

ESTIMATING NONLINEAR FUNCTIONAL RELATIONS
BY MAXIMUM LIKELIHOOD AND LEAST SQUARES

by

J.W. Sawyer, Jr. and K.L.Q. Read

Department of Biostatistics
University of North Carolina at Chapel Hill

Institute of Statistics Mimeo Series No. 1374

January 1982

Estimating Nonlinear Functional Relations
by Maximum Likelihood and Least Squares

J. W. Sawyer and K.L.Q. Read

Institute of Statistics Mimeo Series No. 1374

Errata

Page 1: eq. (1.1) should read

$$x_{i,p+1} = f(x_{i1}, \dots, x_{ip}, \alpha), \quad i = 1, 2, \dots, t$$

Line 4 from bottom: asterisk V as: V^* .

insert footnote at bottom of page:

* If V is taken to refer to arbitrarily chosen vectors

$$\underline{X}_{i1}, \dots, \underline{X}_{it}, \text{ concatenated as } \underline{X}' = (\underline{X}'_{i1}, \dots, \underline{X}'_{it}),$$

the m_i 's need not be equal.

Page 2: line 2: V, not Σ .

para 2 line 6: insert after ... of the data: (that is, one

para 2 line 8: parameters). (insert right parenthesis)

Page 4: 6 lines above eq. (2.1): the pr + r functions

$$\text{eq. (2.1): } \alpha_v = g_v(x_{i1}, \dots, x_{ir}), \quad v = 1, 2, \dots, r$$

Page 6: 2 lines below eq. (2.5): any subset i_1, \dots, i_r of the ...

$$\text{Page 8: } 2^{\text{nd}} \text{ eq: } S_N(x_{i1}, \dots, x_{it}) = \sum_{i=1}^t m_i (\bar{X}_i - x_i)' \hat{\Sigma}^{-1} (\bar{X}_i - x_i)$$

$$\text{Page 12: } 2^{\text{nd}} \text{ eq: } \frac{\partial^2 N}{\partial \sigma_{ij}^2} = -N\sigma^{ij} + \frac{1}{2} + \text{tr}[\underline{\Sigma}^{-1}[\underline{E}_{ij} + \underline{E}_{ji}]\underline{\Sigma}^{-1}\underline{B}]$$

Estimating Nonlinear Functional Relations
by Maximum Likelihood and Least Squares

by

J. W. Sawyer, Jr.¹

K. L. Q. Read²

¹ Department of Biostatistics, School of Hygiene and Public Health,
The Johns Hopkins University, Baltimore, Maryland 21205.

² Department of Mathematical Statistics and Operations Research,
University of Exeter, Exeter EX4 4PU.

SUMMARY

It is assumed that $p + 1$ unknown variables at t points are such that the $(p + 1)^{st}$ variable is given by a nonlinear function, f , of the other p variables and a vector α of structural parameters. Replicate observations on the $p + 1$ variables for each of the t points are assumed to be mutually independent and multivariate normal, with common covariance matrix Σ . For known Σ , obtaining the maximum likelihood estimate $\hat{\alpha}$ of α is equivalent to minimizing a weighted sum of squares with respect to α and pt nuisance parameters. If Σ is replaced by a consistent estimate $\hat{\Sigma}$ in the weighted sum of squares, the weighted least squares estimate $\tilde{\alpha}$ is no longer equal to $\hat{\alpha}$. However, for N total observations, $\sqrt{N}(\hat{\alpha} - \alpha)$ and $\sqrt{N}(\tilde{\alpha} - \alpha)$ have the same asymptotic distribution. This result implies that, even if Σ must be estimated, a close relationship exists between maximum likelihood and the generalized least squares procedure which Dolby (1972) has already applied to nonlinear functional relations with Σ known.

KEYWORDS: Nonlinear functional relation, maximum likelihood, weighted least squares, generalized least squares

1. INTRODUCTION

For $(p + 1)$ -dimensional vectors \underline{x}_i , $i = 1, 2, \dots, t$, for an r -dimensional vector $\underline{\alpha}$, $r \leq t$, and for known f , let

$$\underline{x}_{i,y+1} = f(\underline{x}_i, \dots, \underline{x}_{ip}, \underline{\alpha}), \quad i = 1, 2, \dots, t \quad (1.1)$$

hold for some unknown nonrandom vectors $\underline{x}_1^0, \dots, \underline{x}_t^0$, and $\underline{\alpha}^0$. Let observations

$$\underline{X}_{ik} = \underline{x}_i^0 + \underline{\epsilon}_{ik}, \quad k = 1, 2, \dots, m_i; i = 1, \dots, t, \quad (1.2)$$

with $E(\underline{\epsilon}_{ik}) = \underline{0}$, exist on each \underline{x}_i^0 . Then the observations (1.2) are said to have a functional relationship. Most treatments of the problem of estimating $\underline{\alpha}_0$ have assumed that each \underline{X}_{ik} is multivariate normal.

If t is fixed while each m_i becomes large, then maximum likelihood is an appropriate method of estimation. Dolby and Freeman (1975) discuss maximum likelihood estimation for linear or nonlinear f , with all $m_i = m > 1$. They allow for correlation between \underline{X}_{ik} and \underline{X}_{jk} for $i \neq j$, but require that \underline{X}_{ik} and \underline{X}_{jh} be independent for $k \neq h$, even if $i = j$. Let $\underline{X}'_k = (\underline{X}'_{1k}, \dots, \underline{X}'_{tk})$, and let all \underline{X}_k have variance-covariance matrix \underline{V}_k .^{*} Then another method of estimating $\underline{\alpha}_0$ is to calculate $\tilde{\underline{\alpha}}$ minimizing the weighted sum of squares

$$\sum_{k=1}^m \left(\underline{X}_k - \tilde{\underline{\alpha}} \right)' \underline{V}_k^{-1} \left(\underline{X}_k - \tilde{\underline{\alpha}} \right), \quad (1.3)$$

^{*} If \underline{V} is taken to refer to arbitrarily chosen vectors $\underline{X}_{1i_1}, \dots, \underline{X}_{ti_t}$, concatenated as $\underline{X}' = (\underline{X}'_{1i_1}, \dots, \underline{X}'_{ti_t})$, the m_i 's need not be equal.

where $\tilde{\mathbf{x}}' = (\tilde{x}'_1, \dots, \tilde{x}'_t)$ satisfies the equations (1.1). For known Σ and multivariate normal observations, $\tilde{\alpha}$ is clearly the maximum likelihood estimate. The need to estimate the p t nuisance parameters $x_{11}^0, \dots, x_{1p}^0, x_{21}^0, \dots, x_{tp}^0$ can be eliminated by a modification of (1.3) which involves first order Taylor series approximations of the constraints (1.1). For $p = 1$ and known V , Dolby (1972) demonstrates that this approach, known as generalized least squares, is equivalent to maximum likelihood, except for the error arising from the first order linear approximation to a nonlinear f . For any p and for unknown V , Dolby and Freeman (1975) establish that the maximum likelihood and generalized least square estimates are equivalent if f is linear.

The equivalence to within an approximation of generalized least squares and maximum likelihood for nonlinear f and known V , which Dolby (1972) establishes for $p = 1$, can easily be extended to the case of arbitrary p by a similar development. Now suppose that a consistent estimate \hat{V} of V replaces V in (1.3), and further suppose that \hat{V} is a fixed, never-updated function of the data (that is, one which does not involve iterated estimates of α_0 or the nuisance parameters). The estimate $\tilde{\alpha}$ obtained by minimizing (1.1) will still be equivalent to a generalized least squares estimate based on \hat{V} instead of V , to within a first order approximation. This can be established by viewing the negative of (1.3) as a "likelihood" to be maximized, and employing the same arguments as Dolby applies to the likelihood when V is known.

The primary purpose of this paper is to demonstrate that, even if \underline{V} must be estimated, minimization of (1.3) is still as precise a procedure as maximum likelihood, in the sense that $\tilde{\underline{\alpha}}$ and the maximum likelihood estimate have the same asymptotic distribution. In the light of the discussion above, this means that any asymptotic inefficiency involved in using generalized least squares as opposed to maximum likelihood for nonlinear f will be due only to the approximation inherent in the generalized least squares estimate.

For normal data the equality of the asymptotic distributions of a weighted least squares estimate $\tilde{\underline{\alpha}}$, employing random weights, and the maximum likelihood estimate, will be established in Section 4. Some preliminary results necessary for Section 4 will be presented in Section 2, along with some results necessary to Section 3, which establishes the consistency of $\tilde{\underline{\alpha}}$ without assuming that the data are normal.

2. PRELIMINARIES

The formal model to be treated is as follows: Let (1.1) and (1.2) hold, and further require that X_{ik} and X_{jh} be independent unless $i = j$ and $k = h$. Let all X_{ik} have non-singular covariance matrix $\underline{\Sigma}$. Thus \underline{V} becomes block diagonal, with all diagonal blocks equal to $\underline{\Sigma}$. (This model is intuitively appealing in that it allows sums of squares and cross products to be pooled in order to estimate $\underline{\Sigma}$. With minor modifications, the arguments which follow will work just as well for the model with general \underline{V} treated

in Dolby and Freeman (1975), or for a model which allows a separate matrix Σ_i for each set of replications X_{i1}, \dots, X_{im_i} , $i = 1, 2, \dots, t$). Let all m_i increase in fixed proportions, that is, for $N = \Sigma m_i$, let $m_i = \pi_i N$ for fixed π_i .

Let f possess first partials with respect to all its arguments which are continuous in a neighborhood of the true values. For every choice i_1, \dots, i_r of the integers from 1 to t , let the Jacobian of the matrix with elements

$$\frac{\partial x_{i_\mu, p+1}}{\partial \alpha_\nu}, \quad \mu, \nu = 1, 2, \dots, r$$

be different from 0. It will follow that, in a neighborhood of the true values, the $pt + r$ functions

$$x_{i_\mu, p+1} = f \left(x_{i_\mu 1}, \dots, x_{i_\mu p}, \underline{\beta} \right) \quad \mu = 1, 2, \dots, r$$

and

$$x_{i_\mu j} = x_{i_\mu j} \quad \mu = 1, 2, \dots, r, \quad j = 1, 2, \dots, p$$

will have an inverse possessing continuous first partials, and in particular,

$$\beta_\nu = g_\nu \left(x_{i_1}, \dots, x_{i_r} \right), \quad \nu = 1, 2, \dots, r \quad (2.1)$$

for some vector function g which possesses continuous first partials. It follows that, within some neighborhood of the true values any vectors x_1, \dots, x_t satisfying (1.1) are such that the corresponding α may be obtained by any one of $\binom{t}{r}$ possible functions g .

This result can be used to establish that if some estimates $\bar{x}_1, \dots, \bar{x}_t$ of the true population means satisfy (1.1), and if the differences $\bar{x}_1 - x_1^0, \dots, \bar{x}_t - x_t^0$ are $O_p(N^{-\frac{1}{2}})$, then $\bar{\alpha} - \alpha_0$ is $O_p(N^{-\frac{1}{2}})$ also.

This result in turn will be necessary in order to apply a result of Weiss (1971). Let a data vector Z_n have log-likelihood $l(Z_n, \theta^0)$ for a q -dimensional vector θ_0 . Let $\tilde{\theta}_n$ be an estimate of θ^0 satisfying

$$\lim_{n \rightarrow \infty} P\{\sqrt{n}|\tilde{\theta}_{in} - \theta_i^0| \geq L(n)\} \rightarrow 0 \quad (2.2)$$

for $i = 1, 2, \dots, q$ and all non-random sequences $\{L(n)\}$ with $L(n) \rightarrow \infty$. (Condition (2.2) holds if and only if $\tilde{\theta}_{in} - \theta_i^0$ is $O_p(n^{-\frac{1}{2}})$). Let $A_n(\tilde{\theta}_n)$ be the vector $\partial l / \partial \theta$ evaluated at $\tilde{\theta}_n$, and let $I_n(\tilde{\theta}_n)$ be the $q \times q$ information matrix evaluated at $\tilde{\theta}_n$. Under mild regularity conditions, Weiss demonstrates that

$$\theta_n^* = \tilde{\theta}_n + [I_n(\tilde{\theta}_n)]^{-1} A_n(\tilde{\theta}_n) \quad (2.3)$$

is such that $\sqrt{n}(\theta_n^* - \theta^0)$ has the same asymptotic distribution as $\sqrt{n}(\hat{\theta}_n - \theta^0)$, where $\hat{\theta}_n$ is the maximum likelihood estimate.

Now if n is equated with N for the model described at the beginning of this section, and the data are multivariate normal, θ^0 will be the vector of dimension $pt + r + (p+1)(p+2)/2$ comprising α , the pt free parameters $x_{11}^0, \dots, x_{1p}^0, x_{21}^0, \dots, x_{tp}^0$, and the free parameters of Σ . Further, $I_N(\theta^0)/N$ will be continuous and non-random

for all permissible N , and $I_{\tilde{\theta}_N}(\tilde{\theta}_N)/N$ will converge to $I_{\tilde{\theta}_N}(\theta^0)/N$ in probability. Now (2.3) can be rewritten

$$\sqrt{N}(\hat{\theta}_N^* - \theta^0) - \sqrt{N}(\tilde{\theta}_N - \theta^0) = (I_{\tilde{\theta}_N}(\tilde{\theta}_N)/N)^{-1} (A_{\tilde{\theta}_N}(\tilde{\theta}_N)/\sqrt{N}).$$

It follows by Slutsky's theorem and a theorem of Cramér (1946) that if

$$A_{\tilde{\theta}_N}(\tilde{\theta}_N)/\sqrt{N} = (\partial l / \partial \tilde{\theta}_N) / \sqrt{N} \quad \dots \quad (2.4)$$

converges in probability to 0, $\sqrt{N}(\hat{\theta}_N^* - \theta^0)$ and $\sqrt{N}(\tilde{\theta}_N - \theta^0)$, and hence $\sqrt{N}(\hat{\theta}_N - \theta^0)$ and $\sqrt{N}(\tilde{\theta}_N - \theta^0)$, will have the same asymptotic distribution.

The following simple result establishes that $\tilde{\alpha}$ meets the criteria (2.2) if the corresponding vectors $\tilde{x}_1, \dots, \tilde{x}_t$ satisfying (1.1) do. Let h_1, \dots, h_s be non-negative functions of random variables Y_1, \dots, Y_s , respectively. Then for nonrandom y

$$P\left(\sum_{i=1}^s h_i(Y_i) \geq y\right) \leq \sum_{i=1}^s P\left(h_i(Y_i) \geq y/s\right) \quad (2.5)$$

A simple proof of (2.5) for $s = 2$ is given in Tucker (1967). Now choose any subset i_1, \dots, i_r at the subscripts 1 to t . Then with probability approaching 1, $\tilde{x}_{i_1}, \dots, \tilde{x}_{i_r}$ will lie in the neighborhood for which the associated function g , as defined by (2.1), exists. Thus, with probability approaching 1,

$$\sqrt{N}(\tilde{\theta}_\mu - \alpha_\mu^0) = \sum_{v=1}^r \sqrt{N}(\tilde{x}_{i_v} - x_{i_v}^0)' \left[\frac{\partial g}{\partial x_{i_v}^*} \right]$$

where the vector $(\tilde{x}_{i_1}^{*'}, \dots, \tilde{x}_{i_r}^{*'})$ lies on the line segment joining $(\tilde{x}_{i_1}^0, \dots, \tilde{x}_{i_r}^0)$ and $(\tilde{x}_{i_1}^{\prime}, \dots, \tilde{x}_{i_r}^{\prime})$. Since each vector of partials $\partial g / \partial x_{i_\nu}^*$ is continuous, it will converge in probability to $\partial g / \partial x_{i_\nu}^0$. Thus, with probability approaching 1, the absolute values of the partials $\partial g / \partial x_{i_\nu j}^*$ are bounded by some constant M which does not depend on the data. It follows that, for any nonrandom sequence $\{L(N)\}$ with $L(N) \rightarrow \infty$

$$P\left(\sqrt{N}|\tilde{\alpha}_{\mu} - \alpha_{\mu}^0| \geq L(N)\right) \leq P\left(M \sum_{\nu=1}^r \sum_{j=1}^{p+1} \sqrt{N}|\tilde{x}_{i_\nu j} - x_{i_\nu j}^0| \geq L(N)\right)$$

with probability approaching 1 as $N \rightarrow \infty$. Hence, if the differences $\tilde{x}_{i_1} - x_{i_1}^0, \dots, \tilde{x}_{i_t} - x_{i_t}^0$ are $O_p(N^{-\frac{1}{2}})$, we have

$$\lim_{N \rightarrow \infty} P\left(\sqrt{N}|\tilde{\alpha}_{\mu} - \alpha_{\mu}^0| \geq L(N)\right) = 0 \quad (2.6)$$

by virtue of (2.5).

From the results of this section, it follows that, for Section 3, it suffices to show that $\tilde{x}_1, \dots, \tilde{x}_t$ minimizing

$$\sum_{i=1}^t \sum_{k=1}^{m_i} \left(\tilde{x}_{ik} - x_i \right)' \hat{\Sigma}^{-1} \left(\tilde{x}_{ik} - x_i \right) \quad (2.7)$$

subject to (1.1) are such that $\tilde{x}_1 - x_1^0, \dots, \tilde{x}_t - x_t^0$ are $O_p(N^{-\frac{1}{2}})$ for suitably chosen $\hat{\Sigma}$. In Section 4, it will suffice to show that the functions (2.4) associated with the likelihood of the observations $\tilde{x}_{ik}, i = 1, 2, \dots, t, j = 1, 2, \dots, m_i$, converge to 0 on probability when evaluated at the weighted least squares solutions and $\hat{\Sigma}$.

3. CONSISTENCY OF THE LEAST SQUARES ESTIMATORS

For normal data, (2.7) is equivalent to a likelihood conditional on an estimated $\hat{\Sigma}$. Consistency and possibly efficiency of $\tilde{\alpha}$ might be argued from that viewpoint. Postulating a distribution is not required in order to demonstrate that $\tilde{\alpha} - \alpha^0$ is $O_p(N^{-\frac{1}{2}})$, however. It is only required that $\hat{\Sigma} - \Sigma$ be elementwise $O_p(N^{-\frac{1}{2}})$, that it be positive definite with probability approaching 1, and that its definition not be dependent on the minimization of (2.7).

Now (2.7) is equal to

$$\sum_{i=1}^t \sum_{k=1}^{m_i} (x_{ik} - \bar{x}_i)' \hat{\Sigma}^{-1} (x_{ik} - \bar{x}_i) + S_N(x_1, \dots, x_t)$$

where \bar{x}_i is the mean of the m_i replications and

$$S_N(x_1, \dots, x_t) = \sum_{i=1}^t m_i (\bar{x}_i - x_i)' \hat{\Sigma}^{-1} (\bar{x}_i - x_i)$$

Thus it is sufficient to determine $\tilde{x}_1, \dots, \tilde{x}_t$ minimizing S_N under the constraints (1.1). Let $\Sigma = \{\sigma_{jh}\}$ and $\Sigma^{-1} = \{\sigma^{jh}\}$.

Then $S_N(x_1^0, \dots, x_t^0)$ is the sum of terms of the form

$$m_i (\bar{x}_{ij} - x_{ij}^0)^2 \sigma^{ij}$$

or

$$m_i (\bar{x}_{ij} - x_{ij}^0) (\bar{x}_{ih} - x_{ih}^0) \sigma^{jh}$$

for $j \neq h$.

Since each σ^{jh} is a continuous function of the elements of $\hat{\Sigma}$,

$\hat{\sigma}^{jh} \rightarrow \sigma^{jh}$ in probability for $j, h = 1, 2, \dots, p+1$. Thus, with probability approaching 1, all elements of $\hat{\Sigma}$ have a nonrandom upper bound M . Since $m_i = \pi_i N$, the mean \bar{X}_{ij} is such that $\bar{X}_{ij} - x_{ij}^0$ is $O_p(N^{-\frac{1}{2}})$. It follows that, for any nonrandom sequence $\{L(N)\}$ with $L(N) \rightarrow \infty$,

$$P\left(\sqrt{m_i} |\bar{X}_{ij} - x_{ij}^0| \geq \sqrt{L(N)}\right) = P\left(m_i (\bar{X}_{ij} - x_{ij}^0)^2 \geq L(N)\right) \quad (3.1)$$

approaches 0 as $N \rightarrow \infty$. Since

$$2|\bar{X}_{ij} - x_{ij}^0| |\bar{X}_{ih} - x_{ih}^0| \leq |\bar{X}_{ij} - x_{ij}^0|^2 + |\bar{X}_{ih} - x_{ih}^0|^2 \quad (3.2)$$

a result similar to (3.1) applies to the cross-product terms. In the manner in which we established (2.6), we apply (2.5), (3.1), (3.2), and the boundedness of the elements of $\hat{\Sigma}^{-1}$ with probability approaching 1 to obtain

$$\lim_{N \rightarrow \infty} P\left(S_N(x_{-1}^0, \dots, x_t^0) \geq L(N)\right) = 0$$

for all nonrandom sequences $\{L(N)\}$ with $L(N) \rightarrow \infty$. Since

$$S_N(\tilde{x}_{-1}, \dots, \tilde{x}_t) \leq S_N(x_{-1}^0, \dots, x_t^0)$$

it follows that

$$\lim_{N \rightarrow \infty} P\left(S_N(\tilde{x}_{-1}, \dots, \tilde{x}_t) \geq L(N)\right) = 0 \quad (3.3)$$

Now let $T'T = \hat{\Sigma}$ be a square root decomposition obtained by some well defined algorithm (such as a Cholesky decomposition), such that the elements of T are continuous functions of the elements of $\hat{\Sigma}$.

Let

$$\tilde{w}_i = (T^{-1})' \tilde{\bar{x}}_i$$

and

$$w_i = (T^{-1})' \tilde{x}_i$$

Then

$$S_N(\tilde{x}_1, \dots, \tilde{x}_p) = \sum_{i=1}^t \sum_{j=1}^{p+1} m_i (w_{ij} - \tilde{w}_{ij})^2$$

From (3.3) and the fact that $m_i = \pi_i N$, it follows that, for all nonrandom $\{L(N)\}$ with $L(N) \rightarrow \infty$, and for all i and j ,

$$\lim_{N \rightarrow \infty} P\{\sqrt{N} |w_{ij} - \tilde{w}_{ij}| \geq L(N)\} = 0$$

of necessity.

Now suppose that $\tilde{T}'\tilde{T} = \tilde{\Sigma}$ is a square root decomposition of $\tilde{\Sigma}$ obtained by the same algorithm as T . It follows that every element T_{jh} of T converges in probability to the corresponding element \tilde{T}_{jh} of \tilde{T} . Since

$$\tilde{\bar{x}}_i - \tilde{x}_i = \tilde{T}'(w_i - \tilde{w}_i),$$

it follows that

$$\sqrt{N} |\tilde{\bar{x}}_{ij} - \tilde{x}_{ij}| \leq \sum_{h=1}^{p+1} \sqrt{N} |w_{ih} - \tilde{w}_{ih}| |T_{hj}|.$$

Since the values $|T_{jh}|$ are bounded by some nonrandom constant with probability approaching 1, it follows, using (2.5) that

$$\lim_{N \rightarrow \infty} P\left(\sqrt{N} |\tilde{\bar{x}}_{ij} - \tilde{x}_{ij}| \geq L(N)\right) = 0 \quad (3.4)$$

But since

$$\sqrt{N}|\tilde{x}_{ij} - x_{ij}^0| \leq \sqrt{N}|\bar{x}_{ij} - x_{ij}^0| + \sqrt{N}|\bar{x}_{ij} - \tilde{x}_{ij}|,$$

the assertion

$$\lim_{N \rightarrow \infty} P\left(\sqrt{N}|\tilde{x}_{ij} - x_{ij}^0| \geq L(N)\right) = 0$$

is established by virtue of (2.5).

4. THE MULTIVARIATE NORMAL CASE

Now let the observations be multivariate normal. A natural estimate for Σ is

$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^t \sum_{k=1}^{m_i} (x_{ik} - \bar{x}_i)(x_{ik} - \bar{x}_i)'$$

The log-likelihood of the observation is

$$\ell_N = -\left[\frac{(p+1)N}{2}\right] \ln(2\pi) - (N/2) \ln|\Sigma| - \frac{1}{2} \sum_{i=1}^t \sum_{k=1}^{m_i} (x_{ik} - \bar{x}_i)' \Sigma^{-1} (x_{ik} - \bar{x}_i)$$

Now clearly

$$-2 \frac{\partial \ell_N}{\partial \alpha_\mu} = \frac{\partial S_N}{\partial \alpha_\mu} \tag{4.1}$$

for $\mu = 1, 2, \dots, r$, and

$$-2 \frac{\partial \ell_N}{\partial x_{ij}} = \frac{\partial S_N}{\partial x_{ij}} \tag{4.2}$$

for $i = 1, 2, \dots, t$, $j = 1, 2, \dots, p$, when ℓ_N is evaluated at $\Sigma = \hat{\Sigma}$. When (4.1) and (4.2) are evaluated at $\hat{\alpha}$ and the pt estimated nuisance parameters $\tilde{x}_{11}, \dots, \tilde{x}_{1p}, \tilde{x}_{21}, \dots, \tilde{x}_{tp}$, all these

partials are identically 0. Thus the conditions that $\frac{1}{\sqrt{N}} \frac{\partial \ell_N}{\partial \alpha}$ and $\frac{1}{\sqrt{N}} \frac{\partial \ell_N}{\partial \bar{x}_{ij}}$ approach 0 in probability, when evaluated at the least squares estimates and $\hat{\Sigma}$, are trivially fulfilled. It remains to show that $\frac{1}{\sqrt{N}} \frac{\partial \ell_N}{\partial \sigma_{ij}} \rightarrow 0$ in probability when evaluated at these estimates.

Write ℓ_N as

$$\ell_N = - \left[(p+1)N/2 \right] \ln(2\pi) - \frac{1}{2} N \ln |\Sigma| - \frac{1}{2} \text{tr} \left[\Sigma^{-1} \sum_{ik} (x_{ik} - \bar{x}_i) (x_{ik} - \bar{x}_i)' \right]$$

and let $\sum_{ik} (x_{ik} - \bar{x}_i) (x_{ik} - \bar{x}_i)' = B$ for convenience. Let

$E_{ij} = \{e_{kh}\}$ be defined by $e_{kh} = 1$ for $k = i$ and $h = j$ and $e_{kh} = 0$ otherwise. Then for $i < j$,

$$\frac{\partial \ell_N}{\partial \sigma_{ij}} = -N \sigma^{ij} + \frac{1}{2} \text{tr} \left[\Sigma^{-1} \left[E_{ij} + E_{ji} \right] \Sigma^{-1} B \right]$$

or

$$\frac{\partial \ell_N}{\partial \sigma_{ij}} = -N \sigma^{ij} + \frac{1}{2} \text{tr} \left[\left(E_{ij} + E_{ji} \right) \Sigma^{-1} B \Sigma^{-1} \right].$$

Let

$$\Sigma^{-1} B \Sigma^{-1} = A = \{a_{ij}\}.$$

Then $\frac{\partial \ell_N}{\partial \sigma_{ij}}$ is given by

$$-N \sigma^{ij} + \frac{1}{2} \left[a_{ij} + a_{ji} \right] = -N \sigma^{ij} + a_{ij} \quad (4.3)$$

For $i = j$, a similar derivation gives

$$\frac{\partial \ell_N}{\partial \sigma_{ii}} = -\frac{1}{2} N \sigma^{ii} + \frac{1}{2} a_{ii} \quad (4.4)$$

[The derivations of (4.3) and (4.4) follow closely the usual line for the derivation of the maximum likelihood estimates for the underlying parameters of multivariate normal data. Compare Rao (1973).] For a square matrix \underline{C} , let $\text{diag}(\underline{C})$ be a diagonal matrix with diagonal elements equal to the diagonal elements of \underline{C} . Then (4.3) and (4.4) imply

$$\frac{\partial \ell_N}{\partial \underline{\Sigma}} = -N \underline{\Sigma}^{-1} + \frac{1}{2} N \text{diag}(\underline{\Sigma}^{-1}) + \underline{A} - \frac{1}{2} \text{diag}(\underline{A}) \quad (4.5)$$

But \underline{A} may be expressed as

$$\underline{A} = \underline{\Sigma}^{-1} [N \underline{\hat{\Sigma}}] \underline{\Sigma}^{-1} + \underline{\Sigma}^{-1} \left[\sum_i m_i (\bar{\underline{x}}_i - \tilde{\underline{x}}_i) (\bar{\underline{x}}_i - \tilde{\underline{x}}_i)' \right] \underline{\Sigma}^{-1}$$

so that, when (4.5) is evaluated at $\underline{\Sigma} = \underline{\hat{\Sigma}}$ and the least squares solutions,

$$\frac{\partial \ell_N}{\partial \underline{\Sigma}} = \underline{\hat{\Sigma}}^{-1} \left[\sum_i m_i (\bar{\underline{x}}_i - \tilde{\underline{x}}_i) (\bar{\underline{x}}_i - \tilde{\underline{x}}_i)' \right] \underline{\hat{\Sigma}}^{-1} - \frac{1}{2} \text{diag} \left\{ \underline{\hat{\Sigma}}^{-1} \sum_i m_i (\bar{\underline{x}}_i - \tilde{\underline{x}}_i) (\bar{\underline{x}}_i - \tilde{\underline{x}}_i)' \underline{\hat{\Sigma}}^{-1} \right\}$$

Now since the elements of $\underline{\hat{\Sigma}}^{-1}$ converge to those $\underline{\Sigma}^{-1}$ in probability,

it follows that $\frac{1}{\sqrt{N}} \frac{\partial \ell_N}{\partial \underline{\Sigma}}$ converges to 0 elementwise if

$$\frac{1}{\sqrt{N}} \sum_i m_i (\bar{\underline{x}}_i - \tilde{\underline{x}}_i) (\bar{\underline{x}}_i - \tilde{\underline{x}}_i)' \xrightarrow{p} 0$$

elementwise. That is, it is sufficient to demonstrate that

$$\sqrt{N} \sum_i m_i (\bar{x}_{ij} - \tilde{x}_{ij}) (\bar{x}_{ih} - \tilde{x}_{ih}) \xrightarrow{p} 0, \quad j, h = 1, 2, \dots, p+1 \quad (4.6)$$

Now for any prescribed $\epsilon > 0$, choosing $L(N) = \sqrt{\epsilon} N^{\frac{1}{2}}$ in (3.4) will yield

$$\lim_{N \rightarrow \infty} P\left(\sqrt{N}(\bar{x}_{ij} - \bar{x}_{ij})^2 \geq \epsilon\right) = 0$$

Since $\sqrt{N}|\bar{x}_{ij} - \bar{x}_{ij}||\bar{x}_{ih} - \bar{x}_{ih}| \leq \sqrt{N}(\bar{x}_{ij} - \bar{x}_{ij})^2 + \sqrt{N}(\bar{x}_{ih} - \bar{x}_{ih})^2$

(4.6) must be valid. Hence the asymptotic distributions of $\sqrt{N}(\tilde{\alpha} - \alpha_0)$ and $\sqrt{N}(\hat{\alpha} - \alpha_0)$, where $\hat{\alpha}$ is the ML estimate, are equal.

REFERENCES

- Cramér, H. (1946). Mathematical Methods of Statistics. Princeton: Princeton University Press.
- Dolby, G. R. (1972). Generalized least squares and maximum likelihood estimation of non-linear functional relationships. J. Roy. Statist. Soc., 34, 393-400.
- Dolby, G. R. and Freeman T. G. (1975). Functional relationships having many independent variables and errors with multivariate normal distribution. J. Multivariate Anal. 5, 466-479.
- Rao, C. R. (1973). Linear Statistical Inference and Its Applications. New York: John Wiley and Sons.
- Tucker, M. G. (1967). A Graduate Course in Probability. New York: Academic Press.
- Weiss, L. (1971). Asymptotic properties of maximum likelihood estimators in some nonstandard cases. J. Amer. Statist. Assoc., 66, 345-350.