

A SUFFICIENT STATISTIC AND A NONSTANDARD LINEARIZATION IN NONLINEAR REGRESSION

ANDREJ PÁZMAN

In a nonlinear model $\mathbf{y} = \boldsymbol{\eta}(\boldsymbol{\theta}) + \boldsymbol{\varepsilon}$ a standard linearization consists in linearizing $\boldsymbol{\eta}(\boldsymbol{\theta})$ at a point $\boldsymbol{\theta}^*$, and in computing the M. L. estimate $\boldsymbol{\tau}(\mathbf{y}, \boldsymbol{\theta}^*)$ in the linearized model. We propose to take $\boldsymbol{\tau}(\mathbf{y}) := (\boldsymbol{\tau}(\mathbf{y}, \boldsymbol{\theta}^1), \dots, \boldsymbol{\tau}(\mathbf{y}, \boldsymbol{\theta}^k))^T$ for some $\boldsymbol{\theta}^1, \dots, \boldsymbol{\theta}^k$ (= the sufficient statistic), linearize each $\boldsymbol{\tau}(\mathbf{y}, \boldsymbol{\theta}^i)$ separately, and then to compute the M. L. estimate $\hat{\boldsymbol{\theta}}(\mathbf{y})$. The variable $\hat{\boldsymbol{\theta}}(\mathbf{y})$ has a smaller variance than $\boldsymbol{\tau}(\mathbf{y}, \boldsymbol{\theta}^i)$, and a comparable bias. Further, $\hat{\boldsymbol{\theta}}(\mathbf{y})$ can be used to approximate the posterior density in a Bayesian approach.

The construction of the sufficient statistic has a geometrical background. Possible consequences for nonlinear experimental design are mentioned.

1. INTRODUCTION AND THE GEOMETRICAL BACKGROUND

Let us consider the nonlinear regression model with normal errors

$$(1) \quad \begin{aligned} \mathbf{y} &= \boldsymbol{\eta}(\boldsymbol{\theta}) + \boldsymbol{\varepsilon}; \quad (\boldsymbol{\theta} \in \Theta) \\ \boldsymbol{\varepsilon} &\sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}) \end{aligned}$$

under standard regularity assumptions: the parameter space Θ is an open subset of \mathbb{R}^m , the variance matrix $\boldsymbol{\Sigma}$ is regular, the regression mapping $\boldsymbol{\eta}: \Theta \mapsto \mathbb{R}^N$ ($N > m$) has continuous second order derivatives on Θ , and the vectors $\partial\boldsymbol{\eta}(\boldsymbol{\theta})/\partial\theta_1, \dots, \partial\boldsymbol{\eta}(\boldsymbol{\theta})/\partial\theta_m$ are linearly independent for every $\boldsymbol{\theta} \in \Theta$. The vector $\mathbf{y} \in \mathbb{R}^N$ is observed, the mapping $\boldsymbol{\eta}$ and the set Θ are known, $\boldsymbol{\Sigma}$ is either known, or of the form $\boldsymbol{\Sigma} = c\mathbf{W}$ with $c > 0$ unknown and \mathbf{W} known. Statistical inference on the unknown vector $\boldsymbol{\theta}$ should be performed.

A well known point estimator in model (1) is the maximum likelihood (= M. L.) estimator

$$(2) \quad \hat{\boldsymbol{\theta}} := \hat{\boldsymbol{\theta}}(\mathbf{y}) := \arg \min_{\boldsymbol{\theta} \in \Theta} \|\mathbf{y} - \boldsymbol{\eta}(\boldsymbol{\theta})\|_{\mathbf{W}}^2.$$

Here $\|\mathbf{a}\|_{\mathbf{W}}^2 := \mathbf{a}^T \mathbf{W}^{-1} \mathbf{a}$; ($\mathbf{a} \in \mathbb{R}^N$).

In the particular case when model (1) is linear, the statistic $\mathbf{y} \in \mathbb{R}^N \mapsto \hat{\boldsymbol{\theta}}(\mathbf{y})$ is not only a point estimator, it is also a sufficient statistic. If model (1) is nonlinear (more exactly, if the expectation surface of model (1)

$$\mathcal{E} := \{\boldsymbol{\eta}(\boldsymbol{\theta}); \boldsymbol{\theta} \in \Theta\}$$

is not a “plane”, the statistic $\hat{\boldsymbol{\theta}}(\mathbf{y})$ is no more sufficient. Consequently it contains less information about $\boldsymbol{\theta}$ than the sample vector \mathbf{y} . (For the distributional properties of $\hat{\boldsymbol{\theta}}(\mathbf{y})$ cf. e.g. [4, 5]).

However, it is possible to look for a statistic in model (1) which is a sufficient statistic, and which is somehow related to the M. L. estimator. In particular, we can require that this statistic coincides with $\hat{\boldsymbol{\theta}}(\mathbf{y})$ when model (1) is linear.

In Section 2 we propose such statistics. They have the following geometrical origin:

Consider the expectation surface \mathcal{E} . It is an m -dimensional surface in the N dimensional sample space \mathbb{R}^N . According to (2), the point $\boldsymbol{\eta}(\hat{\boldsymbol{\theta}}) \in \mathcal{E}$ is obtained by the \mathbf{W} -orthogonal projection of the point \mathbf{y} onto \mathcal{E} . Consider now for any $\boldsymbol{\theta}^* \in \Theta$ the set

$$T_{\boldsymbol{\theta}^*} := \left\{ \boldsymbol{\eta}(\boldsymbol{\theta}^*) + \frac{\partial \boldsymbol{\eta}(\boldsymbol{\theta}^*)}{\partial \boldsymbol{\theta}^T} \mathbf{v}; \mathbf{v} \in \mathbb{R}^m \right\}.$$

Geometrically, $T_{\boldsymbol{\theta}^*}$ is the tangent plane to the surface \mathcal{E} at the point $\boldsymbol{\eta}(\boldsymbol{\theta}^*) \in \mathcal{E}$. Statistically, $T_{\boldsymbol{\theta}^*}$ is the expectation surface of a linear model which approximates model (1):

$$(3) \quad \mathbf{y} - \boldsymbol{\eta}(\boldsymbol{\theta}^*) = \frac{\partial \boldsymbol{\eta}(\boldsymbol{\theta}^*)}{\partial \boldsymbol{\theta}^T} (\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \boldsymbol{\varepsilon}$$

$$\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}).$$

The M. L. estimate in this linearized model is

$$(4) \quad \boldsymbol{\tau}(\mathbf{y}, \boldsymbol{\theta}^*) := \arg \min_{\boldsymbol{\theta}} \left\| \mathbf{y} - \boldsymbol{\eta}(\boldsymbol{\theta}^*) - \frac{\partial \boldsymbol{\eta}(\boldsymbol{\theta}^*)}{\partial \boldsymbol{\theta}^T} (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right\|_{\mathbf{W}}^2.$$

It is the result of the \mathbf{W} -orthogonal projection of the point \mathbf{y} onto $T_{\boldsymbol{\theta}^*}$.

The statistic $\mathbf{y} \mapsto \boldsymbol{\tau}(\mathbf{y}, \boldsymbol{\theta}^*)$ is sufficient in model (3), however, it is not in model (1). Therefore, we proceed further by considering not one but many (eventually all) tangent planes to \mathcal{E} , and by projecting \mathbf{W} -orthogonally the sample point \mathbf{y} onto all of them. (The reader which is familiar with differential geometry see that we are using the “tangent space” of \mathcal{E}). Consequently, instead of one random vector $\boldsymbol{\tau}(\mathbf{y}, \boldsymbol{\theta}^*)$ we consider the set of random vectors

$$(5) \quad \{\boldsymbol{\tau}(\mathbf{y}, \boldsymbol{\theta}^*); \boldsymbol{\theta}^* \in D\}$$

for some $D \subset \Theta$. Evidently, this is a (vector-valued) random process defined on D . This process will be shown to have several pleasant structural properties.

a) It is a Gaussian random process having a covariance function which does not depend on $\boldsymbol{\theta}$.

b) Each component $\tau(\mathbf{y}, \theta^*)$ of this process is related to a linear approximative model.

c) When D is adequately chosen, the mapping

$$\mathbf{y} \in \mathbb{R}^N \mapsto \{\tau(\mathbf{y}, \theta^*); \theta^* \in D\}$$

is a sufficient statistic in model (1).

In Section 3–5 we try to demonstrate that such a process is useful. We restrict our attention to the case of a finite $D = \{\theta^1, \dots, \theta^k\}$, and instead of the process we consider a $k \cdot m$ dimensional random vector $\tau(\mathbf{y}) := (\tau^T(\mathbf{y}, \theta^1), \dots, \tau^T(\mathbf{y}, \theta^k))^T$. If we linearize each component of $\tau(\mathbf{y})$ separately, we obtain a new, nonstandard linearization of model (1) which is more efficient than the standard linearization (3) (see Proposition 2). This allows to obtain an approximative expression for the posterior probability density of θ (Proposition 3). Moreover, using quadratic functions of $\tau(\mathbf{y})$ we can discuss some confidence regions for θ , both for the case when Σ is known and when $\Sigma = \sigma^2 \mathbf{I}$ with an unknown σ .

2. SUFFICIENT STATISTICS

As is well known (cf. [1], Chapt. VIII.1.), the M.L. estimate of θ in the linear model (3) can be expressed in the form

$$(6) \quad \tau(\mathbf{y}, \theta^*) = \mathbf{M}^{-1}(\theta^*) \mathbf{F}^T(\theta^*) \mathbf{W}^{-1} [\mathbf{y} - \eta(\theta^*)] + \theta^*$$

where

$$\{\mathbf{F}(\theta)\}_{ij} := \frac{\partial \eta_i(\theta)}{\partial \theta_j}; \quad (i = 1, \dots, N, \quad j = 1, \dots, m),$$

$$\mathbf{M}(\theta) := \mathbf{F}^T(\theta) \mathbf{W}^{-1} \mathbf{F}(\theta).$$

Consequently, (5) is a Gaussian random process with the mean

$$(7) \quad \mathbf{m}_\theta(\theta^*) = \mathbf{M}^{-1}(\theta^*) \mathbf{F}^T(\theta^*) \mathbf{W}^{-1} [\eta(\theta) - \eta(\theta^*)] + \theta^*; \quad (\theta^* \in D)$$

and with the covariance function $c\mathbf{K}(\theta^*, \theta^0)$ where

$$\mathbf{K}(\theta^*, \theta^0) = \mathbf{M}^{-1}(\theta^*) \mathbf{F}^T(\theta^*) \mathbf{W}^{-1} \mathbf{F}(\theta^0) \mathbf{M}^{-1}(\theta^0).$$

We see that $\mathbf{K}(\theta^*, \theta^0)$ does not depend on θ (= the true value of the parameters).

When the set D is finite, $D = \{\theta^1, \dots, \theta^k\}$, it is better to consider the $(m \cdot k)$ -dimensional random vector

$$(8) \quad \tau := \tau(\mathbf{y}) := \begin{pmatrix} \tau(\mathbf{y}, \theta^1) \\ \vdots \\ \tau(\mathbf{y}, \theta^k) \end{pmatrix},$$

instead of the random process (5). Here each component $\tau(\mathbf{y}, \theta^i)$ is defined according to (6).

The mean and the variance matrix of τ are equal to

$$\mathbf{m}_\theta := \mathbf{E}_\theta(\tau) = (\mathbf{m}_\theta^T(\theta^1), \dots, \mathbf{m}_\theta^T(\theta^k))^T$$

$$\text{Var}(\tau) = c \mathbf{S}$$

where

$$(9) \quad \mathbf{S} := \begin{pmatrix} \mathbf{K}(\theta^1, \theta^1), & \dots, & \mathbf{K}(\theta^1, \theta^k) \\ \mathbf{K}(\theta^k, \theta^1), & \dots, & \mathbf{K}(\theta^k, \theta^k) \end{pmatrix}.$$

If \mathbf{A} is any $r \times s$ matrix, we denote by $\mathcal{M}(\mathbf{A}) := \{\mathbf{A}\mathbf{u} : \mathbf{u} \in \mathbb{R}^s\}$ the linear subspace of \mathbb{R}^r spanned by the columns of \mathbf{A} .

Proposition 1. If for every $\theta \in \Theta$

$$\mathcal{M}[\mathbf{F}(\theta)] \subset \mathcal{M}[(\mathbf{F}(\theta^1), \dots, \mathbf{F}(\theta^k))],$$

then the statistic

$$\mathbf{y} \in \mathbb{R}^N \mapsto \tau(\mathbf{y}) \in \mathbb{R}^{mk}$$

is sufficient in model (1).

Proof. Let \mathcal{L} be the linear manifold in \mathbb{R}^N (the “plane”) spanned by the set

$$\bigcup_{\theta \in \Theta} T_\theta.$$

Let us define

$$\mathbf{z}^\wedge := \mathbf{z}^\wedge(\mathbf{y}) := \arg \min_{\mathbf{z} \in \mathcal{L}} \|\mathbf{y} - \mathbf{z}\|_{\mathbf{W}}^2$$

The probability density of \mathbf{y} is equal to

$$f(\mathbf{y} | \theta) = (2\pi)^{-N/2} \det^{-1/2}(\boldsymbol{\Sigma}) \exp\{-\|\mathbf{y} - \boldsymbol{\eta}(\theta)\|_{\mathbf{W}}^2/(2c)\} \approx$$

$$\approx \exp\{-\|\mathbf{y} - \mathbf{z}^\wedge(\mathbf{y})\|_{\mathbf{W}}^2/(2c)\} \exp\{-\|\mathbf{z}^\wedge - \boldsymbol{\eta}(\theta)\|_{\mathbf{W}}^2/(2c)\}; \quad (\theta \in \Theta).$$

Hence, according to the factorisation theorem (cf. [1], Chapt. XV. 5.), the statistic $\mathbf{z}^\wedge(\mathbf{y})$ is sufficient in model (1).

Denote by

$$\mathbf{P}_i := \mathbf{F}(\theta^i) \mathbf{M}^{-1}(\theta^i) \mathbf{F}^T(\theta^i) \mathbf{W}^{-1}$$

the \mathbf{W} -orthogonal projector onto $\mathcal{M}[\mathbf{F}(\theta^i)]$. The mapping $\mathbf{z} \mapsto (\mathbf{P}_1(\mathbf{z} - \boldsymbol{\eta}(\theta^1)), \dots, \mathbf{P}_k(\mathbf{z} - \boldsymbol{\eta}(\theta^k)))$ is one-to one on \mathcal{L} . Indeed, take $\mathbf{z}, \mathbf{z}^* \in \mathcal{L}$ such that

$$\mathbf{P}_i(\mathbf{z} - \boldsymbol{\eta}(\theta^i)) = \mathbf{P}_i(\mathbf{z}^* - \boldsymbol{\eta}(\theta^i)); \quad (i = 1, \dots, k).$$

Multiplying by $\mathbf{F}^T(\theta^i) \mathbf{W}^{-1}$ from the left, we obtain

$$\mathbf{F}^T(\theta^i) \mathbf{W}^{-1}(\mathbf{z} - \mathbf{z}^*) = 0; \quad (i = 1, \dots, k),$$

i.e. $(\mathbf{z} - \mathbf{z}^*)$ is \mathbf{W} -orthogonal to $\mathcal{M}[(\mathbf{F}(\theta^1), \dots, \mathbf{F}(\theta^k))]$, hence to \mathcal{L} . Consequently, $(\mathbf{z} - \mathbf{z}^*)^T \mathbf{W}^{-1}(\mathbf{z} - \mathbf{z}^*) = 0$ hence, $\mathbf{z} = \mathbf{z}^*$.

It follows that $\mathbf{y} \in \mathbb{R}^N \mapsto (\mathbf{P}_1(\mathbf{z}^\wedge(\mathbf{y}) - \boldsymbol{\eta}(\theta^1)), \dots, \mathbf{P}_k(\mathbf{z}^\wedge(\mathbf{y}) - \boldsymbol{\eta}(\theta^k)))$ is a sufficient statistic in model (1).

Since $\mathbf{z}^\wedge(\mathbf{y})$ is the \mathbf{W} -orthogonal projection of \mathbf{y} onto \mathcal{L} we have

$$\mathbf{P}_i(\mathbf{z}^\wedge(\mathbf{y}) - \boldsymbol{\eta}(\theta^i)) = \mathbf{P}_i(\mathbf{y} - \boldsymbol{\eta}(\theta^i)); \quad (i = 1, \dots, k).$$

Further, the equality

$$\mathbf{F}(\theta^i) \boldsymbol{\tau}(\mathbf{y}, \theta^i) = \mathbf{P}_i(\mathbf{y} - \boldsymbol{\eta}(\theta^i)) + \mathbf{F}(\theta^i) \theta^i$$

which follows from Eq. (6), specifies $\boldsymbol{\tau}(\mathbf{y}, \theta^i)$ uniquely, since $\mathbf{F}(\theta^i)$ is of full rank. Consequently the mapping $\boldsymbol{\tau}(\mathbf{y}) \mapsto (\mathbf{P}_1(\mathbf{z}^\wedge(\mathbf{y}) - \boldsymbol{\eta}(\theta^1)), \dots, \mathbf{P}_k(\mathbf{z}^\wedge(\mathbf{y}) - \boldsymbol{\eta}(\theta^k)))$ is one-to-one. It follows that $\boldsymbol{\tau}(\mathbf{y})$ is a sufficient statistic in model (1). \square

Corollary 1. If $D \subset \Theta$ is such that

$$\mathcal{M}[\mathbf{F}(\theta)] \subset \mathcal{M}[(\mathbf{F}(\theta^1), \dots, \mathbf{F}(\theta^k))]; \quad (\theta \in \Theta)$$

for some finite set $\{\theta^1, \dots, \theta^k\} \subset D$, then

$$\mathbf{y} \in \mathbb{R}^N \mapsto \{\boldsymbol{\tau}(\mathbf{y}, \theta^*); \theta^* \in D\}$$

is sufficient. Particularly

$$\mathbf{y} \in \mathbb{R}^N \mapsto \{\boldsymbol{\tau}(\mathbf{y}, \theta^*); \theta^* \in \Theta\}$$

is always sufficient.

Corollary 2. Let $\pi(\theta)$ be a probability density on Θ (the prior density) such that

$$\mathcal{M}[\mathbf{F}(\theta)] \subset \mathcal{M}[(\mathbf{F}(\theta^1), \dots, \mathbf{F}(\theta^k))]; \quad (\theta \in \text{supp}(\pi)).$$

Then

$$\pi(\theta | \boldsymbol{\tau}(\mathbf{y})) = \pi(\theta | \mathbf{y})$$

where $\pi(\theta | \mathbf{u})$ denotes the posterior density of θ given \mathbf{u} .

Proof. As in Proposition 1, we can prove that $\boldsymbol{\tau}(\mathbf{y})$ is sufficient in the model

$$\mathbf{y} = \boldsymbol{\eta}(\theta) + \boldsymbol{\varepsilon}; \quad (\theta \in \text{supp}(\pi)).$$

Hence $f(\mathbf{y} | \theta)$ can be factorized, i.e. we can write

$$f(\mathbf{y} | \theta) = h(\mathbf{y}) g(\boldsymbol{\tau}(\mathbf{y}), \theta)$$

for some functions h and g . It follows that

$$\pi(\theta | \mathbf{y}) = \frac{f(\mathbf{y} | \theta) \pi(\theta)}{\int_{\text{supp}(\pi)} f(\mathbf{y} | \mathbf{t}) \pi(\mathbf{t}) d\mathbf{t}} = \frac{g(\boldsymbol{\tau}(\mathbf{y}), \theta) \pi(\theta)}{\int_{\text{supp}(\pi)} g(\boldsymbol{\tau}(\mathbf{y}), \mathbf{t}) \pi(\mathbf{t}) d\mathbf{t}}$$

Hence $\mathbf{y} \in \mathbb{R}^N \mapsto \pi(\theta | \mathbf{y})$ is a function of $\boldsymbol{\tau}(\mathbf{y})$. According to the definition of conditional distributions (cf. [7], Chapt. V. 1.) it means that $\pi(\theta | \mathbf{y}) = \pi(\theta | \boldsymbol{\tau}(\mathbf{y}))$. \square

3. A NONSTANDARD LINEARIZATION

Let us consider the random vector $\boldsymbol{\tau}(\mathbf{y})$ (the sufficient statistic) defined in Eq. (8). We have

$$(10) \quad \boldsymbol{\tau}(\mathbf{y}) \sim \mathcal{N}(\mathbf{m}_\theta, c \mathbf{S}); \quad (\theta \in \Theta)$$

where \mathbf{m}_θ and \mathbf{S} are given by (9).

Instead of taking the linearization (3) we propose to linearize (10), i.e. to take

approximately

$$(11) \quad \tau(\mathbf{y}) \sim \mathcal{N}(\mathbf{J}\theta, c \mathbf{S}); \quad (\theta \in \mathbb{R}^m)$$

where

$$\mathbf{J} := (\mathbf{1}, \dots, \mathbf{1})^T$$

and $\mathbf{1}$ is the $m \times m$ identity matrix. The linearization (11) is the linearization (3) applied separately to each component $\tau(\mathbf{y}, \theta^i)$ of the vector $\tau(\mathbf{y})$.

To compare the standard linearization (3) with (11) take for θ^* any point of the set $\{\theta^1, \dots, \theta^k\}$, say $\theta^* = \theta^1$. Then consider the BLUE-s (= best linear unbiased estimates) of θ in both models. The BLUE in model (3) is equal to $\tau(\mathbf{y}, \theta^1)$, and is expressed in Eq. (6). Although the matrix \mathbf{S} is singular (in general), and $\mathcal{M}(\mathbf{J}) \not\subset \mathcal{M}(\mathbf{S})$, the vector θ can be estimated without bias in model (11), say by the estimate

$$\frac{1}{k} \mathbf{J}^T \tau(\mathbf{y}).$$

Hence the BLUE exists also in model (11). Let us denote it by $\tilde{\theta}(\mathbf{y})$. We refer to [3], Theorems 5.2.2 and 5.2.5 for explicit expressions for $\tilde{\theta}(\mathbf{y})$ and $\text{Var } \tilde{\theta}(\mathbf{y})$. We have

$$\tilde{\theta}(\mathbf{y}) = \mathbf{Q} \tau(\mathbf{y}), \quad \text{Var } \tilde{\theta}(\mathbf{y}) = c \mathbf{V}$$

where

$$(12) \quad \begin{aligned} \mathbf{Q} &:= [\mathbf{J}^T(\mathbf{S} + \mathbf{J}\mathbf{J}^T)^{-1} \mathbf{J}]^{-1} \mathbf{J}^T(\mathbf{S} + \mathbf{J}\mathbf{J}^T)^{-1} \\ \mathbf{V} &:= [\mathbf{J}^T(\mathbf{S} + \mathbf{J}\mathbf{J}^T)^{-1} \mathbf{J}]^{-1} - \mathbf{1} \end{aligned}$$

We note that $\mathbf{J}^T(\mathbf{S} + \mathbf{J}\mathbf{J}^T)^{-1} \mathbf{J}$ is nonsingular, since \mathbf{J} is of full rank and $\mathcal{M}[\mathbf{J}] = \mathcal{M}[\mathbf{J}\mathbf{J}^T] \subset \mathcal{M}[\mathbf{S} + \mathbf{J}\mathbf{J}^T]$. In the particular case that \mathbf{S} is regular, we have simpler formulae

$$(13) \quad \begin{aligned} \mathbf{Q} &= (\mathbf{J}^T \mathbf{S}^{-1} \mathbf{J})^{-1} \mathbf{J}^T \mathbf{S}^{-1} \\ \mathbf{V} &= (\mathbf{J}^T \mathbf{S}^{-1} \mathbf{J})^{-1} \end{aligned}$$

Hence in the linearized model (11) we have

$$(14) \quad \tilde{\theta}(\mathbf{y}) \sim \mathcal{N}(\theta, c \mathbf{V}); \quad (\theta \in \mathbb{R}^m)$$

but in the linearized model (3) we have

$$(15) \quad \tau(\mathbf{y}, \theta^1) \sim \mathcal{N}(\theta, c \mathbf{M}^{-1}(\theta^1)); \quad (\theta \in \mathbb{R}^m).$$

To compare what linearization is better, we shall compare the exact distributions of $\tilde{\theta}(\mathbf{y})$ and $\tau(\mathbf{y}, \theta^1)$.

Proposition 2. The random vectors $\tilde{\theta}(\mathbf{y})$ and $\tau(\mathbf{y}, \theta^1)$ are exactly distributed according to

$$(16) \quad \tilde{\theta}(\mathbf{y}) \sim \mathcal{N}(\mathbf{Q}\mathbf{m}_\theta, c \mathbf{V}); \quad (\theta \in \Theta)$$

$$(17) \quad \tau(\mathbf{y}, \theta^1) \sim \mathcal{N}(\mathbf{m}_\theta(\theta^1), c \mathbf{M}^{-1}(\theta^1)); \quad (\theta \in \Theta).$$

The vectors expressing the bias

$$\mathbf{Q}\mathbf{m}_\theta - \theta$$

and the bias

$$m_{\theta}(\theta^1) - \theta$$

are of the sample order of magnitude. The estimator $\tilde{\theta}(\mathbf{y})$ is more efficient since the matrix $\text{Var} [\tau(\mathbf{y}, \theta^1)] - \text{Var} [\tilde{\theta}(\mathbf{y})]$ is positive semidefinite.

Proof. Both variables $\tilde{\theta}(\mathbf{y})$ and $\tau(\mathbf{y}, \theta^1)$ are linear in \mathbf{y} , hence they are normally distributed. The mean and the variance of $\tau(\mathbf{y}, \theta^1)$ is given in Eq. (7). The mean and the variance of $\tilde{\theta}(\mathbf{y})$ follow from Eq. (12) and from the mean and the variance of $\tau(\mathbf{y})$ in Eq. (9).

The bias of $\tilde{\theta}(\mathbf{y})$ is

$$\mathbf{Q}m_{\theta} - \theta = \mathbf{Q} \begin{pmatrix} m_{\theta}(\theta^1) \\ \vdots \\ m_{\theta}(\theta^k) \end{pmatrix} - \theta.$$

The bias of $\tau(\mathbf{y}, \theta^1)$ can be written in the form

$$m_{\theta}(\theta^1) - \theta = \mathbf{QJ} m_{\theta}(\theta^1) - \theta = \mathbf{Q} \begin{pmatrix} m_{\theta}(\theta^1) \\ \vdots \\ m_{\theta}(\theta^1) \end{pmatrix} - \theta$$

since $\mathbf{QJ} = \mathbf{I}$, according to (12) and (13). Thus if $m_{\theta}(\theta^i) - \theta$ is of the same order for every $i = 1, \dots, k$, then $\tau(\mathbf{y}, \theta^1)$ and $\tilde{\theta}(\mathbf{y})$ have the bias of the same order as well.

The random variable $\tau(\mathbf{y}, \theta^1)$ can be written in the form

$$\tau(\mathbf{y}, \theta^1) = (\mathbf{1}, \mathbf{0}, \dots, \mathbf{0}) \tau(\mathbf{y})$$

hence it is a linear unbiased estimator of θ in model (11). Since $\tilde{\theta}(\mathbf{y})$ is the BLUE in the same model, it follows that $\text{Var} [\tau(\mathbf{y}, \theta^1)] - \text{Var} [\tilde{\theta}(\mathbf{y})]$ is positive semidefinite. \square

Note 1. According to Eq. (7) we can write the bias in the form

$$m_{\theta}(\theta^1) - \theta := r(\theta, \theta^1), \quad \mathbf{J}m_{\theta} - \theta = \mathbf{J} \begin{pmatrix} r(\theta, \theta^1) \\ \vdots \\ r(\theta, \theta^k) \end{pmatrix}$$

where from the Taylor formula for $\eta(\theta)$ at θ^i we obtain

$$(18) \quad \begin{aligned} r(\theta, \theta^i) &:= \mathbf{M}^{-1}(\theta^i) \mathbf{F}^T(\theta^i) \mathbf{W}^{-1} [\eta(\theta) - \eta(\theta^i)] + \theta^i - \theta = \\ &= \frac{1}{2} \mathbf{M}^{-1}(\theta^i) \mathbf{F}^T(\theta^i) \mathbf{W}^{-1} \left[(\theta - \theta^i)^T \left[\frac{\partial^2 \eta(\theta)}{\partial \theta \partial \theta^T} \right]_{\lambda \theta + (1-\lambda)\theta^i} (\theta - \theta^i) \right] \end{aligned}$$

for some number $\lambda \in (0, 1)$ depending on θ and on θ^i .

The expression for $r(\theta, \theta^i)$ is small either if $[\theta - \theta^i]^T \mathbf{M}(\theta^i) [\theta - \theta^i]$ is small or if model (1) is not too much curved, since

$$\sup \left\{ \left\| \mathbf{v}^T \frac{\partial^2 \eta(\theta)}{\partial \theta \partial \theta^T} \mathbf{v} \right\|_{\mathbf{W}} / \mathbf{v}^T \mathbf{M}(\theta^i) \mathbf{v}; \quad 0 \neq \mathbf{v} \in \mathbb{R}^m \right\}$$

is related to the curvatures of Bates and Watts [2] in model (1). We used here the notation

$$\mathbf{v}^T \frac{\partial^2 \boldsymbol{\eta}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \mathbf{v} := \sum_{ij} v_i \frac{\partial^2 \boldsymbol{\eta}(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} v_j.$$

It is important to note here that $E_{\theta}[\tilde{\boldsymbol{\theta}}(\mathbf{y})] = \mathbf{Qm}_{\theta}$ is a “mixture” of the means of $\boldsymbol{\tau}(\mathbf{y}, \boldsymbol{\theta}^1), \dots, \boldsymbol{\tau}(\mathbf{y}, \boldsymbol{\theta}^k)$. In some cases the “mixture” is such that the bias of $\tilde{\boldsymbol{\theta}}(\mathbf{y})$ is much smaller than the bias of every $\boldsymbol{\tau}(\mathbf{y}, \boldsymbol{\theta}^i)$. This depends on the choice off $\boldsymbol{\theta}^1, \dots, \boldsymbol{\theta}^k$.

Note 2. When $\boldsymbol{\tau}(\mathbf{y})$ is a sufficient statistic (Proposition 1) we arrive to $\tilde{\boldsymbol{\theta}}(\mathbf{y})$ according to the scheme

$$\begin{aligned} (1) \mapsto \text{sufficient statistic } \boldsymbol{\tau} \mapsto (10) \mapsto \text{linearization of } \boldsymbol{\tau} \mapsto (11) \quad \mapsto \\ \mapsto \text{sufficient statistic } \tilde{\boldsymbol{\theta}} \mapsto (14) \end{aligned}$$

Example 1. We shall consider the simple nonlinear model $\mathbf{y} = \boldsymbol{\eta}(\boldsymbol{\theta}) + \boldsymbol{\varepsilon}$ with $N = 2, m = 1, \Theta = (0, \pi), \boldsymbol{\eta}(\boldsymbol{\theta}) = (\cos \theta, \sin \theta)^T, \boldsymbol{\Sigma} = \mathbf{W} = \mathbf{I}$. (The expectation surface is a halfcircle). In this case we have $\mathbf{F}^T(\boldsymbol{\theta}) = (-\sin \theta, \cos \theta), \mathbf{M}(\boldsymbol{\theta}) = 1; (\boldsymbol{\theta} \in \Theta)$. To construct $\boldsymbol{\tau}(\mathbf{y})$ take two points $\theta^1 = \theta^* - \delta, \theta^2 = \theta^* + \delta$ for some fixed $\delta > 0, \theta^* \in \Theta$. By simple computations we obtain

$$\begin{aligned} \boldsymbol{\tau}(\mathbf{y}) &= \begin{pmatrix} -y_1 \sin(\theta^* - \delta) + y_2 \cos(\theta^* - \delta) + \theta^* - \delta \\ -y_1 \sin(\theta^* + \delta) + y_2 \cos(\theta^* + \delta) + \theta^* + \delta \end{pmatrix} \\ \mathbf{S} &= \begin{pmatrix} 1, & \cos(2\delta) \\ \cos(2\delta), & 1 \end{pmatrix}, \quad \mathbf{S}^{-1} = \begin{pmatrix} 1, & -\cos(2\delta) \\ -\cos(2\delta), & 1 \end{pmatrix} \Big/ \sin^2(2\delta) \\ \mathbf{V} &= (\mathbf{J}^T \mathbf{S}^{-1} \mathbf{J})^{-1} = (1 + \cos 2\delta)/2 = \cos^2 \delta < 1 \\ \tilde{\boldsymbol{\theta}}(\mathbf{y}) &= (\mathbf{J}^T \mathbf{S}^{-1} \mathbf{J})^{-1} \mathbf{J}^T \mathbf{S}^{-1} \boldsymbol{\tau}(\mathbf{y}) = \cos \delta [-y_1 \sin \theta^* + y_2 \cos \theta^*] + \theta^* \end{aligned}$$

When $\delta \mapsto 0$ we obtain $\boldsymbol{\tau}(\mathbf{y}, \theta^*)$:

$$\boldsymbol{\tau}(\mathbf{y}, \theta^*) = [-y_1 \sin \theta^* + y_2 \cos \theta^*] + \theta^*.$$

Further

$$\begin{aligned} E_{\theta}[\tilde{\boldsymbol{\theta}}(\mathbf{y})] - \boldsymbol{\theta} &= \cos \delta \sin(\boldsymbol{\theta} - \theta^*) + (\theta^* - \boldsymbol{\theta}), \\ E_{\theta}[\boldsymbol{\tau}(\mathbf{y}, \theta^*)] - \boldsymbol{\theta} &= \sin(\boldsymbol{\theta} - \theta^*) + (\theta^* - \boldsymbol{\theta}). \end{aligned}$$

Hence for δ not very large, the bias of $\tilde{\boldsymbol{\theta}}(\mathbf{y})$ and of $\boldsymbol{\tau}(\mathbf{y}, \theta^*)$ is approximatively the same. The mean square error of $\tilde{\boldsymbol{\theta}}(\mathbf{y})$ is equal to

$$E_{\theta}[\tilde{\boldsymbol{\theta}}(\mathbf{y}) - \boldsymbol{\theta}]^2 = \cos^2 \delta + [\cos \delta \sin(\boldsymbol{\theta} - \theta^*) + (\theta^* - \boldsymbol{\theta})]^2 := \psi(\delta).$$

We have

$$\frac{d\psi}{d\delta} = -2 \cos \delta \sin \delta [1 + \sin^2(\boldsymbol{\theta} - \theta^*)] - 2 \sin \delta [\sin(\boldsymbol{\theta} - \theta^*)] (\theta^* - \boldsymbol{\theta}).$$

Hence $d\psi/d\delta|_{\delta=0} = 0$. Further

$$\begin{aligned} \frac{d^2\psi}{d\delta^2}\Big|_{\delta=0} &= -2[-\sin^2\delta + \cos^2\delta][1 + \sin^2(\theta - \theta^*)]\Big|_{\delta=0} - \\ &\quad - 2\cos\delta[\sin(\theta - \theta^*)](\theta^* - \theta)\Big|_{\delta=0} = \\ &= -2[1 + \sin^2(\theta - \theta^*)] + 2(\theta - \theta^*)\sin(\theta - \theta^*). \end{aligned}$$

If $\theta > \theta^*$ then

$$\frac{d^2\psi}{d\delta^2}\Big|_{\delta=0} \leq -2 + 2[\theta - \theta^* - \sin(\theta - \theta^*)].$$

If $\theta < \theta^*$ then

$$\frac{d^2\psi}{d\delta^2}\Big|_{\delta=0} \leq -2 + 2[\theta^* - \theta - \sin(\theta^* - \theta)].$$

Hence, if θ^* is so near to θ that

$$|\theta - \theta^*| - \sin(|\theta - \theta^*|) < 1$$

then $E_\theta[\tilde{\theta}(\mathbf{y}) - \theta]^2$ attains its maximum at $\delta = 0$. Consequently

$$E_\theta[\tilde{\theta}(\mathbf{y}) - \theta]^2 < E_\theta[\tau(\mathbf{y}, \theta^*) - \theta]^2.$$

4. THE POSTERIOR PROBABILITY DENSITY OF θ

Consider a normal prior density $\pi(\theta)$ in model (1),

$$\pi(\theta) = (2\pi)^{-m/2} \det^{-1/2}(\mathbf{H}) \exp\left\{-\frac{1}{2}(\theta - \theta^0)^T \mathbf{H}^{-1}(\theta - \theta^0)\right\},$$

where \mathbf{H} is a given matrix and $\theta^0 \in \Theta$ is a given vector. Denote by $\pi(\theta | \mathbf{y})$ the corresponding posterior density. If $\tau(\mathbf{y})$ is a sufficient static (Corollary 2 to Proposition 1) then

$$\pi(\theta | \mathbf{y}) = \pi(\theta | \tau(\mathbf{y})).$$

This is not a normal density. However, using the linearization described in Section 3, we can write approximately

$$\pi(\theta | \mathbf{y}) \doteq \pi_{\text{lin}}(\theta | \tilde{\theta}(\mathbf{y}))$$

where $\tilde{\theta}(\mathbf{y})$ is supposed to be distributed according to Eq. (14).

Proposition 3. $\pi_{\text{lin}}(\theta | \tilde{\theta}(\mathbf{y}))$ is a normal probability density with the mean equal to

$$(19) \quad \theta^0 + \mathbf{H}(c\mathbf{V} + \mathbf{H})^{-1}(\tilde{\theta}(\mathbf{y}) - \theta^0)$$

and with the variance matrix equal to

$$(20) \quad \mathbf{H} - \mathbf{H}(c\mathbf{V} + \mathbf{H})^{-1}\mathbf{H}$$

where $\tilde{\theta}(\mathbf{y})$ and $\mathbf{V} = \text{Var } \tilde{\theta}(\mathbf{y})$ are defined by Eqs. (12) resp. (13).

Proof. Denote by $h_{\text{lin}}(\tilde{\theta} | \theta)$ the probability density of $\tilde{\theta}$ corresponding to Eq. (14).

Consider the vector

$$\begin{pmatrix} \tilde{\theta} \\ \theta \end{pmatrix}$$

as a random vector with the joint density $\pi(\theta) h_{\text{lin}}(\tilde{\theta} | \theta)$ and denote by $E(\cdot)$ the operator of taking the mean with respect to this density. By simple computations we obtain

$$\begin{aligned} E(\theta) &= \theta^0 \\ E(\tilde{\theta}) &= E[E_{\theta}(\tilde{\theta})] = \theta^0 \\ E[(\theta - \theta^0)(\theta - \theta^0)^T] &= \mathbf{H} \\ E[(\theta - \theta^0)(\tilde{\theta} - \theta^0)^T] &= E[(\theta - \theta^0) E_{\theta}(\tilde{\theta} - \theta^0)^T] = \mathbf{H} \\ E[(\tilde{\theta} - \theta^0)(\tilde{\theta} - \theta^0)^T] &= E[E_{\theta}[(\tilde{\theta} - \theta^0)(\tilde{\theta} - \theta^0)^T]] = c \mathbf{V} + \mathbf{H}. \end{aligned}$$

Hence

$$\begin{pmatrix} \tilde{\theta} \\ \theta \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \theta^0 \\ \theta^0 \end{pmatrix}, \begin{pmatrix} c \mathbf{V} + \mathbf{H}, & \mathbf{H} \\ \mathbf{H}, & \mathbf{H} \end{pmatrix} \right).$$

According to [6], Chapt. 8. a 2, (V), the conditional density of θ given $\tilde{\theta}$ is normal with the mean

$$\theta^0 + \mathbf{H}(c \mathbf{V} + \mathbf{H})^{-1} (\tilde{\theta} - \theta^0)$$

and with the variance

$$\mathbf{H} - \mathbf{H}(c \mathbf{V} + \mathbf{H})^{-1} \mathbf{H}. \quad \square$$

Note. The statistic $\tilde{\theta}(\mathbf{y})$ is sufficient in the linearized model (11). Therefore we can write (compare with Corollary 2)

$$\pi_{\text{lin}}(\theta | \tilde{\theta}(\mathbf{y})) = \pi_{\text{lin}}(\theta | \tau(\mathbf{y})).$$

On the other hand, the exact posterior density is

$$\pi(\theta | \mathbf{y}) = \pi(\theta | \tau(\mathbf{y})).$$

Hence we can compare the approximative and the exact posterior density from

$$\frac{\pi_{\text{lin}}(\theta | \tilde{\theta}(\mathbf{y}))}{\pi(\theta | \mathbf{y})} = \frac{\pi_{\text{lin}}(\theta | \tau(\mathbf{y}))}{\pi(\theta | \tau(\mathbf{y}))}.$$

In $\pi_{\text{lin}}(\theta | \tau(\mathbf{y}))$ we take $E_{\theta}[\tau(\mathbf{y})] = \mathbf{J}\theta$, in $\pi(\theta | \tau(\mathbf{y}))$ we take $E_{\theta}[\tau(\mathbf{y})] = \mathbf{m}_{\theta}$, otherwise the Bayes formulae for computing $\pi_{\text{lin}}(\theta | \tau(\mathbf{y}))$ and $\pi(\theta | \tau(\mathbf{y}))$ are the same.

5. A NOTE ON CONFIDENCE REGIONS FOR θ

We consider confidence regions for θ which are based on $\tilde{\theta}(\mathbf{y})$. We note that they are of restricted importance, since they are influenced by the choice of θ^* in (3), resp. by the choice of the points $\theta^1, \dots, \theta^k$ in (8) which in fact represents a prior knowledge about θ .

From (16) we obtain that the set

$$(21) \quad \{\theta: \|\tilde{\theta}(\mathbf{y}) - \mathbf{Qm}_\theta\|_{\mathbf{V}}^2 < c\chi_m^2(\alpha)\}$$

is a confidence region for θ in the case that $c\mathbf{W}$ is known. α is the exact confidence level, and $\chi_m^2(\alpha)$ is the α -quantile of the χ^2 distribution with m degrees of freedom.

Example 2. Take the set-up from Example 1. We have

$$\mathbf{Qm}_\theta = \mathbb{E}_\theta[\tilde{\theta}(\mathbf{y})] = \cos \delta [-\cos \theta \sin \theta^* + \sin \theta \cos \theta^*] + \theta^*.$$

Hence

$$\|\tilde{\theta}(\mathbf{y}) - \mathbf{Qm}_\theta\|_{\mathbf{V}}^2 = [-(y_1 - \cos \theta) \sin \theta^* + (y_2 - \sin \theta) \cos \theta^*]^2.$$

We see that the confidence region (21) does not depend on δ , hence the standard and the nonstandard linearizations are equivalent as regard to the confidence regions. This is by no way in contradiction to Proposition 2; the random variable $\tilde{\theta}(\mathbf{y})$ has a small variance, however, this has no importance for confidence reasoning. On the other hand, the obtained confidence region depends very much on θ^* .

To understand the situation geometrically, let us write $\tilde{\theta}(\mathbf{y})$ in the form

$$\tilde{\theta}(\mathbf{y}) = \mathbf{L}\mathbf{y} + \mathbf{l}$$

for some matrix \mathbf{L} and some vector \mathbf{l} (This is possible, since $\tilde{\theta}(\mathbf{y})$ is linear in \mathbf{y}). Further we have

$$c\mathbf{V} = \text{Var } \tilde{\theta}(\mathbf{y}) = c\mathbf{LW}^{-1}\mathbf{L}^T$$

hence

$$\mathbf{P} := \mathbf{L}^T\mathbf{V}^{-1}\mathbf{LW}^{-1}$$

is a \mathbf{W} -orthogonal projector. We can verify that

$$(22) \quad \|\mathbf{P}[\mathbf{y} - \boldsymbol{\eta}(\theta)]\|_{\mathbf{W}}^2 = \|\tilde{\theta}(\mathbf{y}) - \mathbf{Qm}_\theta\|_{\mathbf{V}}^2.$$

Hence the confidence region (21) has the form

$$\{\theta: \|\mathbf{P}[\mathbf{y} - \boldsymbol{\eta}(\theta)]\|_{\mathbf{W}}^2 < c\chi_m^2(\alpha)\}.$$

This confidence region, although exact, gives poor results (it is too large) if the value $\|\mathbf{P}[\boldsymbol{\eta}(\theta_{\text{true}}) - \boldsymbol{\eta}(\theta^*)]\|_{\mathbf{W}}^2$ is large. (We note, that this is zero if model (1) is linear.)

Another consequence of (22) is that $\|\tilde{\theta}(\mathbf{y}) - \mathbf{Qm}_\theta\|_{\mathbf{V}}^2$ and $\|(\mathbf{I} - \mathbf{P})[\mathbf{y} - \boldsymbol{\eta}(\theta)]\|_{\mathbf{W}}^2$ are independent random variables. Hence another confidence region (of the exact confidence level α) is of the form

$$\left\{ \theta: \frac{(N-m)\|\tilde{\theta}(\mathbf{y}) - \mathbf{Qm}_\theta\|_{\mathbf{V}}^2}{m\|(\mathbf{I} - \mathbf{P})[\mathbf{y} - \boldsymbol{\eta}(\theta)]\|_{\mathbf{W}}^2} < F_{m, N-m}(\alpha) \right\}$$

where $F_{m, N-m}(\alpha)$ is the α -quantile of the F-distribution with m and $N - m$ degrees of freedom. The advantage of this region comparing with (21) is that it can be used in the case when c is unknown.

6. CONSEQUENCES FOR NONLINEAR EXPERIMENTAL DESIGN

The covariance matrix of $\tilde{\theta}(\mathbf{y})$ (Eq. (16)), and the approximative a posteriori covariance matrix (Eq. (20)) do not depend on the observed vector \mathbf{y} , and are smaller than the corresponding variances in the standard linearization. Therefore they are adequate to construct optimality criteria for optimum experimental design in nonlinear models.

(Received March 28, 1989.)

REFERENCES

- [1] J. Anděl: *Matematická statistika (Mathematical Statistics)*. SNTL/ALFA, Praha 1978.
- [2] D. M. Bates and D. G. Watts: Relative curvature measures of nonlinearity. *J. Roy. Statist. Soc. Ser. B* 42 (1980), 1–25.
- [3] L. Kubáček: *Foundations of Estimation Theory*. Elsevier, Amsterdam 1988. (Slovak edition: *Základy teórie odhadu*. Veda, Bratislava 1982.)
- [4] A. Pázman: Probability distribution of the multivariate nonlinear least squares estimates. *Kybernetika* 20 (1984), 209–230.
- [5] A. Pázman: On information matrices in nonlinear experimental design. *J. Statist. Plann. Inference* 21 (1989), 253–263.
- [6] R. C. Rao: *Linear Statistical Inference and Its Applications*. J. Wiley, New York 1973. (Czech translation: *Lineární metody statistické indukce a jejich aplikace*. Academia, Praha 1978.)
- [7] A. Rényi: *Teorie pravděpodobnosti*. Academia, Praha 1972. (German original: *Wahrscheinlichkeitsrechnung mit einem Anhang über Informationstheorie*. VEB Deutscher Verlag d. Wissenschaften, Berlin 1962.)

Andrej Pázman, Matematický ústav SAV (Mathematical Institute — Slovak Academy of Sciences), Obrancov mieru 49, 814 73 Bratislava. Czechoslovakia.