

# Information geometry of estimating functions in semi-parametric statistical models

SHUN-ICHI AMARI<sup>1,2\*</sup> and MOTOAKI KAWANABE<sup>2</sup>

<sup>1</sup>*Department of Mathematical Engineering and Information Physics, Faculty of Engineering, University of Tokyo, 7-3-1 Hongo, Bunkyo-ku, Tokyo 113, Japan*

<sup>2</sup>*Frontier Research Program, Institute of Physical and Chemical Research (RIKEN), Wako-shi, 351-01 Japan*

For semi-parametric statistical estimation, when an estimating function exists, it often provides an efficient or a good consistent estimator of the parameter of interest against nuisance parameters of infinite dimensions. The present paper elucidates the structure of estimating functions, based on the dual differential geometry of statistical inference and its extension to fibre bundles. The paper studies the following problems. First, when does an estimating function exist and what is the set of all the estimating functions? Second, how are the asymptotic variances of the estimators derived from estimating functions and when are the estimators efficient? Third, how do we adaptively choose a practically good (quasi-)estimating function from the observed data? The concept of  $m$ -curvature freeness plays a fundamental role in solving the above problems.

*Keywords:* dual geometry; dual parallel transport; efficient score function; estimating function; Hilbert fibred structure;  $m$ -curvature free; semi-parametric model

## 1. Introduction

A semi-parametric statistical model treats a family of probability distributions  $\{p(x, \theta, k)\}$  specified by a finite-dimensional parameter  $\theta$  of interest and an infinite-dimensional nuisance parameter  $k$ . Estimation of the parameter of interest in such a model has attracted statisticians for many years, because various important problems are formulated in terms of semi-parametric models. When the nuisance parameter is finite-dimensional, a fundamental role is played by the efficient or projected score function, which is the projection of the score function on the space orthogonal to the score functions of the nuisance parameters. The Cramér–Rao-type inequality has been established in terms of the efficient Fisher information and the bound is asymptotically attainable.

It is not easy to generalize these results to the semi-parametric case. Levit (1978), Begun *et al.* (1983) and Small and McLeish (1988; 1989) defined the efficient Fisher information in

\* To whom correspondence should be addressed.

the semi-parametric model by using the projected score, and showed that the Cramér–Rao-type inequality holds. Only recently has its asymptotic attainability been elucidated under certain regularity conditions (Bickel *et al.* 1993) by various efforts based on functional analysis (for example, Ritov and Bickel 1990; van der Vaart 1991; see also Groeneboom and Wellner 1992; Pfanzagl 1990). It is also known that there exist cases where the bound is not asymptotically attainable (see Hasminskii and Ibragimov 1983; Ritov and Bickel 1990). Estimation procedures are generally very complicated because of the infinite dimensionality of the nuisance parameter. Moreover, as is clearly shown in Bickel *et al.* (1993), the rigorous analytical foundation involves difficulties, which we do not state in the present paper.

The present paper aims to elucidate the differential geometrical structure underlying the semi-parametric model from the point of view of the dual differential geometry of statistical inference (Amari 1985; Barndorff-Nielsen 1986; Murray and Rice 1993). This may complement the analytical and geometrical treatments of Bickel *et al.* (1993). It should be noted that our mathematical treatments are not rigorous in the sense of pure mathematics, because we do not give rigorous regularity conditions for theorems. Also, we do not enter into the rigorous formulation of the dual differential geometry of the relevant function space, since the appropriate pure mathematical background has not yet been developed. The present paper belongs to the field of ‘theoretical mathematics’ in the sense of Jaffe and Quinn (1993). The paper proposes interesting geometrical ideas and some heuristic arguments which lead to useful statistical procedures but are not necessarily mathematically rigorous.

There is a class of estimators obtained by solving a simple equation of the type

$$\sum_{i=1}^n \mathbf{y}(x_i, \boldsymbol{\theta}) = 0, \quad (1.1)$$

where  $x_1, \dots, x_n$  are independent observations from an identical distribution  $p(x, \boldsymbol{\theta}, k)$  of a semi-parametric model. Here, the vector function  $\mathbf{y}(x, \boldsymbol{\theta})$  should satisfy

$$E_{\boldsymbol{\theta}, k}[\mathbf{y}(x, \boldsymbol{\theta})] = 0$$

for all  $k$ , where  $E_{\boldsymbol{\theta}, k}$  denotes the expectation with respect to the distribution specified by  $\boldsymbol{\theta}$  and  $k$ . Such a function  $\mathbf{y}(x, \boldsymbol{\theta})$  which does not depend on the unknown nuisance parameter  $k$  is called an *estimating function* (Godambe 1976) and such an estimator is also called an *M*-estimator (Huber 1981; Bickel *et al.* 1993). It gives a practically tractable method of estimation (see, for example, McLeish and Small 1988; Godambe 1991; Godambe and Heyde 1987). Further regularity conditions will be imposed later. The present paper aims to elucidate estimating functions and their efficiency from the geometrical point of view by introducing Hilbert bundles and their dual parallel transports.

The class of *M*-estimators does not necessarily include efficient estimators. This is one of the points discussed by Pfanzagl (1990), who criticized the method of estimating functions for this reason. On the other hand, it is widely known that the class includes efficient estimators in many practically important cases, and moreover, it gives tractable and robust estimating procedures. Therefore, it is not wise to ignore this class of estimators. On the contrary, it is important to study such fundamental problems as the following:

- (1) When does an estimating function exist?
- (2) What is the set of all the estimating functions?
- (3) What is the best estimator derived from estimating functions?
- (4) When does the best estimating function give an efficient estimator? In the case where the best one is not efficient, what is the amount of loss of information caused by using the estimation function method?
- (5) How to choose a good (quasi-)estimating function based on the observed data?

The present paper treats these problems. Here, a quasi-estimating function implies that it is one obtained from the observed data adaptively. The concept of  $m$ -information-curvature freeness plays a fundamental role.

The method of the present paper is based on informal notes by Amari (1987a) and is an extension of Amari and Kumon (1988). The theory is motivated by the information geometry (Amari 1985; Nagaoka and Amari 1982; Barndorff-Nielsen 1986; Murray and Rice 1993), which studies the structure of the manifold of probability distributions or a statistical model by introducing a Riemannian metric due to the Fisher information and a pair of dual affine connections. It has been proved to be a powerful method in various fields of information science (Amari 1985; 1987b; Amari and Han 1989; Amari and Kumon 1988; Okamoto *et al.* 1991; Amari *et al.* 1992), although the pure mathematical background has been established mostly in finite-dimensional cases. Infinite-dimensional cases are being studied by mathematicians. See Friedrich (1991), Lafferty (1988) and Kanbayashi (1994).

## 2. Semi-parametric statistical models and estimating functions

Let  $p(x, \boldsymbol{\theta}, k)$  be a probability density function of a random variable  $x$  with respect to a common dominating measure  $\mu(dx)$ , specified by two kinds of parameters  $\boldsymbol{\theta} = (\theta^1, \dots, \theta^m)$  and  $k$ , where  $\boldsymbol{\theta} \in \Theta$  is a finite-dimensional vector,  $\Theta$  is an open set of  $\mathbb{R}^m$  and  $k \in K$  is an infinite-dimensional parameter, typically occupying a space of functions. The set of distributions  $S = \{p(x, \boldsymbol{\theta}, k)\}$  is called a semi-parametric statistical model, where  $\boldsymbol{\theta}$  is called the parameter of interest and  $k$  is called the nuisance parameter.

Let  $\mathbf{y}(x, \boldsymbol{\theta}) = [y_i(x, \boldsymbol{\theta}), i = 1, \dots, m]$ , be a vector-valued smooth function of  $\boldsymbol{\theta}$ , not depending on  $k$ , of the same dimension as  $\boldsymbol{\theta}$ . Such a function is called an estimating function when it satisfies the following conditions,

$$E_{\boldsymbol{\theta}, k}[\mathbf{y}(x, \boldsymbol{\theta})] = \mathbf{0}, \quad (2.1)$$

$$\det |E_{\boldsymbol{\theta}, k}[\partial_{\boldsymbol{\theta}} \mathbf{y}(x, \boldsymbol{\theta})]| \neq 0, \quad (2.2)$$

$$E_{\boldsymbol{\theta}, k}[\|\mathbf{y}(x, \boldsymbol{\theta})\|^2] < \infty, \quad E_{\boldsymbol{\theta}, k}[\|\partial_{\boldsymbol{\theta}} \mathbf{y}(x, \boldsymbol{\theta})\|^2] < \infty, \quad (2.3)$$

for all  $\boldsymbol{\theta}$  and  $k$ , where  $E_{\boldsymbol{\theta}, k}$  denotes the expectation with respect to the distribution  $p(x, \boldsymbol{\theta}, k)$ ,  $\partial_{\boldsymbol{\theta}} \mathbf{y}$  is the gradient of  $\mathbf{y}$  with respect to  $\boldsymbol{\theta}$ , i.e., the matrix whose elements are  $(\partial y_i / \partial \theta^j)$  in the component form,  $\det |\cdot|$  denotes the determinant of a matrix, and  $\|\mathbf{y}\|^2$  is the squared norm of the vector  $\mathbf{y}$ ,  $\|\mathbf{y}\|^2 = \sum (y_i)^2$ . We further need that  $\int \mathbf{y} p d\mu$  is differentiable with respect to  $\boldsymbol{\theta}$  (Godambe 1976). Condition (2.1) is essential, as is shown

below. Condition (2.2) guarantees that  $\mathbf{y}$  depends substantially on  $\boldsymbol{\theta}$ . This excludes trivial functions. Condition (2.3) guarantees the applicability of the law of large numbers and the central limit theorem shown below.

When a function satisfying (2.1)–(2.3) exists, by replacing the expectation in (2.1) by the empirical sum, we have an estimator  $\hat{\boldsymbol{\theta}}$  of  $\boldsymbol{\theta}$  by solving

$$\sum_{i=1}^n \mathbf{y}(x_i, \hat{\boldsymbol{\theta}}) = \mathbf{0}, \quad (2.4)$$

where  $x_1, \dots, x_n$  are  $n$  independently and identically distributed observations. This is called the estimating equation and such an estimator is called an  $M$ -estimator.

The asymptotic behaviour of the  $M$ -estimator  $\hat{\boldsymbol{\theta}}$  is obtained from the expansion

$$0 = \sum_{i=1}^n \mathbf{y}(x_i, \hat{\boldsymbol{\theta}}) = \sum_{i=1}^n \mathbf{y}(x_i, \boldsymbol{\theta}) + \sum_{i=1}^n \partial_{\boldsymbol{\theta}} \mathbf{y}(x_i, \boldsymbol{\theta})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) + O_p(\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2), \quad (2.5)$$

by applying the law of large numbers to  $(1/n) \sum \partial_{\boldsymbol{\theta}} \mathbf{y}(x_i, \boldsymbol{\theta})$  and the central limit theorem to  $(1/\sqrt{n}) \sum \mathbf{y}(x_i, \boldsymbol{\theta})$ ,

$$\frac{1}{n} \sum \partial_{\boldsymbol{\theta}} \mathbf{y}(x_i, \boldsymbol{\theta}) \approx A, \quad (2.6)$$

$$\frac{1}{\sqrt{n}} \sum \mathbf{y}(x_i, \boldsymbol{\theta}) \sim \boldsymbol{\varepsilon},$$

where  $A = E_{\boldsymbol{\theta}, k}[\partial_{\boldsymbol{\theta}} \mathbf{y}(x, \boldsymbol{\theta})]$ ,  $\boldsymbol{\varepsilon}$  is a normal random variable subject to  $N(0, V)$  with

$$V = E_{\boldsymbol{\theta}, k}[\mathbf{y}\mathbf{y}^T], \quad (V_{ij} = E_{\boldsymbol{\theta}, k}[y_i y_j]), \quad (2.8)$$

$\mathbf{y}$  is a column vector and  $\mathbf{y}^T$  is its transposition, and  $\approx$  and  $\sim$  denote convergence in probability and in distribution, respectively. From this, by neglecting higher-order terms, we have

$$\hat{\boldsymbol{\theta}} - \boldsymbol{\theta} \sim A^{-1} \boldsymbol{\varepsilon}. \quad (2.9)$$

Condition (2.2) guarantees the existence of  $A^{-1}$  and (2.3) guarantees the existence of  $V$ . The following proposition holds.

**Proposition 1.** *Under the ordinary regularity conditions, the estimator  $\hat{\boldsymbol{\theta}}$  obtained from an estimating function  $\mathbf{y}(x, \boldsymbol{\theta})$  is consistent and asymptotically normally distributed, with the asymptotic covariance matrix*

$$AV[\hat{\boldsymbol{\theta}}; \mathbf{y}] = A^{-1} E_{\boldsymbol{\theta}, k}[\mathbf{y}\mathbf{y}^T] (A^T)^{-1}, \quad (2.10)$$

where the asymptotic covariance matrix is defined by

$$AV[\hat{\boldsymbol{\theta}}; \mathbf{y}] = \lim_{n \rightarrow \infty} n E_{\boldsymbol{\theta}, k}[(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^T]. \quad (2.11)$$

Let  $T(\boldsymbol{\theta})$  be a non-singular  $m \times m$  matrix smoothly depending on  $\boldsymbol{\theta}$ . It should be noted that  $\mathbf{y}^*(x, \boldsymbol{\theta}) = T(\boldsymbol{\theta})\mathbf{y}(x, \boldsymbol{\theta})$  gives an estimating function, equivalent to  $\mathbf{y}$  in the sense of yielding the same estimates.

We give some examples of semi-parametric statistical models for later use. See Bickel *et al.* (1993) and Groeneboom and Wellner (1992) for many interesting models.

**Example 1. Linear dependence.** We consider the following simple but interesting example of linear regression known as one of the Neyman–Scott problems (Neyman and Scott 1948). This example is analysed in Section 7. Let  $x = (\alpha, \beta)$  be composed of two components  $\alpha$  and  $\beta$ . It is supposed that the two variables satisfy the linear relation

$$\beta = \theta\alpha$$

in the ideal case, but their observations are contaminated by normal noises. Let  $\xi^* = (\xi_1, \xi_2, \dots)$  be an unknown infinite sequence, and we assume that the  $i$ th observation  $x_i = (\alpha_i, \beta_i)$  is a noisy version of  $(\xi_i, \theta\xi_i)$  given by

$$\begin{aligned}\alpha_i &= \xi_i + n_i, \\ \beta_i &= \theta\xi_i + n'_i\end{aligned}\tag{2.12}$$

where  $n_i$  and  $n'_i$  ( $i = 1, 2, \dots$ ) are independent and are subject to the normal distribution  $N(0, \sigma^2)$ .

There are many estimators of  $\theta$ . A very simple one is

$$\hat{\theta}' = \frac{\sum \alpha_i \beta_i}{\sum \alpha_i^2},$$

which minimizes the regression error along the  $\beta$ -axis  $\sum (\beta_i - \theta\alpha_i)^2$ . However, this  $\hat{\theta}'$  is not good, and is not even a consistent estimator. Another estimator is

$$\hat{\theta}_\infty = \frac{\sum \beta_i}{\sum \alpha_i},$$

which is the solution of the estimating function

$$y_\infty(x, \theta) = \beta - \theta\alpha.$$

The estimator given by the estimating function

$$y_0(x, \theta) = (\beta - \theta\alpha)(\theta\beta + \alpha)$$

is the  $\theta$  component of the joint maximum likelihood estimator of all the parameters  $(\theta, \xi_1, \xi_2, \dots, \xi_n)$  and is known to be equal to the least-squares estimator with respect to the distance of the observed data to the regression line in the orthogonal direction. Which of the two latter consistent estimators is better? There is no definite answer because it depends on the sequence  $\xi^*$ .

When  $\xi^*$  is considered to consist of independently and identically distributed observations from an unknown distribution  $k(\xi)$ , the present problem can be written in the semi-parametric form

$$p(x, \theta, k) = \int q(x, \theta, \xi)k(\xi)d\xi,$$

where

$$q(x, \theta, \xi) = c \exp \left[ -\frac{1}{2\sigma^2} \{(\alpha - \xi)^2 + (\beta - \theta\xi)^2\} \right].$$

This problem is popular but is still very interesting. We propose later an interesting class of estimators. It is the class of  $c$ -estimators  $\hat{\theta}_c$ , obtained from the following family of estimating functions

$$y_c(x, \theta) = (\beta - \theta\alpha)(\theta\beta + \alpha + c)$$

where  $c$  denotes a real constant. A good estimator is obtained by determining  $c$  adaptively from the observed data.

The problem is a special case of general mixture models discussed in the following.

**Example 2. Mixture models.** Let  $\{q(x, \theta, \xi)\}$  be a regular statistical model, where both the parameter of interest  $\theta$  and the nuisance parameter  $\xi$  are of finite dimensions. Let  $x_i$ ,  $i = 1, 2, \dots, n$ , be  $n$  independent observations from  $q(x, \theta, \xi_i)$ , where  $\theta$  is common but  $\xi_i$  takes a different value at each observation. Moreover, we assume that the unknown  $\xi_i$  are independently generated subject to a common but unknown probability distribution having a density function  $k(\xi)$ . Then, the  $x_i$  are regarded as independent observations from the semi-parametric model

$$p(x, \theta, k) = \int q(x, \theta, \xi)k(\xi)d\xi, \quad (2.13)$$

where  $k(\xi)$  is the infinite-dimensional nuisance parameter. This model is called the mixture model. This type of problem was studied by Neyman and Scott (1948) and has attracted many researchers (among them Andersen 1970; Lindsay 1982; Kumon and Amari 1984; Amari and Kumon 1988; Pfanzagl 1990). There are a lot of interesting and important examples in this class. Most researchers have treated the distributions of the following exponential form as examples,

$$q(x, \theta, \xi) = \exp\{\xi \cdot \mathbf{s}(x, \theta) + r(x, \theta) - \psi(\theta, \xi)\}, \quad (2.14)$$

where  $\mathbf{s}(x, \theta)$  is a vector not depending on  $\xi$  and  $\cdot$  is the inner product. Here, the distribution is of exponential type for  $\xi$  when  $\theta$  is fixed. Models of this type have the  $m$ -flat nuisance structure (Amari 1987a) or the convex structure (Bickel *et al.* 1993), to be explained later, so that they possess nice properties. The linear dependence model (Example 1) is a special case of this type.

**Example 3. Semi-parametric additive regression.** Let  $(\mathbf{x}, t)$  be covariates to which  $y$  is connected by

$$y = \theta \cdot \mathbf{x} + k(t) + \epsilon, \quad (2.15)$$

where  $\theta$  is the parameter of interest,  $k(t)$  is an unknown nuisance smooth function of  $t$ , and  $\epsilon$  is a noise term. We assume here that  $\epsilon$  is subject to  $N(0, 1)$ . The problem is to estimate  $\theta$  based on  $n$  observations  $(y_i, \mathbf{x}_i, t_i)$ ,  $i = 1, 2, \dots, n$ , where we assume that  $(\mathbf{x}, t)$  is subject to a

known joint distribution  $q(\mathbf{x}, t)$  (see Cuzick 1992). This model is analysed later to show that our theory is applicable also to non- $m$ -flat models.

### 3. Hilbert fibre space and score function

Given a probability density function  $p(x)$ , its small deviation in the direction of  $a(x)$  can be represented by a curve  $p(x, t)$  starting from  $p(x)$ ,

$$p(x, t) = p(x)\{1 + ta(x)\}, \quad (3.1)$$

where  $t(0 \leq t < \epsilon)$  is the parameter of the curve. Here

$$E[a(x)] = 0$$

holds, where  $E$  is the expectation with respect to  $p(x)$ , because of

$$\int p(x)\{1 + ta(x)\}d\mu(x) = 1. \quad (3.2)$$

In order to be specific, we consider the linear space of functions which satisfy

$$E[a(x)] = 0, \quad E[\{a(x)\}^2] < \infty. \quad (3.3)$$

This set is a Hilbert space  $H_p$ , often denoted by  $L_2^0(p)$ , with the inner product of  $a(x)$  and  $b(x)$  defined by

$$\langle a(x), b(x) \rangle = E[a(x)b(x)]. \quad (3.4)$$

We call the random variable

$$a(x) = \left. \frac{d}{dt} \log p(x, t) \right|_{t=0} \quad (3.5)$$

the tangent vector of the curve (3.1). This is the score function for the one-dimensional statistical model (3.1) parametrized by  $t$ . Refer to Pfanzagl (1990), van der Vaart (1991), Groeneboom and Wellner (1992) and Bickel *et al.* (1993) for mathematical details on the rigorous construction of the tangent space.

Given a semi-parametric model  $S = \{p(x, \boldsymbol{\theta}, k)\}$ , the Hilbert space  $H_p = H_{\boldsymbol{\theta}, k}$  is associated with each point  $(\boldsymbol{\theta}, k)$ , that is, with each distribution  $p(x) = p(x, \boldsymbol{\theta}, k)$  specified by  $(\boldsymbol{\theta}, k)$ . A collection of such  $H_{\boldsymbol{\theta}, k}$  is called a *fibred structure*, where the fibres are the Hilbert spaces.

We first define the tangent directions along the parameter of interest. Let

$$u_i(x, \boldsymbol{\theta}, k) = \frac{\partial}{\partial \theta^i} \log p(x, \boldsymbol{\theta}, k) \quad (3.6)$$

be the score function with respect to the  $i$ th component  $\theta^i$  of  $\boldsymbol{\theta}$ . Obviously,

$$E_{\boldsymbol{\theta}, k}[u_i] = 0 \quad (3.7)$$

and we further assume that  $u_i$  is square-integrable. Then it belongs to  $H_{\theta,k}$ . We call the subspace spanned by these  $u_i$  the tangent subspace  $T_{\theta,k}^I$  along the parameter of interest. The vector score function is  $\mathbf{u} = (u_1, \dots, u_m)$ .

We next define the tangent directions along the nuisance parameter. Let us assume that, for any  $\tilde{k}$  in a small neighbourhood of  $k$  in the set  $K$  of the nuisance parameter, there exists a curve  $c(t)$  connecting them, such that  $c(0) = k$  and  $c(\epsilon) = \tilde{k}$ , and that the score function for the one-dimensional statistical model  $p\{x, \theta, c(t)\}$  parametrized by  $t$ ,

$$v(x, \theta, k, c) = \left. \frac{d}{dt} \log p\{x, \theta, c(t)\} \right|_{t=0} \quad (3.8)$$

belongs to  $H_{\theta,k}$ . This  $v$  is the tangent vector along  $c(t)$  of the nuisance parameter. Let  $T_{\theta,k}^N$  be the smallest closed subspace including all such  $v$ s. We call it the nuisance tangent space.

Now, let us project the score function  $u_i$  onto the subspace  $(T_{\theta,k}^N)^\perp$  which is the orthogonal complement of  $T_{\theta,k}^N$ . The result is the function  $u_i^E = u_i - v$  that minimizes  $E[|u_i - v|^2]$ ,  $v \in T_{\theta,k}^N$ . The vector function  $\mathbf{u}^E = (u_i^E)$  is called the efficient score function and the  $u_i^E$  are called the components of the efficient or projected score functions (see Begun *et al.* 1983; Amari and Kumon 1988; Small and McLeish 1989). Let  $T_{\theta,k}^E$  be the subspace of  $H_{\theta,k}$  spanned by the components  $u_i^E$  of the efficient score function.

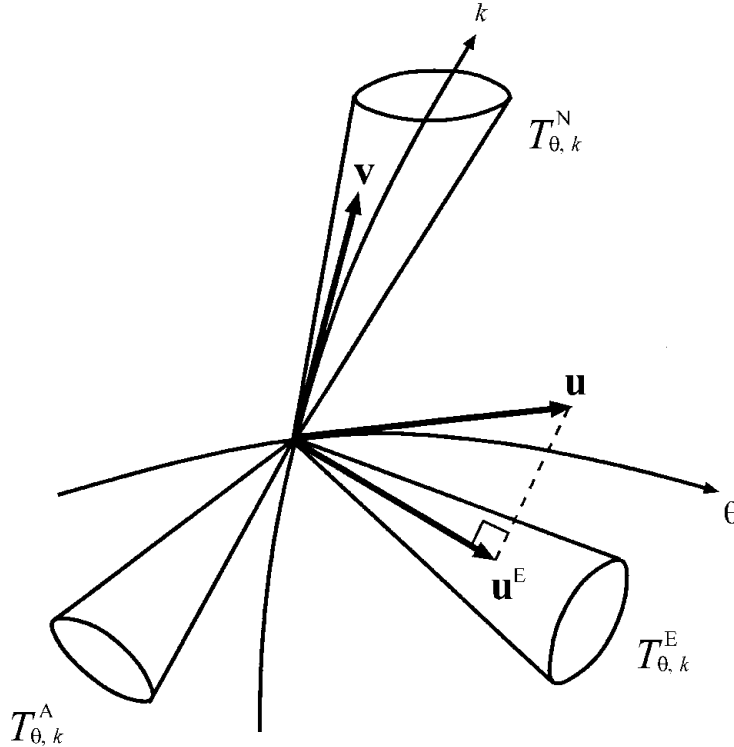


Figure 1. Orthogonal decomposition of  $H_{\theta,k}$



Let  $T_{\theta,k}^A$  be the orthogonal complement of  $T_{\theta,k}^N \oplus T_{\theta,k}^E$ . This is called the ancillary subspace and spans directions orthogonal to any changes in the parameter of interest and the nuisance parameter. We thus have the orthogonal decomposition of the Hilbert fibre space (Fig. 1; see Amari 1987a; Amari and Kumon 1988; see also Small and McLeish 1988),

$$H_{\theta,k} = T_{\theta,k}^E \oplus T_{\theta,k}^A \oplus T_{\theta,k}^N. \quad (3.9)$$

The matrix  $G^E = (g_{ij}^E)$  defined by using the efficient score function

$$g_{ij}^E(\theta, k) = E_{\theta,k}[u_i^E u_j^E] \quad (3.10)$$

is called the efficient Fisher information matrix. Begun *et al.* (1983) proved that  $G^E$  gives the Cramér–Rao bound of the asymptotic covariance of estimators  $\hat{\theta}$ ,

$$\lim_{n \rightarrow \infty} nE[(\hat{\theta} - \theta)(\hat{\theta} - \theta)^T] \geq (G^E)^{-1} \quad (3.11)$$

for any asymptotically normally distributed unbiased estimator in a semi-parametric model. Ritov and Bickel (1990) showed that the bounds may be unattainable in some cases. However, the above bounds are attainable in many important cases under mild regularity conditions (Bickel *et al.* 1993).

## 4. Invariant decomposition of Hilbert fibres due to dual parallel transports

An estimating function  $\mathbf{y}(x, \theta)$  satisfies the unbiasedness condition (2.1) for all  $k$ . Such a global structure is elucidated by introducing two parallel transports of the Hilbert fibres along the nuisance space.

Let  $a(x)$  be a random variable belonging to  $H_{\theta,k}$ . Let us fix  $\theta$ , and consider the subset  $S_{\theta} = \{p(x, \theta, k) | k \in K\}$ . We define two parallel transports of a vector  $a(x)$  from  $H_{\theta,k}$  to  $H_{\theta,k'}$  (Amari 1987a). Then

$$\prod_k^{(e)} a(x) = a(x) - E_{\theta,k'}[a(x)], \quad (4.1)$$

$$\prod_k^{(m)} a(x) = \frac{p(x, \theta, k)}{p(x, \theta, k')} a(x) \quad (4.2)$$

are called the  $e$ -parallel transport and the  $m$ -parallel transport of  $a(x)$  from  $(\theta, k)$  to  $(\theta, k')$ , respectively. It should be noted that the  $e$ -parallel transport exists only when the expectation of  $a(x)$  at  $(\theta, k')$  exists. It is easy to show that

$$E_{\theta,k'} \left[ \prod_k^{(m)} a(x) \right] = 0$$

always holds, and that

$$E_{\theta, k'} \left[ \prod_k^{(e) k'} a(x) \right] = 0$$

holds when  $E_{\theta, k'}[a(x)]$  exists. However, the  $e$ - and/or  $m$ -parallel transports of  $a(x)$  do not necessarily belong to  $H_{\theta, k'}$ . They belong to  $H_{\theta, k'}$  only when they are square-integrable at  $(\theta, k')$  with respect to  $p(x, \theta, k')$ .

The parallel transports are generalizations of the dual geometrical structures derived from the underlying  $e$ - and  $m$ -connections or  $e$ - and  $m$ -covariant derivatives (Amari 1985; see also Amari and Kumon 1988), but we will not go into mathematical details of differential geometry.

The following lemma shows the most important property connecting the two parallel transports. The proof is immediate and hence omitted.

**Lemma 1.** *The two parallel transports are dual in the sense that, for any two  $a(x), b(x) \in H_{\theta, k}$ , the inner product*

$$\langle a, b \rangle_{\theta, k} = \left\langle \prod_k^{(e) k'} a, \prod_k^{(m) k'} b \right\rangle_{\theta, k'}, \quad (4.3)$$

where the suffix  $(\theta, k)$  denotes that the expectation is taken with respect to  $p(x, \theta, k)$ , is kept invariant when their parallel transports belong to  $H_{\theta, k'}$ .

It is remarked that an estimating function is  $e$ -invariant,

$$\prod_k^{(e) k'} \mathbf{y}(x, \theta) = \mathbf{y}(x, \theta)$$

because of (2.1) where  $\prod^{(e)}$  operates componentwise. Let us consider a curve  $k = k(t)$ ,  $k_0 = k(0)$ , in the nuisance space. By differentiating (2.1) with respect to  $t$  along the curve  $k = k(t)$  and exchanging the integral and differentiation, we have

$$\begin{aligned} & \frac{d}{dt} \int p\{x, \theta, k(t)\} \mathbf{y}(x, \theta) d\mu(x) \Big|_{t=0} \\ &= \int v\{x, \theta, k_0\} p\{x, \theta, k_0\} \mathbf{y}(x, \theta) d\mu(x) \\ &= \langle v, \mathbf{y}(x, \theta) \rangle_{\theta, k_0} = 0, \end{aligned}$$

where

$$v = \frac{d}{dt} \log p\{x, \theta, k(t)\} \Big|_{t=0}$$

is the nuisance tangent direction at  $k_0$ . This holds for any  $k_0$  so that any estimating function  $\mathbf{y}(x, \boldsymbol{\theta})$  is orthogonal to  $v$  at any point  $(\boldsymbol{\theta}, k)$ . However, from

$$\begin{aligned} \langle v, \mathbf{y}(x, \boldsymbol{\theta}) \rangle_{\boldsymbol{\theta}, k'} &= \left\langle \prod_{k'}^{(m)} v, \prod_{k'}^{(e)} \mathbf{y} \right\rangle_{\boldsymbol{\theta}, k} \\ &= \left\langle \prod_{k'}^{(m)} v, \mathbf{y} \right\rangle_{\boldsymbol{\theta}, k}, \end{aligned} \quad (4.4)$$

where  $v$  is a nuisance tangent direction at  $k'$ , the orthogonality condition at  $k'$  is transferred to that at  $k$  by the  $m$ -parallel transports of  $v$  at  $k'$ . This shows that an estimating function  $\mathbf{y}$  is orthogonal, not only to the nuisance tangent direction at any  $k$ , but to the  $m$ -parallel transports from  $k'$  to  $k$  of the nuisance tangent directions at any  $k'$ .

Motivated by the above discussion, we now reorganize the decomposition (3.9) of  $H_{\boldsymbol{\theta}, k}$  by taking account of the global structure induced by the parallel transports. The information fibre space  $F_{\boldsymbol{\theta}, k}^I$  at  $(\boldsymbol{\theta}, k)$  is constructed from  $T_{\boldsymbol{\theta}, k}^E$  such that its elements are orthogonal not only to the nuisance tangent space  $T_{\boldsymbol{\theta}, k}^N$  at  $(\boldsymbol{\theta}, k)$  but also to the  $m$ -parallel transports from  $k'$  to  $k$  of  $T_{\boldsymbol{\theta}, k'}^N$  at points  $(\boldsymbol{\theta}, k')$  for all  $k' \in K$ . To this end, we first consider a vector

$$r(x) \in T_{\boldsymbol{\theta}, k}^E \oplus T_{\boldsymbol{\theta}, k}^A$$

whose  $e$ -transport exists in the Hilbert space  $H_{\boldsymbol{\theta}, k'}$  for every  $k'$ , that is,

$$\mathbb{E}_{\boldsymbol{\theta}, k'}[\{r(x)\}^2] < \infty. \quad (4.5)$$

Suppose its  $e$ -transport to  $(\boldsymbol{\theta}, k')$  is orthogonal to  $T_{\boldsymbol{\theta}, k'}^N$  for every  $k' \in K$ , that is,

$$\left\langle v, \prod_{k'}^{(e)} r(x) \right\rangle_{\boldsymbol{\theta}, k'} = \left\langle \prod_{k'}^{(m)} v, r(x) \right\rangle_{\boldsymbol{\theta}, k} = 0, \quad v \in T_{\boldsymbol{\theta}, k'}^N. \quad (4.6)$$

We express this by saying that  $r(x)$  is free of any nuisance tangent directions at every  $k'$  when it is  $e$ -transported. In this case  $r(x)$  is a candidate for an estimating function. This implies that  $r(x)$  is orthogonal to not only  $T_{\boldsymbol{\theta}, k}^N$  but also all the  $m$ -transports from  $k'$  to  $k$  of  $T_{\boldsymbol{\theta}, k'}^N$ . This suggests that we enlarge the nuisance space  $T_{\boldsymbol{\theta}, k}^N$  to

$$\text{span} \left\{ \bigcup_{k'} \prod_{k'}^{(m)} T_{\boldsymbol{\theta}, k'}^N \right\},$$

and define the directions orthogonal to it. Any estimating functions should be orthogonal to the enlarged subspace.

To define such a space formally, we consider the closed subspace of  $H_{\boldsymbol{\theta}, k}$  consisting of the vectors satisfying the above conditions (4.5) and (4.6) and refer to it tentatively by  $F_{\boldsymbol{\theta}, k}^{IA}$ ,

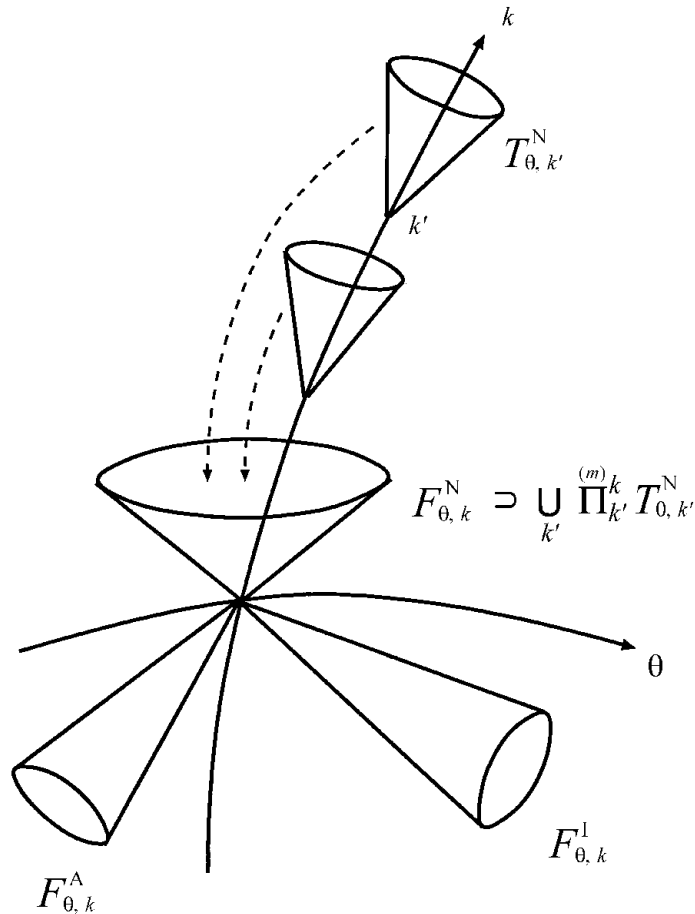
where I denotes the information part and A denotes the ancillary part. Obviously, the following hold:

$$F_{\theta,k}^{\text{IA}} \subset T_{\theta,k}^{\text{E}} \oplus T_{\theta,k}^{\text{A}},$$

$$F_{\theta,k}^{\text{IA}} \perp \bigcup_{k'} \prod_{k'}^{(m)} T_{\theta,k'}^{\text{N}}.$$

The nuisance fibre space  $F_{\theta,k}^{\text{N}}$  is defined as the orthogonal complement of  $F_{\theta,k}^{\text{IA}}$  in  $H_{\theta,k}$ . Obviously

$$F_{\theta,k}^{\text{N}} \supset \prod_{k'}^{(m)} T_{\theta,k'}^{\text{N}}. \quad (4.7)$$



**Figure 2.** Alternative orthogonal decomposition of  $H_{\theta,k}$  based on  $m$ -parallel transports

We have

$$H_{\theta,k} = F_{\theta,k}^{IA} \oplus F_{\theta,k}^N. \quad (4.8)$$

We now decompose  $F_{\theta,k}^{IA}$ . Let  $u_i^I$  be the projection of the score  $u_i$  onto  $F_{\theta,k}^{IA}$ . This is called the information score. The information fibre space denoted by  $F_{\theta,k}^I$  is the subspace spanned by the information score functions  $u_i^I(x, \theta, k)$ . The ancillary fibre space is the orthogonal complement of  $F_{\theta,k}^I$  in  $F_{\theta,k}^{IA}$  and is denoted by  $F_{\theta,k}^A$ ,

$$F_{\theta,k}^{IA} = F_{\theta,k}^I \oplus F_{\theta,k}^A.$$

We thus have another orthogonal decomposition of  $H_{\theta,k}$ ,

$$H_{\theta,k} = F_{\theta,k}^I \oplus F_{\theta,k}^A \oplus F_{\theta,k}^N, \quad (4.9)$$

which represents a more global structure of  $H_{\theta,k}$  (Fig. 2). The subspace  $F_{\theta,k}^N$  includes all the  $m$ -parallel transports of  $T_{\theta,k'}^N$  from  $k'$  to  $k$  because of the relation (4.6). It is the information fibre  $F_{\theta,k}^I$  that plays an important role.

It is helpful to compare the two decompositions (3.9) and (4.9).  $F_{\theta,k}^N$  is enlarged from  $T_{\theta,k}^N$  to include all the  $m$ -parallel transports of  $T_{\theta,k'}^N$ . The  $F_{\theta,k}^I$  is constructed from the projections of the  $\theta$ -score functions to the orthogonal complement of the enlarged  $F_{\theta,k}^N$ , while  $T_{\theta,k}^E$  is constructed from the projections to the orthogonal complement of smaller  $T_{\theta,k}^N$ . Therefore, when  $T_{\theta,k}^N = F_{\theta,k}^N$  holds,  $F_{\theta,k}^I = T_{\theta,k}^E$  holds. This is an important special case.

## 5. Estimating functions and their efficiency

Based on the decomposition (4.9) of the Hilbert space  $H_{\theta,k}$ , we can now characterize the set of all the estimating functions. We first answer the two important questions when an estimating function exists and what is the set of all the estimating functions. To this end, we prove the following two important lemmas.

**Lemma 2.** *Any component of an estimating function  $\mathbf{y}(x, \theta)$  belongs to  $F_{\theta,k}^I \oplus F_{\theta,k}^A$  for any  $k$ . Let  $y_i^I(x, \theta, k)$  be the projection of the  $i$ th component of  $\mathbf{y}(x, \theta)$  onto  $F_{\theta,k}^I$ . Then, the  $y_i^I(x, \theta, k), i = 1, \dots, m$ , span  $F_{\theta,k}^I$ .*

**Proof.** Let  $\mathbf{y}(x, \theta)$  be an estimating function. Then, its  $e$ -transport always exists in the corresponding Hilbert space because of (2.3), and it is  $e$ -invariant:

$$\prod^{(e)} \mathbf{y}(x, \theta) = \mathbf{y}(x, \theta). \quad (5.1)$$

We have already shown in (4.4) that

$$\left\langle \prod_{k'}^{(m)} v(x, \theta, k'), \mathbf{y}(x, \theta) \right\rangle_{\theta,k} = 0 \quad (5.2)$$

for all  $k$  and  $k'$ . Therefore,  $\mathbf{y}$  is included in  $F_{\theta,k}^1 \oplus F_{\theta,k}^A$ . Moreover, by differentiating (2.1) with respect to  $\theta$ , we have

$$E_{\theta,k}[\partial_{\theta}\mathbf{y}(x, \theta)] + \langle \mathbf{u}, \mathbf{y}(x, \theta) \rangle = 0, \quad (5.3)$$

where  $\langle \mathbf{u}, \mathbf{y} \rangle$  denotes a matrix whose elements are  $\langle u_i, y_j \rangle$ . This shows that

$$E_{\theta,k}[\partial_{\theta}\mathbf{y}] = -\langle \mathbf{u}, \mathbf{y} \rangle.$$

Since  $\mathbf{y}$  belongs to  $F_{\theta,k}^1 \oplus F_{\theta,k}^A$  and the projection of  $\mathbf{u}$  onto the space  $F_{\theta,k}^1 \oplus F_{\theta,k}^A$  includes no  $F_{\theta,k}^A$  part, we have

$$\langle \mathbf{u}, \mathbf{y} \rangle = \langle \mathbf{u}^1, \mathbf{y} \rangle,$$

where  $\mathbf{u}^1$  is the projection of  $\mathbf{u}$  onto  $F_{\theta,k}^1$ . Therefore, (2.2) implies that the determinant of  $\langle u_i^1, y_j \rangle$  does not vanish and that the projections of vectors  $y_i$  onto  $F_{\theta,k}^1$  span  $F_{\theta,k}^1$ . This also shows that  $F_{\theta,k}^1$  is non-degenerate, that is, its dimension is the same as that of the parameter  $\theta$  of interest.  $\square$

We now prove what is in essence a converse of Lemma 2.

**Lemma 3.** *Any vector  $\mathbf{w}(x, \theta)$  belonging to  $F_{\theta,k}^1 \oplus F_{\theta,k}^A$  for some  $k \in K$  is an estimating function provided  $F_{\theta,k'}^1$  is non-degenerate for every  $k'$  and the projections of the components  $w_i$  of  $\mathbf{w}$  onto  $F_{\theta,k'}^1$  span  $F_{\theta,k'}^1$  for every  $k'$ .*

**Remark.** The first part of the proof, for which the conditions of the lemma are not needed, will show that the space  $F_{\theta,k}^1 \oplus F_{\theta,k}^A$  is  $e$ -invariant.

**Proof.** Let  $c(t)$  be a curve connecting two points  $k$  and  $k'$ ,  $k = c(0)$ , and put

$$\mathbf{f}(t) = E_{\theta,c(t)}[\mathbf{w}(x, \theta)].$$

Since the  $e$ -parallel transport of  $w$  from  $k$  to  $c(t)$  is written as

$$\prod_k^{(e)} c(t) \mathbf{w} = \mathbf{w} - \mathbf{f}(t),$$

$\mathbf{w}$  is  $e$ -invariant if  $\mathbf{f}(t) = 0$  holds. Hence, we need to prove  $\mathbf{f}(t) = 0$ . Obviously  $\mathbf{f}(0) = 0$ . By differentiation, we have

$$\begin{aligned} \frac{d}{dt}\mathbf{f}(t) &= \int \frac{d}{dt} p\{x, \theta, c(t)\} \mathbf{w}(x, \theta) d\mu(x) \\ &= E_{\theta,c(t)}[v(t)\mathbf{w}(x, \theta)] = \langle v, \mathbf{w} \rangle_{\theta,c(t)}, \end{aligned}$$

where

$$v(t) = \frac{d}{dt} \log p\{x, \theta, c(t)\}.$$

Since  $\mathbf{w} \in F_{\boldsymbol{\theta},k}^{\text{IA}}$ , its  $e$ -parallel transport is orthogonal to  $v(t)$ , and hence we have

$$\begin{aligned} 0 &= \left\langle v(t), \prod_k^{(e)}{}^{c(t)} \mathbf{w} \right\rangle_{\boldsymbol{\theta},c(t)} = \langle v(t), \mathbf{w} \rangle_{\boldsymbol{\theta},c(t)} - \mathbf{f}(t) \mathbf{E}_{\boldsymbol{\theta},c(t)}[v(t)], \\ &= \langle v(t), \mathbf{w} \rangle_{\boldsymbol{\theta},c(t)}. \end{aligned}$$

Hence  $d\mathbf{f}(t)/dt = 0$ , so that  $\mathbf{f}(t) = 0$ . This proves the  $e$ -invariance of  $\mathbf{w}$ ,

$$\prod_k^{(e)}{}^{k'} \mathbf{w} = \mathbf{w}$$

and

$$\mathbf{E}_{\boldsymbol{\theta},k'}[\mathbf{w}(x, \boldsymbol{\theta})] = 0$$

for all  $k'$ . Since the projections  $w_i^1$  of the components  $w_i$  of  $\mathbf{w}$  onto  $F_{\boldsymbol{\theta},k'}^1$  do, by assumption, span  $F_{\boldsymbol{\theta},k'}^1$  for any  $k'$ , then  $\mathbf{E}_{\boldsymbol{\theta},k'}[\partial_{\boldsymbol{\theta}} \mathbf{w}(x, \boldsymbol{\theta})] = -\langle \mathbf{u}^1, \mathbf{w} \rangle_{\boldsymbol{\theta},k'} = -\langle \mathbf{u}^1, \mathbf{w}^1 \rangle_{\boldsymbol{\theta},k'}$  is of full rank. Therefore,  $\mathbf{w}$  is an estimating function satisfying (2.1)–(2.3).  $\square$

By combining Lemmas 2 and 3, we have the following proposition.

**Proposition 2.** *Any estimating function  $\mathbf{y}(x, \boldsymbol{\theta}) = \{y_i(x, \boldsymbol{\theta})\}$  can be decomposed for any  $k$  as a sum*

$$\mathbf{y}(x, \boldsymbol{\theta}) = T(\boldsymbol{\theta}, k) \mathbf{u}^1(x, \boldsymbol{\theta}, k) + \mathbf{a}(x, \boldsymbol{\theta}, k), \quad (5.4)$$

where the component  $a_i(x, \boldsymbol{\theta}, k)$  of  $\mathbf{a}$  belongs to  $F_{\boldsymbol{\theta},k}^{\Delta}$  and  $T(\boldsymbol{\theta}, k)$  is a non-singular matrix. Conversely, any function  $\mathbf{y}(x, \boldsymbol{\theta})$  defined in the form of (5.4) at a fixed  $k_0$  gives an estimating function provided the projections of the components  $y_i(x, \boldsymbol{\theta})$  onto  $F_{\boldsymbol{\theta},k'}^1$  span  $F_{\boldsymbol{\theta},k'}^1$  for every  $k'$ .

It is possible to choose a basis for the efficient scores such that  $T(\boldsymbol{\theta}, k)$  becomes the identity at some  $k$ . The proposition also gives a simple sufficient condition for the existence of an estimating function. This is shown by putting  $\mathbf{a}(x, \boldsymbol{\theta}, k) = 0$  in Proposition 2. Lemma 2 gives a necessary condition. We summarize these.

**Proposition 3.** *A necessary condition for the existence of an estimating function is that  $F_{\boldsymbol{\theta},k}^1$  is non-degenerate, that is,  $m$ -dimensional. A sufficient condition is that, for a fixed  $k_0$ , the projections of  $u_i^1(x, \boldsymbol{\theta}, k_0)$  ( $i = 1, \dots, m$ ) onto  $F_{\boldsymbol{\theta},k'}^1$  span  $F_{\boldsymbol{\theta},k'}^1$  for every  $k'$ . The vector  $\mathbf{u}^1(x, \boldsymbol{\theta}, k_0)$  is an estimating function in this case.*

**Remark.** We discuss informally a necessary and sufficient condition. It is clear that any estimating function can be written in the form

$$y_i(x, \boldsymbol{\theta}) = \int \sum_{j=1}^m \varphi_{ij}(k) u_j^1(x, \boldsymbol{\theta}, k) d\mu(k) + a_i(x, \boldsymbol{\theta}),$$

where  $a_i(x, \boldsymbol{\theta})$  is purely ancillary, that is orthogonal to  $F_{\boldsymbol{\theta},k}^1$  at any  $k$ . Therefore, a necessary and sufficient condition for the existence of estimating functions is the existence of

functionals  $\varphi_{ij}(k)$  such that the projections  $y_i^I(x, \boldsymbol{\theta}, k)$  of the above  $y_i$  onto  $F_{\boldsymbol{\theta}, k}^I$  ( $i = 1, \dots, m$ ) span  $y_i^I(x, \boldsymbol{\theta}, k)$  for all  $k$ . Let us define

$$\rho_{ij}(k, k') = \langle u_i^I(x, \boldsymbol{\theta}, k), u_j^I(x, \boldsymbol{\theta}, k') \rangle_{\boldsymbol{\theta}, k'}.$$

Then, the condition is equivalent to the existence of  $\varphi_{ij}(k)$  such that the matrix

$$B_{ij} = \sum_{n=1}^m \int \varphi_{in}(k) \rho_{nj}(k, k') d\mu(k)$$

is non-degenerate for all  $k'$ . The sufficient condition of Proposition 3 guarantees the existence of such  $\varphi_{ij}$ . When  $u_j^I(x, \boldsymbol{\theta}, k')$  ( $i = 1, \dots, m$ ) are non-degenerate, because of continuity, there exists a neighbourhood  $N_K(k)$  of  $k$  such that the projections of  $u_i^I(x, \boldsymbol{\theta}, k)$ , ( $i = 1, \dots, m$ ) onto  $F_{\boldsymbol{\theta}, k'}$  span  $F_{\boldsymbol{\theta}, k'}$  for all elements  $k'$  of  $N_K(k)$ . If a function  $\mathbf{y}(x, \boldsymbol{\theta})$  satisfies the conditions of estimating functions where the global condition (2.2) is replaced by the local one,

$$\det |E_{\boldsymbol{\theta}, k'}[\partial_{\boldsymbol{\theta}} \mathbf{y}(x, \boldsymbol{\theta})]| \neq 0,$$

at least in a neighbourhood of a point  $k$ , such a function is called a local estimating function. The non-degeneracy of  $\mathbf{u}^I(x, \boldsymbol{\theta}, k)$  is a necessary and sufficient condition of the existence of local estimating functions at  $k$ .

We next calculate the asymptotic covariance matrix of an estimating function. This calculation leads us to the optimal estimating function. As discussed above, an estimating function  $\mathbf{y}(x, \boldsymbol{\theta})$  may be decomposed as

$$\mathbf{y}(x, \boldsymbol{\theta}) = \mathbf{u}^I(x, \boldsymbol{\theta}, k) + \mathbf{a}(x, \boldsymbol{\theta}, k),$$

where  $\mathbf{u}^I = (u_i^I) \in F_{\boldsymbol{\theta}, k}^I$  and  $\mathbf{a} = (a_i) \in F_{\boldsymbol{\theta}, k}^A$ . It is easy to show that

$$E[\partial_{\boldsymbol{\theta}} \mathbf{a}] = -\langle \mathbf{u}, \mathbf{a} \rangle = 0,$$

by differentiating  $E_{\boldsymbol{\theta}, k}[\mathbf{a}(x, \boldsymbol{\theta}, k)] = 0$ . Therefore, we have

$$-E[\partial_{\boldsymbol{\theta}} \mathbf{y}] = -E[\partial_{\boldsymbol{\theta}} \mathbf{u}^I] = \langle \mathbf{u}, \mathbf{u}^I \rangle = \langle \mathbf{u}^I, \mathbf{u}^I \rangle, \quad (5.5)$$

$$E[\mathbf{y}\mathbf{y}^T] = E[\mathbf{u}^I(\mathbf{u}^I)^T] + E[\mathbf{a}\mathbf{a}^T] = G^I + G^A, \quad (5.6)$$

where we put

$$G^I = E[\mathbf{u}^I(\mathbf{u}^I)^T], \quad G^A = E[\mathbf{a}\mathbf{a}^T]. \quad (5.7)$$

So, by Proposition 1, we have the following result.

**Proposition 4.** *The asymptotic covariance matrix derived from an estimating function  $\mathbf{y}(x, \boldsymbol{\theta})$  is given by*

$$\text{AV}[\hat{\boldsymbol{\theta}}; \mathbf{y}] = (G^I)^{-1} + (G^I)^{-1} G^A (G^I)^{-1}. \quad (5.8)$$

The estimating function given by

$$\mathbf{y}(x, \boldsymbol{\theta}) = \mathbf{u}^I(x, \boldsymbol{\theta}, k_0) \quad (5.9)$$



where  $k_0$  is a fixed point in  $K$ , is optimal among all the  $M$ -estimators when the true distribution happens to be specified by  $(\boldsymbol{\theta}, k_0)$  and  $\mathbf{u}_i^1$  span  $F_{\boldsymbol{\theta}, k}^1$  for every  $k$ . The optimal asymptotic covariance is given by  $(G^1)^{-1}$  in this case.

It is shown that the information fibre  $F_{\boldsymbol{\theta}, k}^1$  plays the fundamental role, being defined through the  $m$ -parallel transports of  $T_{\boldsymbol{\theta}, k'}^N$ . However, in many important cases,  $F_{\boldsymbol{\theta}, k}^1$  is equal to the simpler  $T_{\boldsymbol{\theta}, k}^E$ , which is easier to obtain. We show when such a simplification occurs.

We fix  $\boldsymbol{\theta}$ , and consider the statistical submodel  $S_{\boldsymbol{\theta}} = \{p(x, \boldsymbol{\theta}, k)\}$ , where  $k \in K$  is the only free parameter. The tangent vectors of  $S_{\boldsymbol{\theta}}$  compose the nuisance tangent space  $T_{\boldsymbol{\theta}, k}^N$ . Let us consider the  $m$ -parallel transport of  $T_{\boldsymbol{\theta}, k'}^N$  from  $(\boldsymbol{\theta}, k')$  to  $(\boldsymbol{\theta}, k)$  and see how it is different from  $T_{\boldsymbol{\theta}, k}^N$ . A manifold in general is said to be flat or curvature-free when its tangent directions are the same at all the points. In the present case, we can compare two tangent spaces  $T_{\boldsymbol{\theta}, k}^N$  and  $T_{\boldsymbol{\theta}, k'}^N$  by the  $m$ -parallel transport of one to the other. We give formal definitions of  $m$ -flatness and  $m$ -information-curvature-freeness.

**Definition 1.** A semi-parametric statistical model  $S$  is said to be  $m$ -flat or  $m$ -convex, when the  $T_{\boldsymbol{\theta}, k'}^N$  are invariant under the  $m$ -parallel transports, that is

$$\prod_{k'}^{(m)} T_{\boldsymbol{\theta}, k'}^N \subset T_{\boldsymbol{\theta}, k}^N$$

for any  $k, k'$  and  $\boldsymbol{\theta}$ . When the  $m$ -parallel transport of  $T_{\boldsymbol{\theta}, k'}^N$  from  $(\boldsymbol{\theta}, k')$  to  $(\boldsymbol{\theta}, k)$  does not include the  $T_{\boldsymbol{\theta}, k}^E$  components for any  $k, k'$  and  $\boldsymbol{\theta}$ , that is,

$$\prod_{k'}^{(m)} T_{\boldsymbol{\theta}, k'}^N \subset T_{\boldsymbol{\theta}, k}^N \oplus T_{\boldsymbol{\theta}, k}^A, \quad (5.10)$$

the model  $S$  is said to be  $m$ -curvature free in the information directions, or for short,  $m$ -information-curvature free.

It is easy to see that, when  $S$  is  $m$ -flat, it is  $m$ -information-curvature free. When  $S_{\boldsymbol{\theta}}$  is not  $m$ -flat,  $S_{\boldsymbol{\theta}}$  is curved in general, because its tangent directions change as  $k$  changes. However, when  $S_{\boldsymbol{\theta}}$  is  $m$ -information-curvature free, the changes in the tangent directions are restricted to the  $T_{\boldsymbol{\theta}, k}^A$  (and  $T_{\boldsymbol{\theta}, k}^N$ ) directions, showing that  $S_{\boldsymbol{\theta}}$  is not curved in the direction of  $T_{\boldsymbol{\theta}, k}^E$ . Now, by using  $m$ -information-curvature freeness, we show a necessary and sufficient condition that  $F_{\boldsymbol{\theta}, k}^1$  coincides with  $T_{\boldsymbol{\theta}, k}^E$  at every  $k$ .

**Lemma 4.** If and only if  $S_{\boldsymbol{\theta}}$  is  $m$ -information-curvature free, we have, for every  $k$ ,

$$F_{\boldsymbol{\theta}, k}^1 = T_{\boldsymbol{\theta}, k}^E. \quad (5.11)$$

**Proof.** By definition, when  $S_{\boldsymbol{\theta}}$  is  $m$ -information-curvature free

$$\bigcup_{k'}^{(m)} \prod_{k'} T_{\boldsymbol{\theta}, k'}^N \subset T_{\boldsymbol{\theta}, k}^N \oplus T_{\boldsymbol{\theta}, k}^A \quad (5.12)$$

holds for every  $k$  and vice versa. At any  $k$ , (5.12) is equivalent to

$$F_{\boldsymbol{\theta}, k}^{1A} = (F_{\boldsymbol{\theta}, k}^N)^\perp \supset T_{\boldsymbol{\theta}, k}^E, \quad (5.13)$$

or, in other words,

$$\left\langle \prod_{k'}^{(m)} v(x, \theta, k'), u_i^E(x, \theta, k) \right\rangle_{\theta, k} = 0$$

holds for every  $k', v \in T_{\theta, k'}^N$  and for all  $i = 1, \dots, m$ , because of (4.6). The information fibre  $F_{\theta, k}^I$  is constructed by projecting the components of the  $\theta$ -score function  $\mathbf{u}$  onto  $F_{\theta, k}^{IA}$ . It is easy to show that  $F_{\theta, k}^I$  is the projection of  $T_{\theta, k}^E$  onto  $F_{\theta, k}^{IA}$  because  $T_{\theta, k}^E \oplus T_{\theta, k}^A \supset F_{\theta, k}^{IA}$ . Therefore, (5.13) leads to  $F_{\theta, k}^I = T_{\theta, k}^E$ . The converse is obvious because (5.11) implies (5.13) which is equivalent to (5.12).  $\square$

We can now show the amount of loss of information caused by adopting the simple method of estimating functions ( $M$ -estimators), relative to the Cramér–Rao-type bound. We also give a necessary and sufficient condition for the best  $M$ -estimator to be lossless, that is, efficient in the sense that it attains the Cramér–Rao-type bound. Here, we define the optimal estimating function. Let the true distribution be  $p(x, \theta, k_0)$ . An estimating function  $\mathbf{y}(x, \theta)$  is said to be optimal at a point  $k_0$  when the asymptotic variance  $\text{AV}[\hat{\theta}; \mathbf{y}]$  is minimal (in the sense of matrices) among the estimators obtained from all the estimating functions. The optimal estimating function depends on  $k_0$ , and there is no guarantee that we can find it because  $k_0$  is unknown.

**Proposition 5.** *When  $S_{\theta}$  is  $m$ -information-curvature free,  $\mathbf{y}(x, \theta) = \mathbf{u}^I(x, \theta, k_0) = \mathbf{u}^E(x, \theta, k_0)$  is the optimal estimating function at  $k_0$  and is efficient at  $k_0$ . When  $S_{\theta}$  is not  $m$ -information-curvature free and  $\mathbf{u}^I(x, \theta, k_0)$  is not equal to  $\mathbf{u}^E(x, \theta, k_0)$ , the minimal loss of information by using the  $M$ -estimators is*

$$G^E - G^I = E_{\theta, k_0}[(\mathbf{u}^E - \mathbf{u}^I)(\mathbf{u}^E - \mathbf{u}^I)^T]. \quad (5.14)$$

**Proof.** When  $S_{\theta}$  is  $m$ -information-curvature free, we have  $\mathbf{u}^I = \mathbf{u}^E$  from (5.11). Therefore,  $G^I = G^E$ , showing that  $\mathbf{u}^I(x, \theta, k_0)$  is optimal and efficient at  $k_0$ . On the other hand,  $\mathbf{u}^E = \mathbf{u}^I + (\mathbf{u}^E - \mathbf{u}^I)$  is an orthogonal decomposition so that the loss of information is given by (5.14).  $\square$

It should be noted that most semi-parametric models so far treated by many researchers are  $m$ -flat. The important role of  $m$ -flatness in the estimation function method is noted by Amari and Kumon (1988), Amari (1987a), and also by Bickel *et al.* (1993) under the name of convexity. The present result shows that the  $m$ -information-curvature freeness is essential, establishing a necessary and sufficient condition for the estimating function method to be lossless. We later give an example in which the model is not  $m$ -flat but  $m$ -information-curvature free, and the optimal estimating function is efficient. However, the optimal estimating function depends on the true  $k$  so that there is still a serious problem of choosing a good  $k_0$  from observed data to derive a good estimating function. It is a merit of estimating functions that, even if we misspecify the true  $k$  and choose an incorrect  $k_0$ , the estimator is still  $\sqrt{n}$ -consistent. The next section proposes a method of choosing a simple but good quasi-estimating function based on the observed data.

## 6. Simple method of choosing good quasi-estimating functions

The final problem is how to choose a good estimating function. It is clear that when  $\mathbf{u}^I(x, \theta, k)$  does not depend on  $k$ , this  $\mathbf{u}^I(x, \theta)$  gives the best estimating function without any loss of Fisher information. There are many examples belonging to this class (Amari and Kumon 1988).

When  $\mathbf{u}^I(x, \theta, k)$  includes an unknown  $k$ , one orthodox idea is first to find a consistent estimator  $\hat{k} = \hat{k}(x_1, \dots, x_n, \theta)$  of  $k$  where  $\theta$  is fixed. Given  $x_1, \dots, x_n$  and hence  $\hat{k}$ , we have that  $\mathbf{u}^I(x, \theta, \hat{k})$  belongs to  $F_{\theta, k}^I \oplus F_{\theta, k}^A$  for any  $k$ . This is the optimal estimating function at  $k = \hat{k}$  when it is applied to observations  $x_i$  other than those used for estimating  $\hat{k}$ . However, when we apply this estimating function to  $x_i$  from which  $\hat{k}$  is obtained, we encounter a problem in analysing the estimating equation

$$\sum_{i=1}^n \mathbf{u}^I(x_i, \theta, \hat{k}) = 0,$$

because the  $\mathbf{u}^I(x_i, \theta, \hat{k})$  are not independent but are dependent through random variables  $\hat{k}(x_1, \dots, x_n, \theta)$ . We call  $\mathbf{u}^I(x, \theta, \hat{k})$  a quasi-estimating function. Theoretically, we can often avoid the difficulty by dividing  $n$  observations into  $\sqrt{n}$  and  $n - \sqrt{n}$  disjoint subsets and by using  $\sqrt{n}$  observations for estimating  $k$  and the other  $n - \sqrt{n}$  observations for estimating  $\theta$ . We do not discuss this problem further (see, for example, Pfanzagl 1990; Bhanja and Ghosh 1992). From a practical point of view, computer simulations show that the use of quasi-estimating functions  $u_i(x, \theta, \hat{k})$  is justified.

It is the point of an estimating function  $\mathbf{u}^I(x, \theta, k)$  that even a misspecified  $k$  still gives a  $\sqrt{n}$ -consistent estimator. Therefore, it is wise for practical purposes to choose a simple but good  $k$ . We propose the following new method. Since we know at least in principle the set of all estimating functions, we choose a parametric family  $\mathbf{y}(x, \theta; \boldsymbol{\eta})$  of estimating functions, where  $\boldsymbol{\eta}$  is a finite-dimensional parameter. We then obtain an adequate  $\hat{\boldsymbol{\eta}}$  based on observed data  $x_1, \dots, x_n$ . There are a number of methods for doing so. There again remains the problem of justifying the application of the quasi-estimating function  $y(x, \theta; \hat{\boldsymbol{\eta}})$  to the data  $x_i$  from which  $\hat{\boldsymbol{\eta}}$  is obtained. However, since  $\hat{\boldsymbol{\eta}}$  is finite-dimensional, it is easier to justify this by using expansions similar to (2.5). We do not discuss this point further. In order to choose  $\hat{\boldsymbol{\eta}}$ , one idea is to use a parametrized subset of  $K$ ,

$$M = \{k(\boldsymbol{\eta})\}, \quad M \subset K,$$

where  $\boldsymbol{\eta}$  is a finite-dimensional parameter that specifies  $k$ . Since the true  $k$  is not necessarily included in the subset  $M$ , we have a statistical model

$$S^* = \{p(x, \theta, \boldsymbol{\eta}) = p[x, \theta, k(\boldsymbol{\eta})]\}$$

parametrized by a finite number of parameters  $(\theta, \boldsymbol{\eta})$  which might not include the true distribution. It is not difficult to obtain an estimate  $(\tilde{\theta}, \tilde{\boldsymbol{\eta}})$ , say, by the maximum likelihood method. However, this  $\tilde{\theta}$  is not consistent in general and the idea of using  $\tilde{\theta}$  was dismissed long ago. The new idea is to use the quasi-estimating function

$$\mathbf{u}^I\{x, \theta, k(\tilde{\boldsymbol{\eta}})\}$$

to obtain a good  $\sqrt{n}$ -consistent estimator  $\hat{\theta}$ . This type of idea is similar to one proposed by Lindsay (1985).

## 7. Examples

In order to explain the basic concepts, we use the simple mixture model of (2.13).

**Example 2. Mixture model.** A mixture model is  $m$ -flat and is hence  $m$ -information-curvature free. In particular, when a model is given by (2.13) and (2.14), the  $\theta$ -score, the information score  $\mathbf{u}^I$ , and the nuisance score in the direction of  $a(\xi)$  can be calculated explicitly.

We first calculate the  $\theta$ -score  $\mathbf{u}$ ,

$$\begin{aligned}\mathbf{u} &= \frac{\partial}{\partial \theta} \log p \\ &= \frac{1}{p(x, \theta, k)} \int (\partial_{\theta} \mathbf{s} \cdot \xi + \partial_{\theta} r - \partial_{\theta} \psi) k(\xi) \exp\{\xi \cdot \mathbf{s} - r - \psi\} d\xi.\end{aligned}$$

Noting that the conditional distribution  $p(\xi|\mathbf{s})$  of  $\xi$  conditioned on  $\mathbf{s}$  is written as

$$p(\xi|\mathbf{s}) = \frac{k(\xi) \exp\{\xi \cdot \mathbf{s} - \psi\}}{\int k(\xi) \exp\{\xi \cdot \mathbf{s} - \psi\} d\xi},$$

the  $\theta$ -score may be written as

$$\mathbf{u} = \partial_{\theta} \mathbf{s} \cdot \mathbf{E}[\xi|\mathbf{s}] + \partial_{\theta} r - \mathbf{E}[\partial_{\theta} \psi|\mathbf{s}],$$

where  $\mathbf{E}[\cdot|\mathbf{s}]$  is the conditional expectation. Similarly, the nuisance score in the direction of  $a(\xi)$  is given by

$$v[a] = \mathbf{E} \left[ \frac{a(\xi)}{k(\xi)} \middle| \mathbf{s} \right].$$

Therefore,  $v[a]$  depends on  $x$  only through  $\mathbf{s}$  so that the nuisance subspace  $T_{\theta, k}^N$  is generated by the random variable  $\mathbf{s}(x, \theta)$ .

It is known that the projection of a random variable  $t$  onto the space generated by  $s_i$  is given by the conditional expectation  $\mathbf{E}[t|s_i]$  and the projection onto the orthogonal complement is  $t - \mathbf{E}[t|s_i]$ . Hence, the efficient score, which is the same as the information score in this case, is given by

$$\begin{aligned}\mathbf{u}^I &= \mathbf{u}^E = \mathbf{u} - \mathbf{E}[\mathbf{u}|\mathbf{s}] \\ &= \{\partial_{\theta} \mathbf{s} - \mathbf{E}[\partial_{\theta} \mathbf{s}|\mathbf{s}]\} \cdot \mathbf{E}[\xi|\mathbf{s}] + \{\partial_{\theta} r - \mathbf{E}[\partial_{\theta} r|\mathbf{s}]\},\end{aligned}$$

where the vector notation should be understood appropriately. This gives the efficient

estimating functions. There is an interesting special case when  $\partial_{\theta}\mathbf{s}$  is a function of  $\mathbf{s}$ . In this case,  $\partial_{\theta}\mathbf{s} = \mathbf{E}[\partial_{\theta}\mathbf{s}|\mathbf{s}]$ . The efficient score is then given by

$$\mathbf{u}^{\mathbf{E}} = \partial_{\theta}r - \mathbf{E}[\partial_{\theta}r|\mathbf{s}].$$

This does not depend on  $k(\boldsymbol{\xi})$ , so that  $\mathbf{u}^{\mathbf{E}}$  is the optimal estimation function at any  $k(\boldsymbol{\xi})$ . There is no information loss.

**Example 1.** *Linear dependence model.* We now apply our general theory to the linear dependence problem of (2.12), where we put  $\sigma^2 = 1$  for simplicity's sake. We further put  $x = (\alpha, \beta)$ , and

$$\begin{aligned} r(x, \theta) &= -\frac{1}{2}(\alpha^2 + \beta^2), \\ s(x, \theta) &= \alpha + \theta\beta, \\ \partial_{\theta}s(x, \theta) &= \beta, \\ \psi(\xi, \theta) &= \frac{1}{2}\xi^2(1 + \theta^2) + \log 2\pi. \end{aligned}$$

The  $\theta$ -score is

$$u = \beta\mathbf{E}[\xi|s] - \theta\mathbf{E}[\xi^2|s].$$

Since

$$\mathbf{E}[\beta|s] = \frac{\theta}{1 + \theta^2}s,$$

we have

$$u^{\mathbf{I}} = u^{\mathbf{E}} = \frac{1}{1 + \theta^2}\{\beta - \theta\alpha\}\mathbf{E}[\xi|s].$$

In order to obtain a good quasi-estimating function, we need to estimate  $\mathbf{E}[\xi|s]$ , which might not be easy.

We study a special case first. When  $k(\boldsymbol{\xi})$  is a normal distribution with mean  $\mu_{\xi}$  and variance  $\sigma_{\xi}^2$ , we have

$$p(s, \xi; \theta, k) = c(s) \exp\left\{\xi s - \frac{(\xi - \mu_{\xi})^2}{2\sigma_{\xi}^2} - \frac{1}{2}(1 + \theta^2)\xi^2\right\}.$$

Hence, after some calculations, we have

$$\mathbf{E}[\xi|s] = \text{const.} \left(s + \frac{\mu_{\xi}}{\sigma_{\xi}^2}\right).$$

Therefore, the optimal estimating function is

$$\begin{aligned} u^{\mathbf{I}} &= \left(s + \frac{\mu_{\xi}}{\sigma_{\xi}^2}\right)(\beta - \theta\alpha) \\ &= \left(\alpha + \theta\beta + \frac{\mu_{\xi}}{\sigma_{\xi}^2}\right)(\beta - \theta\alpha). \end{aligned}$$

Motivated by this (see also Pfanzagl 1990), we consider a family of estimating functions parametrized by  $c$ ,

$$y(x, \theta; c) = (\alpha + \theta\beta + c)(\beta - \theta\alpha).$$

The estimator obtained from  $y(x, \theta; c)$  is called the  $c$ -estimator and is denoted by  $\hat{\theta}_c$ . This class includes various estimators. When  $c \rightarrow \infty$ , the estimator is

$$\hat{\theta}_\infty = \frac{\sum \beta_i}{\sum \alpha_i}.$$

We calculate the maximum likelihood estimators  $(\hat{\theta}; \hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_n)$  of all the parameters. Here  $\hat{\theta}$  is the least-squares estimator in the orthogonal direction, and is given by the 0-estimator  $\hat{\theta}_0$ . This estimator is not optimal. When the sequence  $(\xi_1, \xi_2, \dots)$  can be regarded as a realization from  $N(\mu_\xi, \sigma_\xi^2)$ , the estimator obtained from  $c = \mu_\xi/\sigma_\xi^2$  is optimal.

The idea of the  $c$ -estimator springs from the following observation. Let us temporarily assume that  $k(\xi)$  is subject to a normal distribution  $N(\mu_\xi, \sigma_\xi^2)$ . Then, the nuisance function  $k$  is parametrized by  $\eta = (\mu_\xi, \sigma_\xi^2)$ . From the observed data, we have an estimator  $\hat{\eta} = (\hat{\mu}_\xi, \hat{\sigma}_\xi^2)$  in which

$$\begin{aligned}\hat{\mu}_\xi &= \frac{1}{n} \sum \alpha_i \\ \hat{\sigma}_\xi^2 &= \frac{1}{n} \sum \alpha_i^2 - \hat{\mu}_\xi^2 - 1.\end{aligned}$$

We use the  $\hat{c}$ -quasi-estimating function, where  $\hat{c}$  is given by

$$\hat{c} = \frac{\hat{\mu}_\xi}{\hat{\sigma}_\xi^2}.$$

The estimator  $\hat{\theta}_{\hat{c}}$  is consistent and improves the maximum likelihood estimator even when  $k(\xi)$  is not normal.

We now analyse the  $c$ -estimators. Given a sequence  $(\xi_1, \xi_2, \dots)$ , the asymptotic variance of the  $c$ -estimator  $\hat{\theta}_c$  is calculated from

$$V_c = \frac{(1 + \theta^2)\{(c + (1 + \theta^2)\bar{\xi})^2 + (1 + \theta^2)^2(\bar{\xi}^2 - \bar{\xi}^2) + (1 + \theta^2)\}}{\{c\bar{\xi} + (1 + \theta^2)\bar{\xi}^2\}^2},$$

where

$$\begin{aligned}\bar{\xi} &= \frac{1}{n} \sum_{i=1}^n \xi_i, \\ \bar{\xi}^2 &= \frac{1}{n} \sum_{i=1}^n \xi_i^2.\end{aligned}$$

The optimal value of  $c$  is given by

$$c_{\text{opt}} = \frac{\bar{\xi}}{\bar{\xi}^2 - (\bar{\xi})^2}.$$

The unknown  $\bar{\xi}$  and  $\bar{\xi}^2$  are estimated by

$$\hat{\xi} = \frac{1}{n} \sum \alpha_i,$$

$$\hat{\xi}^2 = \frac{1}{n} \sum \alpha_i^2 - 1.$$

The estimator  $\hat{\theta}_c$  is not necessarily efficient, but is close to the optimal among the  $c$ -estimators. We show how it improves the maximum likelihood estimator. For example, when  $\bar{\xi} = 1$ ,  $\bar{\xi}^2 = 2$ ,  $\sigma^2 = 1$ , and  $\theta = 0$ , the variances of the MLE  $\hat{\theta}_0$  and of  $\hat{\theta}_\infty$  are, respectively,

$$V[\bar{\theta}_{\text{MLE}}] \approx \frac{3}{4n},$$

$$V[\hat{\theta}_\infty] \approx \frac{1}{n},$$

while that of the  $c_{\text{opt}}$  estimator is

$$V[\hat{\theta}_c] \approx \frac{2}{3n}.$$

Computer simulations show that  $\hat{\theta}_c$  attains this.

**Example 3.** *Semi-parametric additive regression.* This model, given by (2.15), is analysed here because it is not  $m$ -flat but is  $m$ -information-curvature free. The probability density function is written as

$$p(y, \mathbf{x}, t; \boldsymbol{\theta}, k) = \frac{1}{\sqrt{2\pi}} \exp \left[ -\frac{1}{2} \{y - \boldsymbol{\theta} \cdot \mathbf{x} - k(t)\}^2 \right],$$

which is not linear in  $k$ . The  $\boldsymbol{\theta}$ -score vector at  $(\boldsymbol{\theta}, k)$  is written as

$$\mathbf{u} = \mathbf{x} \{y - \boldsymbol{\theta} \cdot \mathbf{x} - k(t)\},$$

while the nuisance score in the direction of  $a(t)$  of a change in function  $k(t)$  is

$$v = a(t) \{y - \boldsymbol{\theta} \cdot \mathbf{x} - k(t)\}.$$

The efficient score is given by

$$\mathbf{u}^{\text{E}} = \{\mathbf{x} - \mathbf{E}[\mathbf{x}|t]\} \{y - \boldsymbol{\theta} \cdot \mathbf{x} - k(t)\}.$$

Since the model is not  $m$ -flat, the  $m$ -parallel transport of the nuisance tangent space  $T_{\boldsymbol{\theta}, k'}^{\text{N}}$  to  $(\boldsymbol{\theta}, k)$  is different from  $T_{\boldsymbol{\theta}, k}^{\text{N}}$ . The  $m$ -parallel transport of  $v$  from  $k'$  to  $k$  is given by

$$\begin{aligned} \prod_{k'}^{(m)} v &= \frac{p(y, \mathbf{x}, t; \boldsymbol{\theta}, k')}{p(y, \mathbf{x}, t; \boldsymbol{\theta}, k)} v \\ &= a(t) \{y - \boldsymbol{\theta} \cdot \mathbf{x} - k'(t)\} \exp \left\{ (k' - k) \left( y - \boldsymbol{\theta} \cdot \mathbf{x} - \frac{k + k'}{2} \right) \right\}. \end{aligned}$$

This does not belong to  $T_{\theta,k}^N$ . The information score  $\mathbf{u}^I$  is the projection of  $\mathbf{u}$  onto the space orthogonal to that spanned by all of  $\prod_{k'}^{(m)k} T_{\theta,k'}^N$ .

In the present case, we can show that

$$\prod_{k'}^{(m)k} T_{\theta,k'}^N \subset T_{\theta,k}^N \oplus T_{\theta,k}^A.$$

This is proved by showing, for any  $v \in T_{\theta,k'}^N$ ,

$$\begin{aligned} & \left\langle \mathbf{u}^E, \prod_{k'}^{(m)k} v \right\rangle_{\theta,k} \\ &= \mathbb{E} \left[ \{ \mathbf{x} - \mathbb{E}[\mathbf{x}|t] \} a(t) (y - \boldsymbol{\theta} \cdot \mathbf{x} - k) (y - \boldsymbol{\theta} \cdot \mathbf{x} - k') \exp \left\{ (k' - k) \left( y - \boldsymbol{\theta} \cdot \mathbf{x} - \frac{k+k'}{2} \right) \right\} \right] \\ &= 0 \end{aligned}$$

To prove this, we first calculate the expectation of the above with respect to  $y$ , giving

$$\begin{aligned} & \left\langle \mathbf{u}^E, \prod_{k'}^{(m)k} v \right\rangle_{\theta,k} \\ &= \mathbb{E}_{\mathbf{x},t} \left[ (\mathbf{x} - \mathbb{E}[\mathbf{x}|t]) a(t) \mathbb{E}_{\epsilon} \left[ \epsilon (\epsilon - k' + k) \exp \left\{ (k' - k) \left( \epsilon - \frac{k' - k}{2} \right) \right\} \right] \right], \end{aligned}$$

where  $\mathbb{E}_{\epsilon}$  is the expectation with respect to the normal random variable

$$\epsilon = y - \boldsymbol{\theta} \cdot \mathbf{x} - k(t)$$

and  $\mathbb{E}_{\mathbf{x},t}$  is the expectation with respect to  $q(\mathbf{x}, t)$ . Since the term in  $\mathbb{E}_{\epsilon}$  depends on  $t$  but not on  $\mathbf{x}$ , by taking the conditional expectation  $\mathbb{E}_{\mathbf{x}|t}$ , we have

$$\left\langle \mathbf{u}^E, \prod_{k'}^{(m)k} v \right\rangle_{\theta,k} = 0$$

which shows that  $S_{\theta}$  is curved only in the direction of  $T_{\theta,k}^A$ . Hence, we have

$$\mathbf{u}^I = \mathbf{u}^E.$$

Any estimating function is equivalently, or to within a matrix  $T(\boldsymbol{\theta}, k)$ , written as

$$\begin{aligned} \mathbf{y} &= \mathbf{u}^I + \mathbf{a} \\ &= \{ \mathbf{x} - \mathbb{E}[\mathbf{x}|t] \} \{ y - \boldsymbol{\theta} \cdot \mathbf{x} - k(t) \} + \mathbf{a} \end{aligned}$$

for any  $k(t)$  and  $\mathbf{a} \in F_{\theta,k}^A$ . From (5.8),  $G^A = 0$  for the estimator with  $\mathbf{a} = 0$ , so that this gives an efficient estimator. It is noteworthy that, even when  $k(t)$  is misspecified, this  $\mathbf{y}$  gives a  $\sqrt{n}$ -consistent estimator provided  $\mathbb{E}[\mathbf{x}|t]$  is known or estimated well from the data. Cuzick (1992) discussed efficient estimation of  $\mathbb{E}[\mathbf{x}|t]$  and  $k(t)$  by smoothing.



## Acknowledgements

The authors thank the editor and reviewers for their constructive comments, suggestions and criticisms, which led to significant improvements to the paper.

## References

- Amari, S. (1985) *Differential-Geometrical Method in Statistics*. Lecture Notes in Statistics 28. New York: Springer-Verlag.
- Amari, S. (1987a). Dual connections on the Hilbert bundles of statistical models. In C.T.J. Dodson (ed.), *Geometrization of Statistical Theory*, pp. 123–152. Lancaster: University of Lancaster Department of Mathematics.
- Amari, S. (1987b) Differential geometry of a parametric family of invertible linear systems – Riemannian metric, dual affine connections and divergence. *Math. Systems Theory*, **20**, 53–82.
- Amari, S. and Han, T.S. (1989) Statistical inference under multi-terminal rate restrictions. *IEEE Trans. Inform. Theory*, **35**, 217–227.
- Amari, S. and Kumon, M. (1988) Estimation in the presence of infinitely many nuisance parameters – geometry of estimating functions. *Ann. Statist.*, **16**, 1044–1068.
- Amari, S., Kurata, K. and Nagaoka, H. (1992) Information geometry of Boltzmann machines. *IEEE Trans. Neural Networks*, **3**, 260–271.
- Andersen, E.B. (1970) Asymptotic properties of conditional maximum likelihood estimators. *J. Roy. Statist. Soc. B*, **32**, 283–301.
- Barndorff-Nielsen, O.E. (1986) Likelihood and observed geometries. *Ann. Statist.*, **14**, 856–873.
- Begun, J.M., Hall, W.J., Huang, W.M. and Wellner, J.A. (1983) Information and asymptotic efficiency in parametric-nonparametric models. *Ann. Statist.*, **11**, 432–452.
- Bhanja, J. and Ghosh, J.K. (1992) Efficient estimation with many nuisance parameters. *Sankhyā Ser. A*, **54**, 1–39, 135–156.
- Bickel, P.J., Klaassen, C.A.J., Ritov, Y. and Wellner, J.A. (1993) *Efficient and Adaptive Estimation for Semiparametric Models*. Baltimore, MD: Johns Hopkins University Press.
- Cuzick, J. (1992) Semiparametric additive regression. *J. Roy. Statist. Soc. B*, **54**, 831–843.
- Friedrich, T. von (1991) Die Fisher-Information und symplektische Strukturen. *Math. Nachr.*, **153**, 273–296.
- Godambe, V.P. (1976) Conditional likelihood and unconditional optimum estimating equations. *Biometrika*, **63**, 277–284.
- Godambe, V.P. (ed.) (1991) *Estimating Functions*. New York: Oxford University Press.
- Godambe, V.P. and Heyde, C.C. (1987) Quasi-likelihood and optimal estimation. *Internat. Statist. Rev.*, **55**, 231–244.
- Groeneboom, P. and Wellner, J.A. (1992) *Information Bounds and Nonparametric Maximum Likelihood Estimation*. Basel: Birkhäuser.
- Hasminskii, R.Z. and Ibragimov, I.A. (1983) Efficient estimation in the presence of infinite dimensional incidental parameters. In *Probability Theory and Mathematical Statistics*. Lecture Notes in Math. 1021, pp. 195–229, New York: Springer-Verlag.
- Huber, P.J. (1981) *Robust Statistics*. New York: Wiley.
- Jaffe, A. and Quinn, F. (1993) Theoretical mathematics. *Bull. Amer. Math. Soc.*, **29**, 1–13.
- Kanbayashi, T. (1994) Statistical manifold of infinite dimensions (in Japanese). *Trans. Japan SIAM*, **4**, 211–228.

- Kumon, M. and Amari, S. (1984) Estimation of a structural parameter in the presence of a large number of nuisance parameters. *Biometrika*, **71**, 445–459.
- Lafferty, J.D. (1988) The density manifold and configuration space, quantization. *Trans. Amer. Math. Soc.*, **305**, 699–741.
- Levit, B.Y. (1978) Infinite-dimensional information inequalities. *Theory Probab. Appl.*, **23**, 371–377.
- Lindsay, B.G. (1982) Conditional score functions: Some optimality results. *Biometrika*, **69**, 503–512.
- Lindsay, B.G. (1985) Using empirical partially Bayes inference for increased efficiency. *Ann. Statist.*, **13**, 914–931.
- McLeish, D.L. and Small, C.G. (1988) *The Theory and Applications of Statistical Inference Functions*. Lecture Notes in Statistics 44. New York: Springer-Verlag.
- Murray, M.K. and Rice, J.W. (1993) *Differential Geometry and Statistics*. New York: Chapman & Hall.
- Nagaoka, H. and Amari, S. (1982) Differential geometry of smooth families of probability distributions. Technical Report 82-7, University of Tokyo.
- Neyman, J. and Scott, E.L. (1948) Consistent estimates based on partially consistent observations. *Econometrica*, **32**, 1–32.
- Okamoto, I., Amari, S. and Takeuchi, K. (1991) Asymptotic theory of sequential estimation: Differential geometrical approach. *Ann. Statist.*, **19**, 961–981.
- Pfanzagl, J. (1990) *Estimation in Semiparametric Models: Some Recent Developments*. Lecture Notes in Statistics 63. New York: Springer-Verlag.
- Ritov, Y. and Bickel, P.J. (1990) Achieving information bounds in non and semiparametric models. *Ann. Statist.*, **18**, 925–938.
- Small, C.G. and McLeish, D.L. (1988) Generalization of ancillarity, completeness and sufficiency in an inference function space. *Ann. Statist.*, **16**, 534–551.
- Small, C.G. and McLeish, D.L. (1989) Projection as a method for increasing sensitivity and eliminating nuisance parameters. *Biometrika*, **73**, 693–703.
- van der Vaart, A.W. (1991) On differentiable functionals. *Ann. Statist.*, **19**, 178–205.

Received July 1994 and revised April 1996