. S. Karlin, *A First Course In Stochastic Processes*, Academic Press, New York and London (1968).