

ON DUAL EXPRESSION OF PRIOR INFORMATION IN BAYESIAN PARAMETER ESTIMATION*

R. Kulhavý^{1,2} and L. Tesář²

¹ Honeywell Technology Center—Prague, 182 08 Prague, Czech Republic

² Institute of Information Theory and Automation, Acad. Sci. CR, 182 08 Prague, Czech Rep.

Abstract: In Bayesian parameter estimation, *a priori* information can be used to shape the prior density of unknown parameters of the model. When chosen in a conjugate, self-reproducing form, the prior density of *parameters* is nothing but a model-based transform of a certain “prior” density of observed *data*. This observation suggests two possible ways of expressing *a priori* knowledge—in terms of parameters of a particular model and in terms of data entering the model. The latter way turns out useful when dealing with statistical models whose parameters lack a direct physical interpretation. In practice, the amount of *a priori* information is usually not sufficient for complete specification of the prior density of data. The paper shows an information-based way of converting such incomplete information into the prior density of unknown parameters.

Keywords: Parameter estimation, Bayesian methods, prior information, information measures, approximation.

1. INTRODUCTION

The role of *a priori* information in parameter estimation has become widely acknowledged in the control community. It is well known that information about the system dynamics may significantly affect parameter estimates. Even rudimentary knowledge of the identified system can effectively regularize parameter estimation. This is crucial namely in poorly conditioned cases when data carry little information.

The use of prior information depends on the estimation methodology used. The Bayesian paradigm provides a framework which is particularly well suited for merging of prior and data-based information. *A priori* information is used here to shape the prior density of unknown parameters while Bayes’s rule serves as a tool for updating prior knowledge by information contained in observed data.

Indeed, when speaking about *a priori* information, one usually has in mind the choice of the prior density, in other words, the choice of the “initial conditions” for the parameter estimator (Peterka, 1981). More generally, when parameters vary in time, prior information

may differ at every time instant. When such information is available, it can be used to shape reference densities in the generalized forgetting schemes (Kulhavý and Kraus, 1996). Finally, there are situations when no data are available and decision is made solely on the basis of prior information. A careful specification of the prior density is of vital importance then.

Provided the model is structured so that its parameters admit a direct physical interpretation, the choice of a particular prior density is a rather intuitive way of expressing available knowledge. The examples of such parameters include rates and orders of chemical reactions, growth rates of biomass or characteristics of valves.

The notorious difficulty with physically parametrized models is their nonlinearity which makes estimation of parameters a non-trivial task. It is why in on-line identification one meets statistically-rooted models more often than first-principle models. A typical example of statistical models is the autoregressive model with external input (ARX). While the parameters of the ARX model are relatively easy to estimate, they usually lack a simple physical interpretation. The regression coefficients describe nothing but a linearized dependence of the system output on the regressor entries (which may be nonlinear functions of past data). The variance

*The work was supported in part by the grants 102/97/0466 and 402/96/0902 of the Grant Agency of the Czech Republic and grant A2075603 of the Academy of Sciences of the Czech Republic.

of the stochastic component is an integral expression of the stochastic behaviour of the system, sensor imprecisions and approximation (mismodelling) errors. As a result, one can rarely specify the prior density of parameters of statistical models on a purely intuitive basis.

It may be easier in such cases to express prior knowledge directly in terms of data observed on the system rather than parameters of a specific model. This alternative approach follows from the observation that the prior density of parameters, when chosen in a self-reproducing form, is a model-based transform of a certain “prior” density of data entering the model. Section 3 shows that this connection holds for any regression-type models. The difference between the two—basically dual—approaches is illustrated on the ARX model.

The amount of information available on-line is typically insufficient to specify the “prior” density of data completely. The major result of Section 4 is a method of building the prior density of model parameters from incomplete information about observed data. The solution suggested here is based upon a natural “information geometry” of probability distributions of observed data (Kulhavý, 1996).

2. BAYESIAN ESTIMATION REVISITED

The basic problem of system identification is to fit a proper model to a dynamic, possibly controlled system. The purpose of this section is to resume the classical view of Bayesian estimation of parameter estimation and to compare it with an alternative viewpoint based upon measuring an information “distance” between the empirical and model-based densities of observed data.

2.1 Probabilistic Model of Dynamic System

The models used in statistical methods of system identification describe—in probabilistic terms—how the system output depends on its past values and possible external inputs.

Sample of Data. Consider a system on which two sequences of continuous random variables are measured

$$Y^{N+m} = (Y_1, \dots, Y_{N+m}), \quad U^{N+m} = (U_1, \dots, U_{N+m})$$

which take values in subsets \mathcal{Y} and \mathcal{U} of $\mathbb{R}^{\dim y}$ and $\mathbb{R}^{\dim u}$, respectively. U_k is defined as a directly manipulated input to the system at time k while Y_k is its output—response of the system at time k to the past history of data represented by the sequences Y^{k-1} and U^k . Both the above sequences form together a *sample* of data.

A sequence of observed (measured) values

$$y^{N+m} = (y_1, \dots, y_{N+m}), \quad u^{N+m} = (u_1, \dots, u_{N+m})$$

is called a *realization* of the sample Y^{N+m} , U^{N+m} or an *observed sample*.

General Regression Model. We shall suppose that the output values Y_k depend on the past data Y_{k-m}^{k-1} , U_{k-m}^k only through a known vector function $Z_k = z(U^k, Y^{k-1})$ taking values in a subset \mathcal{Z} of $\mathbb{R}^{\dim z}$. Hence, we assume that

$$s_k(y_k | y^{k-1}, u^k) = s_k(y_k | z_k) \quad (1)$$

for $k = m + 1, \dots, N + m$. Moreover, we suppose that the conditional density of Y_k given $Z_k = z_k$ is identical for all k , i.e., $s_k(y | z) = s(y | z)$. Finally, we assume that (y_N, z_N) is recursively computable given its last value (y_{N-1}, z_{N-1}) and the latest data (y_N, u_N) , i.e., there exists a map F such that

$$(y_N, z_N) = F((y_{N-1}, z_{N-1}), (y_N, u_N)).$$

Model Family. The density $s(y | z)$ is considered to be a member of a given family

$$\mathcal{S} = \{s_\theta(y | z) : \theta \in \mathcal{T}\}$$

parametrized by the parameter θ taking values in a subset \mathcal{T} of $\mathbb{R}^{\dim \theta}$. We restrict ourselves to the case that $s_\theta(y | z) > 0$ for all $(y, z) \in \mathcal{Y} \times \mathcal{Z}$ and all $\theta \in \mathcal{T}$.

The objective of parameter estimation is to find a proper value of the parameter θ given the observed sample y^{N+m} , u^{N+m} .

Natural Conditions of Control. In general, the dependence of the input U_k on the past data Y^{k-1} , U^{k-1} and the parameter θ is expressed through a conditional density $\gamma_k(u_k | y^{k-1}, u^{k-1}, \theta)$. In most cases of practical interest, we may adopt a simplifying assumption (Peterka, 1981) that the only information about θ used in computation of the input is the information contained in the past data. More precisely, we assume that the equality

$$\gamma_k(u_k | y^{k-1}, u^{k-1}, \theta) = \gamma_k(u_k | y^{k-1}, u^{k-1}) \quad (2)$$

holds at $k = m + 1, \dots, N + m$.

2.2 Bayesian Estimation: The Classical View

Joint Density of Sample. By chain rule, the joint density q_θ^N of Y_{m+1}^{N+m} and U_{m+1}^{N+m} conditional on m initial values of Y_k and U_k can be rewritten as follows

$$\begin{aligned} q_\theta^N(y_{m+1}^{N+m}, u_{m+1}^{N+m} | y^m, u^m) \\ = \prod_{k=m+1}^{N+m} s_\theta(y_k | y^{k-1}, u^k) \gamma_k(u_k | y^{k-1}, u^{k-1}, \theta). \end{aligned}$$

Taking into account the model assumption (1) and the natural conditions of control (2), we can write

$$\begin{aligned} q_\theta^N(y_{m+1}^{N+m}, u_{m+1}^{N+m} | y^m, u^m) \\ = \prod_{k=m+1}^{N+m} s_\theta(y_k | z_k) \gamma_k(u_k | y^{k-1}, u^{k-1}). \end{aligned} \quad (3)$$

Posterior Density. When the unknown parameter θ is treated as a random variable Θ , its uncertainty is naturally described by the *posterior* density conditional on the observed sample y^{N+m}, u^{N+m}

$$p_N(\theta) \triangleq p(\theta|y^{N+m}, u^{N+m}).$$

The subscript N indicates conditioning on N data points $(y_{m+1}, z_{m+1}), \dots, (y_{N+m}, z_{N+m})$.

Given a prior density conditional on *a priori* information and possibly m initial values y^m, u^m

$$p_0(\theta) \triangleq p(\theta|y^m, u^m),$$

the posterior density $p_N(\theta)$ follows by Bayes's theorem. When substituting for the joint density $q_\theta^N(y_{m+1}^{N+m}, u_{m+1}^{N+m}|y^m, u^m)$ from (3) and taking the natural conditions of control (2) for granted, we obtain

$$\begin{aligned} p_N(\theta) &\propto p_0(\theta) q_\theta^N(y_{m+1}^{N+m}, u_{m+1}^{N+m}|y^m, u^m) \\ &\propto p_0(\theta) \prod_{k=m+1}^{N+m} s_\theta(y_k|z_k) \gamma_k(u_k|y^{k-1}, u^{k-1}) \end{aligned}$$

where the symbol \propto stands for equality up to a normalizing factor. As a result, we have the formula

$$p_N(\theta) \propto p_0(\theta) \prod_{k=m+1}^{N+m} s_\theta(y_k|z_k). \quad (4)$$

2.3 Bayesian Estimation: A View via Inaccuracy

Introducing the notion of conditional inaccuracy, we can reformulate probability-based estimation as an explicit approximation problem.

Empirical Density. Given the sample y^{N+m}, u^{N+m} , the *joint empirical density* of (Y, Z) is defined as

$$r_N(y, z) \triangleq \frac{1}{N} \sum_{k=m+1}^{N+m} \delta(y - y_k, z - z_k)$$

where $\delta(\cdot)$ is a Dirac function satisfying $\delta(y, z) = 0$ for $y \neq 0$ or $z \neq 0$ and $\iint_{y \times z} \delta(y, z) dy dz = 1$. The subscript N is used again to indicate the number of data points $(y_{m+1}, z_{m+1}), \dots, (y_{N+m}, z_{N+m})$ the empirical density is based on.

Conditional Inaccuracy. Given the joint empirical density $r_N(y, z)$ and a conditional theoretical density $s_\theta(y|z)$, we define *conditional inaccuracy* of r_N relative to s_θ as

$$\bar{K}(r_N:s_\theta) \triangleq \iint r_N(y, z) \log \frac{1}{s_\theta(y|z)} dy dz.$$

The concept of conditional inaccuracy is generalization of Kerridge's inaccuracy (Kerridge, 1961) introduced originally for the case of independent and identically distributed data.

Joint Density of Sample. The joint density of sample (3) can be rewritten as

$$q_\theta^N(y_{m+1}^{N+m}, u_{m+1}^{N+m}|y^m, u^m) = \Gamma_{N+m} \prod_{k=m+1}^{N+m} s_\theta(y_k|z_k)$$

where Γ_{N+m} is a factor independent of θ . Further, using conditional inaccuracy, we can rewrite the θ -dependent part as follows

$$\begin{aligned} &\prod_{k=m+1}^{N+m} s_\theta(y_k|z_k) \\ &= \exp\left(N \frac{1}{N} \sum_{k=1}^N \log s_\theta(y_k|z_k)\right) \\ &= \exp\left(-N \iint r_N(y, z) \log \frac{1}{s_\theta(y|z)} dy dz\right) \\ &= \exp(-N \bar{K}(r_N:s_\theta)). \end{aligned}$$

Putting both together, we obtain the following expression

$$q_\theta^N(y_{m+1}^{N+m}, u_{m+1}^{N+m}|y^m, u^m) = \Gamma_{N+m} \exp(-N \bar{K}(r_N:s_\theta)).$$

Posterior Density. Applying Bayes's theorem and substituting for the joint density of sample from the last formula, we get the *posterior* density of Θ conditional on the observed sample y^{N+m}, u^{N+m} in the form

$$p_N(\theta) \propto p_0(\theta) \exp(-N \bar{K}(r_N:s_\theta)) \quad (5)$$

which separates explicitly the key ingredients of Bayesian estimation—the amount of data, the empirical and theoretical densities of observed data and the prior density of unknown parameters.

3. CONJUGATE PRIORS

In practical estimation, it is convenient if the prior density $p_0(\theta)$ is chosen from a *conjugate* family—closed under conditioning on the observed data (Robert, 1989).

General Form of Conjugate Prior. The expression (5) suggests the general form of self-reproducing priors

$$p_\nu(\theta) \propto \exp(-\nu \bar{K}(r_\nu:s_\theta)). \quad (6)$$

Here $r_\nu(y, z)$ stands for a joint “prior” density of (Y, Z) which is based upon prior information and possibly m initial values y^m and u^m . The scalar ν describes the degree of belief in r_ν ; the bigger ν , the more concentrated the prior density $p_\nu(\theta)$ is. The factor ν is supposed to be nonnegative but not necessarily integer.

By formal analogy with (5), ν can be regarded as the number of actual or fictitious observations $r_\nu(y, z)$ is built on. The density (6) then can be seen as a “posterior” density given a uniform prior density $p_0(\theta) \propto 1$

and \mathbf{v} data with the empirical density $r_{\mathbf{v}}(y, z)$ (cf. the use of fictitious data in Kárný *et al.*, 1985).

Hence, instead of expressing prior knowledge in terms of probability distribution of parameters of a specific model, we suggest to express it directly through distribution of observed data. The latter way enables us to keep a single description of prior information even when considering different model families. The formula (6) gives a clue how the prior density of (Y, Z) is to be converted, given a particular model $s_{\theta}(y|z)$, into the prior density of Θ .

Posterior Density. It is easy to verify that given the conjugate prior (6), the posterior density (5) preserves its form

$$\begin{aligned} p_{\mathbf{v}+\mathbf{N}}(\theta) &\propto p_{\mathbf{v}}(\theta) \exp(-N\bar{K}(r_{\mathbf{N}}:s_{\theta})) \\ &\propto \exp(-\mathbf{v}\bar{K}(r_{\mathbf{v}}:s_{\theta})) \exp(-N\bar{K}(r_{\mathbf{N}}:s_{\theta})). \end{aligned}$$

Indeed, the posterior density can be written as

$$\boxed{p_{\mathbf{v}+\mathbf{N}}(\theta) \propto \exp(-(\mathbf{v} + N)\bar{K}(r_{\mathbf{v}+\mathbf{N}}:s_{\theta}))} \quad (7)$$

where

$$r_{\mathbf{v}+\mathbf{N}}(y, z) = \frac{\mathbf{v}}{\mathbf{v} + N} r_{\mathbf{v}}(y, z) + \frac{N}{\mathbf{v} + N} r_{\mathbf{N}}(y, z) \quad (8)$$

stands for a *mixture* of the prior density $r_{\mathbf{v}}(y, z)$ and the empirical density $r_{\mathbf{N}}(y, z)$. In accordance with our intuition, the weight on $r_{\mathbf{v}}(y, z)$ tends to zero as $N \rightarrow \infty$.

Exponential Family. An important special case occurs when the model family $\{s_{\theta}(y|z)\}$ is imbedded in an exponential family, i.e., when for every θ the density $s_{\theta}(y|z)$ can be expressed as

$$s_{\theta}(y|z) = s_0(y|z) \exp(\lambda^T(\theta) h(y, z) - \psi(\lambda(\theta))) \quad (9)$$

where $s_0(y|z)$ is a fixed ‘‘origin’’, $h: \mathcal{Y} \times \mathcal{Z} \rightarrow \mathbb{R}^n$ defines a canonical statistic and $\lambda \in \mathbb{R}^n$ is a natural parameter of the enveloping exponential family.

Note that the factor $\psi(\lambda(\theta))$ does *not* depend on z . This is a crucial assumption which essentially says that the logarithm of $s_{\theta}(y|z)$ can always be factorized so that the data (y, z) and the parameter θ are separated. To stress that the above definition is stronger than the usual definition of the exponential family of conditional densities, we shall call the enveloping exponential family *jointly exponential* (in both y and z).

Given the model (9), the prior density (6) takes the form

$$p_{\mathbf{v}}(\theta) \propto \exp(-\mathbf{v}\lambda^T(\theta) \bar{h}_{\mathbf{v}})$$

where

$$\bar{h}_{\mathbf{v}} = \mathcal{E}_{\mathbf{v}}(h(Y, Z)) = \iint r_{\mathbf{v}}(y, z) h(y, z) dy dz$$

denotes the expectation of the canonical statistic $h(Y, Z)$ of the enveloping exponential family with respect to the prior density $r_{\mathbf{v}}(y, z)$.

The posterior density $p_{\mathbf{v}+\mathbf{N}}(\theta)$ preserves the form of $p_{\mathbf{v}}(\theta)$

$$p_{\mathbf{v}+\mathbf{N}}(\theta) \propto \exp(-(\mathbf{v} + N)\lambda^T(\theta) \bar{h}_{\mathbf{v}+\mathbf{N}})$$

where

$$\bar{h}_{\mathbf{v}+\mathbf{N}} = \mathcal{E}_{\mathbf{v}+\mathbf{N}}(h(Y, Z)) = \iint r_{\mathbf{v}+\mathbf{N}}(y, z) h(y, z) dy dz$$

denotes the expectation of $h(Y, Z)$ with respect to the prior-modified empirical density $r_{\mathbf{v}+\mathbf{N}}(y, z)$ defined by (8).

ARX Model. Consider the model

$$Y_k = \theta^T Z_k + E_k, \quad E_k \sim N(0, \sigma^2)$$

where θ and Z_k are column vectors and E_1, E_2, \dots are independent, normally distributed random variables with zero mean and known variance σ^2 . The conditional density function of Y_k given $Z_k = z_k$ thus takes the form

$$s_{\theta}(y|z) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y - \theta^T z)^2\right).$$

After straightforward arrangements we find that the prior density (6) can be written as

$$p_{\mathbf{v}}(\theta) \propto \exp\left(-\frac{1}{2\sigma^2} \mathbf{v}(\theta - \hat{\theta}_{\mathbf{v}})^T P_{\mathbf{v}}^{-1} (\theta - \hat{\theta}_{\mathbf{v}})\right) \quad (10)$$

with

$$\begin{aligned} \hat{\theta}_{\mathbf{v}} &= \mathcal{E}_{\mathbf{v}}(ZZ^T)^{-1} \mathcal{E}_{\mathbf{v}}(ZY), \\ P_{\mathbf{v}}^{-1} &= \mathcal{E}_{\mathbf{v}}(ZZ^T) \end{aligned}$$

where $\mathcal{E}_{\mathbf{v}}(\cdot)$ stands for the expectation with respect to $r_{\mathbf{v}}(y, z)$. The matrix $P_{\mathbf{v}}^{-1}$ is supposed positive definite.

The expression (10) indicates that we can provide prior information about the ARX model parameters in two different ways—either via the expectations of model-specific functions of data (Y, Z)

$$\mathcal{E}_{\mathbf{v}}(ZY), \quad \mathcal{E}_{\mathbf{v}}(ZZ^T) \quad (11)$$

or via the mean and covariance of the parameter vector θ

$$\mathcal{E}_{\mathbf{v}}(\theta) = \hat{\theta}_{\mathbf{v}}, \quad \mathcal{E}_{\mathbf{v}}(\theta - \hat{\theta}_{\mathbf{v}})(\theta - \hat{\theta}_{\mathbf{v}})^T = \frac{\sigma^2 P_{\mathbf{v}}}{\mathbf{v}}. \quad (12)$$

Note that the subscript \mathbf{v} in (12) indicates expectation with respect to $p_{\mathbf{v}}(\theta)$.

The above raises the natural question what is a better way of expressing prior information. Although (12) is the standard choice in the Bayesian literature, (11) may often be a simpler and more natural way of expressing available knowledge. Prior information about data can often be accumulated quite cheaply, e.g., by processing archived process data or data generated by a simulation model.

4. PRIOR RESTORATION

In practice, we are often unable to specify the density $r_v(y, z)$ completely. The results for the ARX model suggest a more realistic problem formulation when, in addition to v , the expectation of a certain statistic $h(Y, Z)$ with respect to $r_v(y, z)$

$$\mathcal{E}_v(h(Y, Z)) = \bar{h}_v$$

is available.

Clearly, as $r_v(y, z)$ is not known now, the formula (6) cannot be used directly. In the following paragraphs we show how the inaccuracy $\bar{K}(r_v:s_\theta)$ can be decomposed into sum of two terms one of which can be made (approximately at least) independent of θ while the other depends on r_v only through \bar{h}_v .

Joint Exponential Family. Given any model density $s_\theta(y|z)$, we can construct a *joint* exponential family $\mathcal{S}_{\theta;h}$ composed of the joint densities

$$s_{\theta,\lambda}(y, z) = s_\theta(y|z) \exp(\lambda^T h(y, z) - \psi(\theta, \lambda)) \quad (13)$$

where $h: \mathcal{Y} \times \mathcal{Z} \mapsto \mathbb{R}^n$ is a canonical statistic of the family, $\lambda \in \mathbb{R}^n$ is a natural parameter and

$$\psi(\theta, \lambda) = \log \iint s_\theta(y|z) \exp(\lambda^T h(y, z)) dy dz$$

is logarithm of the normalizing divisor. The functions $h_0(y, z) \equiv 1, h_1(y, z), \dots, h_n(y, z)$ are assumed linearly independent.

We define a *h-projection* $s_{\theta,\hat{\lambda}}(y, z)$ of the density $r_v(y, z)$ onto the exponential family $\mathcal{S}_{\theta;h}$ by the equality

$$\iint s_{\theta,\hat{\lambda}}(y, z) h(y, z) dy dz = \iint r_v(y, z) h(y, z) dy dz \quad (14)$$

and denote the set of all densities $r(y, z)$ with the identical *h-projection* as

$$\mathcal{R}_v \triangleq \left\{ r(y, z) : \iint r(y, z) h(y, z) dy dz = \bar{h}_v, \right. \\ \left. \iint r(y, z) dy dz = 1, r(y, z) \geq 0 \right\}.$$

Pythagorean Relationship. The conditional inaccuracy $\bar{K}(r_v:s_\theta)$ can be regarded as an *unnormalized joint inaccuracy* of $r_v(y, z)$ relative to the function $s_\theta(y|z)$

$$\bar{K}(r_v:s_\theta) = \iint r_v(y, z) \log \frac{1}{s_\theta(y|z)} dy dz = K(r_v:s_\theta).$$

Let $s_{\theta,\lambda}(y, z)$ be exponential (13) and $\hat{\lambda}$ satisfy (14). It holds then

$$K(r_v:s_\theta) - K(r_v:s_{\theta,\hat{\lambda}}) \\ = \iint r_v(y, z) \log \frac{s_{\theta,\hat{\lambda}}(y, z)}{s_\theta(y|z)} dy dz$$

$$= \hat{\lambda}^T \left(\iint r_v(y, z) h(y, z) dy dz \right) - \psi(\theta, \hat{\lambda}) \\ = \hat{\lambda}^T \left(\iint s_{\theta,\hat{\lambda}}(y, z) h(y, z) dy dz \right) - \psi(\theta, \hat{\lambda}) \\ = \iint s_{\theta,\hat{\lambda}}(y, z) \log \frac{s_{\theta,\hat{\lambda}}(y, z)}{s_\theta(y|z)} dy dz \\ = D(s_{\theta,\hat{\lambda}} \| s_\theta)$$

where $D(s_{\theta,\hat{\lambda}} \| s_\theta)$ stands for an *unnormalized joint Kullback-Leibler distance* (Kullback and Leibler, 1951). This implies the following Pythagorean-like relationship that links together the information measures (cf. Fig. 1)

$$K(r_v:s_\theta) = K(r_v:s_{\theta,\hat{\lambda}}) + D(s_{\theta,\hat{\lambda}} \| s_\theta). \quad (15)$$

The formula can be seen as analogy of the Pythagorean theorem that holds for Kullback-Leibler distances between probability distributions (Čencov, 1982).

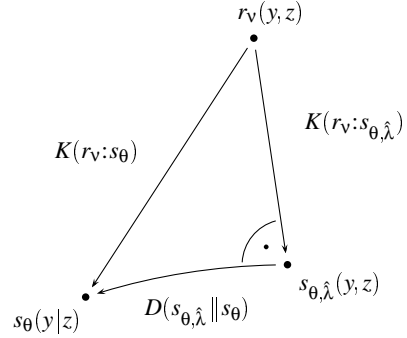


Fig. 1. Pythagorean-like decomposition of conditional inaccuracy.

Approximation of Inaccuracy. It follows directly from the Pythagorean relationship (15) that if $K(r_v:s_{\theta,\hat{\lambda}})$ can be made—by a proper choice of $h(y, z)$ —“almost independent” of θ

$$K(r_v:s_{\theta,\hat{\lambda}}) \approx \text{const.}, \quad (16)$$

then the inaccuracy $K(r_v:s_\theta)$ of the unknown $r_v(y, z)$ relative to $s_\theta(y|z)$ can be approximated as follows

$$K(r_v:s_\theta) \approx D(s_{\theta,\hat{\lambda}} \| s_\theta) + \text{const.} \quad (17)$$

With this approximation, the conjugate prior density (6) can be approximated as follows

$$\hat{p}_v(\theta) \propto \exp(-v D(s_{\theta,\hat{\lambda}} \| s_\theta)). \quad (18)$$

A number of attractive properties of the approximation stem from the use of information measures (for details see Kulhavý, 1996).

The unnormalized Kullback-Leibler distance $D(s_{\theta,\hat{\lambda}} \| s_\theta)$ can be computed from the following identity

$$0 = \min_{\lambda} D(s_{\theta,\hat{\lambda}} \| s_{\theta,\lambda})$$

$$\begin{aligned}
&= \min_{\lambda} \iint s_{\theta, \hat{\lambda}}(y, z) \log \frac{s_{\theta, \hat{\lambda}}(y, z)}{s_{\theta, \lambda}(y, z)} dy dz \\
&= D(s_{\theta, \hat{\lambda}} \| s_{\theta}) - \max_{\lambda} (\lambda^T \hat{h}(\theta, \hat{\lambda}) - \psi(\theta, \lambda)) \\
&= D(s_{\theta, \hat{\lambda}} \| s_{\theta}) - \max_{\lambda} (\lambda^T \bar{h}_v - \psi(\theta, \lambda))
\end{aligned}$$

which implies

$$\boxed{D(s_{\theta, \hat{\lambda}} \| s_{\theta}) = \max_{\lambda} (\lambda^T \bar{h}_v - \psi(\theta, \lambda))}. \quad (19)$$

Choice of h -Statistic. To ensure that (16) holds, we have to choose the statistic $h(y, z)$ in a proper way.

The optimum choice is to set the statistic $h(y, z)$ equal to the canonical statistic of a jointly exponential family (9) enveloping the model family $\{s_{\theta}(y|z)\}$. In such a case $K(r_v; s_{\theta, \hat{\lambda}})$ is independent of θ and (17) holds with equality.

Unfortunately, in many cases of practical interest, the enveloping jointly exponential family has a too large or even infinite dimension which forces us to use a lower-dimensional statistic not sufficient for precise (up to an additive constant) restoration of the inaccuracy $K(r_v; s_{\theta})$. A general class of h -statistics closest to the above optimum is constructed as follows (cf. Zacks, 1971). Consider a vector space Λ that contains functions $\lambda(\theta) = \log s_{\theta}(y|z)$ for all $(y, z) \in \mathcal{Y} \times \mathcal{Z}$. Let $L_i(\cdot)$, $i = 1, \dots, n$ be a set of linear functionals defined on the vector space Λ . Suppose in addition that the linear functionals are normalized so that $L_i(1) = 0$ for $i = 1, \dots, n$. Then define

$$\boxed{h_i(y, z) = L_i(\log s_{\theta}(y|z))} \quad (20)$$

for $i = 1, \dots, n$.

A simple example of the above definition is the logarithm of density ratio: given $n + 1$ points $\theta_1^*, \dots, \theta_{n+1}^*$ in the parameter space \mathcal{T} , we set

$$h_i(y, z) = \log \frac{s_{\theta_{i+1}^*}(y|z)}{s_{\theta_i^*}(y|z)}. \quad (21)$$

5. CONCLUDING REMARKS

The information-based framework presented in the paper offers an appealing possibility of building the prior density $p_v(\theta)$ of unknown parameters via model-based transformation of the prior density $r_v(y, z)$ of observed data. The view of prior information via observed data, which can be seen as dual to the standard Bayesian view through model parameters, suggests a number of recipes for quantification of prior knowledge:

(a) Take N archived process data, set $r_v(y, z)$ equal to the empirical density $r_N(y, z)$ and choose v according to the relevance or reliability of the data used.

(b) Do the same with data produced by a faithful (perhaps first-principle) simulation model.

(c) Express your knowledge about (Y, Z) via the prior expectation $\mathcal{E}_v(h(Y, Z))$. Choose v so as to reflect your confidence in the expectations provided.

(d) Allow yourself to be more vague by saying only

$$\bar{h}_{v, \min} \leq \mathcal{E}_v(h(Y, Z)) \leq \bar{h}_{v, \max}.$$

(e) Play with various definitions of the h -statistic. It may be easier to specify $\mathcal{E}_v(h(Y, Z))$ for some particular $h(Y, Z)$. Compare the difficulty of guessing *a priori* on the second moments of data

$$\mathcal{E}_v(YZ), \quad \mathcal{E}_v(ZZ^T)$$

or the log-likelihood ratios

$$\mathcal{E}_v \left(\log \frac{s_{\theta_{i+1}^*}(Y|Z)}{s_{\theta_i^*}(Y|Z)} \right), \quad i = 1, \dots, n.$$

The latter may be easier, especially when having experience with a set of particular models $s_{\theta_i^*}(y|z)$ on archived data or a simulation model. Since the latter is the empirical expectation of (21), it may also bring more information into estimation, especially for non-linear or non-Gaussian models.

REFERENCES

- Čencov, N. N. (1982). *Statistical Decision Rules and Optimal Inference*, Vol. 53 of *Transl. of Math. Monographs*. Amer. Math. Soc., Providence, RI.
- Kárný, M., A. Halousková, J. Böhm, R. Kulhavý and P. Nedoma (1985). Design of linear quadratic adaptive control: theory and algorithms for practice. *Kybernetika*, Supplement to No. 3–6.
- Kerridge, D. F. (1961). Inaccuracy and inference. *J. Roy. Statist. Soc. Ser. B*, **23**, 284–294.
- Kulhavý, R. (1996). *Recursive Nonlinear Estimation: A Geometric Approach*, Vol. 216 of *Lecture Notes in Control and Information Sciences*. Springer-Verlag, London.
- Kulhavý, R. and F. J. Kraus (1996). On duality of regularized exponential and linear forgetting. *Automatica*, **32**, 1403–1415.
- Kullback, S. and R. A. Leibler (1951). On information and sufficiency. *Ann. Math. Statist.*, **22**, 79–86.
- Peterka, V. (1981). Bayesian approach to system identification. In: *Trends and Progress in System Identification* (P. Eykhoff, Ed.), Chap. 8, pp. 239–304. Pergamon Press, Elmsford, N.Y.
- Robert, C. P. (1989). *The Bayesian Choice: A Decision-Theoretic Motivation*. Springer-Verlag, Berlin.
- Zacks, S. (1971). *The Theory of Statistical Inference*. Wiley, New York.