# ON CONSTRAINED MAXIMUM LIKELIHOOD ESTIMATION
# WITH NON-I.I.D. OBSERVATIONS

Martin Crowder

## Summary

Consistency and asymptotic normality of the m.l.e. are examined in the non-i.i.d. case when the parameters are constrained. Inequality constraints are considered as an application.

## 1. Introduction

Sometimes the parameters in a model are not functionally independent, that is, there exist certain relations connecting them. Nevertheless, it may be undesirable to eliminate redundancies for reasons of symmetry, or because the natural interpretation of the parameters would be lost, or because of mathematical intractability.

The data have joint probability distribution of known form depending on a parameter $\theta \in R^k$. The log-likelihood function $l_n(\theta)$, based on $n$ observations, is twice continuously differentiable in a neighbourhood of $\theta_0$, the true parameter, and the information matrix is $B_n = \mathrm{E}[-l_n''(\theta_0)]$. We assume throughout the "usual regularity conditions": $\mathrm{E}[l_n'(\theta_0)] = 0$ and $\mathrm{Var}[l_n'(\theta_0)] = B_n$.

There are $r$ constraints $(0 < r < k)$ of the form $h(\theta) = 0$, where $h = (h_1, \cdots, h_r)^T$. The corresponding $k \times r$ Jacobian matrix, $\{H\}_{ij} = \partial h_j / \partial \theta_i$, is continuous and of full rank $r$ at $\theta_0$; this prevents reduction to a single constraint such as $\sum h_j^2(\theta) = 0$.

Often the parameter $\theta$ is not identifiable in the model, and some of the constraints are needed to achieve identifiability. In a sense made more precise below this corresponds to singularity of $B_n$, and this is encompassed in the results.

Aitchison and Silvey [1] consider the case of i.i.d. observations, where $B_n = nB_1$ is nonsingular. They prove, under certain conditions,

---

asymptotic existence, consistency and asymptotic normality of a solution of the likelihood equations using an extension of the Cramér [3] method.

Silvey [8] deals with the i.i.d. case, but allowing $B_1$ to be singular. He uses the alternative approach to proof of consistency due to Wald [9]. The paper is mainly concerned with testing the truth of $h(\theta_0) = 0$, and consequently $h(\theta_0) \neq 0$ is allowed in the analysis, unlike the case here. Also, more emphasis is placed here on the precise manner in which the constraints may be used to construct a suitable nonsingular version of $B_n$ and the connection with parameter identifiability.

The formal methods in Section 2, deriving asymptotic results for the constrained m.l.e., stem from those of Aitchison and Silvey [1] and Silvey [8], but the modifications are not automatic, as is illustrated by Example 1 below. Also, attention is drawn to the enhanced efficiency of the constrained m.l.e. over that of the unconstrained m.l.e. Section 3 contains some examples, the last of which provides a resolution of a problem posed by Hudson [5]. In Section 4 the connection between parameter identifiability, the information matrix, and constraints is explored, and auxiliary results are derived to support the foregoing theory.

## 2. Asymptotic results for the constrained m.l.e.

The likelihood equations for maximizing $l_n(\theta)$ subject to the constraints are $l_n'(\theta) + H\lambda = 0$ and $h(\theta) = 0$, where $\lambda$ is a vector of $r$ Lagrange multipliers. Application of the mean value theorem at $\theta_0$ to these equations gives

$$(2.1) \qquad l_n'(\theta_0) + l_n''(\theta_0, \theta)(\theta - \theta_0) + H(\theta, \theta)\lambda = 0 \ ,$$

$$(2.2) \qquad H^T(\theta_0, \theta)(\theta - \theta_0) = 0 \ ;$$

the general notation $H(\theta_1, \theta_2)$ is used to convey that the rows of the matrix are evaluated at possibly different points on the line segment joining $\theta_1$ and $\theta_2$.

Suppose that the constraints $h_i(\theta) = 0$ $(i = 1, \cdots, s)$ are sufficient to identify $\theta$ in the sense detailed in Section 4. Let $H$ be partitioned as $(H_1, H_2)$, where $H_1$ is $k \times s$, and let $G = (H_1, H_2 M)$, where $M$ is an arbitrary $(r-s)$-rowed matrix. Note that $G^T(\theta_0, \theta)(\theta - \theta_0) = 0$ from (2.2). Let $B_n^*(\theta) = B_n + G(\theta_0, \theta)C_n G^T(\theta_0, \theta)$, where $C_n$ is positive definite. By Lemma 5 below $B_n^*(\theta)$ is nonsingular, as is $H^T(\theta_0, \theta)B_n^*(\theta)^{-1}H(\theta, \theta)$ for $\theta$ near $\theta_0$ (since $H$ has full rank at $\theta_0$). Let

$$R_n(\theta) = I_k + B_n^*(\theta)^{-1}\{l_n''(\theta_0, \theta) - G(\theta_0, \theta)C_n G^T(\theta_0, \theta)\}$$
$$= B_n^*(\theta)^{-1}\{B_n + l_n''(\theta_0, \theta)\} \ ,$$

$$Q_n(\theta) = B_n^*(\theta)^{-1} H(\theta, \theta) \{H^T(\theta_0, \theta) B_n^*(\theta)^{-1} H(\theta, \theta)\}^{-1},$$

$$w_n(\theta) = [I_k - Q_n(\theta) H^T(\theta_0, \theta)] \{B_n^*(\theta)^{-1} l_n'(\theta_0) + R_n(\theta)(\theta - \theta_0)\}.$$

If $s=0$, so that the model is identified without constraints and $B_n$ is nonsingular, we can take $C_n = 0$ throughout, so $B_n^* = B_n$.

THEOREM 1. *There exists a consistent solution of the likelihood equations if for each $\Delta > 0 \exists \delta \in (0, \Delta]$ s.t.*

$$\text{Prob} [\sup (\theta - \theta_0)^T w_n(\theta) < \delta^2] \to 1 \qquad as \ n \to \infty \ (\sup \ over \ |\theta - \theta_0| = \delta).$$

PROOF. Equation (2.1) is equivalent to

$$l_n'(\theta_0) + \{l_n''(\theta_0, \theta) - G(\theta_0, \theta) C_n G^T(\theta_0, \theta)\}(\theta - \theta_0) + H(\theta, \theta)\lambda = 0,$$

i.e.,

(2.3) $$l_n'(\theta_0) + B_n^*(\theta) \{R_n(\theta) - I_k\}(\theta - \theta_0) + H(\theta, \theta)\lambda = 0.$$

Multiply (2.3) by $H^T(\theta_0, \theta) B_n^*(\theta)^{-1}$ and use (2.2) to obtain

$$H^T(\theta_0, \theta) B_n^*(\theta)^{-1} l_n'(\theta_0) + H^T(\theta_0, \theta) R_n(\theta)(\theta - \theta_0)$$
$$+ H^T(\theta_0, \theta) B_n^*(\theta)^{-1} H(\theta, \theta)\lambda = 0.$$

(2.4) $$\therefore \quad \lambda = - \{H^T(\theta_0, \theta) B_n^*(\theta)^{-1} H(\theta, \theta)\}^{-1} H^T(\theta_0, \theta)$$
$$\cdot \{B_n^*(\theta)^{-1} l_n'(\theta_0) + R_n(\theta)(\theta - \theta_0)\}.$$

On multiplying (2.3) by $B_n^*(\theta)^{-1}$ and substituting for $\lambda$ one obtains

(2.5) $$w_n(\theta) - (\theta - \theta_0) = 0.$$

The pair (2.4), (2.5) are equivalent to the original equations (2.1), (2.2): if (2.5) has a solution, say $\dot\theta_n$, then $\dot\lambda_n$ is defined by (2.4) and these are solutions of (2.1), (2.2).

Now apply the equivalent of Brouwer's fixed point theorem as in Aitchison and Silvey [1]: a solution $\dot\theta_n$ exists, with $|\dot\theta_n - \theta_0| < \delta$, if $(\theta - \theta_0)^T$ times (2.5) is strictly negative for $|\theta - \theta_0| = \delta$.

THEOREM 2. *There exists a consistent solution of the likelihood equations if for each $\Delta > 0 \exists K < \infty$, $d > 0$, $\delta > 0$, $p \in [0, 1]$ s.t. $\delta < \min(d, \Delta)$ and*
( i ) $|Q_n(\theta) H^T(\theta_0, \theta)| \leq K$ *for* $|\theta - \theta_0| \leq d$ *and all* $n$,
( ii ) $\text{Prob} [|B_n^*(\theta)^{-1} l_n'(\theta_0)| < p\delta(K + \sqrt{k})^{-1}] \to 1$ *as* $n \to \infty$ *when* $|\theta - \theta_0| = \delta$,
(iii) $\text{Prob} [|R_n(\theta)| < (1-p)(K + \sqrt{k})^{-1}] \to 1$ *as* $n \to \infty$ *when* $|\theta - \theta_0| = \delta$.

PROOF. This is based on the following, which holds for $|\theta - \theta_0| = \delta$:

$$|(\theta - \theta_0)^T w_n(\theta)| \leq \delta |w_n(\theta)| \leq \delta |I_k - Q_n(\theta) H^T(\theta_0, \theta)|$$

$$\cdot \{|B_n^*(\theta)^{-1}l_n'(\theta_0)| + \delta |R_n(\theta)|\} \ .$$

From (i) $\sup |I_k - Q_n(\theta)H^T(\theta_0, \theta)| \leq K + \sqrt{k}$, $\sup$ over $|\theta - \theta_0| = \delta$. Hence

$$\text{Prob} [\sup (\theta - \theta_0)^T w_n(\theta) < \delta^2] \geq \text{Prob} [\sup |w_n(\theta)| < \delta]$$

$$(2.6) \qquad \geq \text{Prob} [\{\sup |B_n^*(\theta)^{-1}l_n'(\theta_0)| < p\delta(K + \sqrt{k})^{-1}\}$$

$$\cap \{\sup |R_n(\theta)| < (1-p)(K + \sqrt{k})^{-1}\}] \ .$$

Now $R_n(\theta)$ is $\theta$-continuous and $\{\theta : |\theta - \theta_0| = \delta\}$ is compact, so (ii) implies Prob $[\sup |R_n(\theta)| < (1-p)(K + \sqrt{k})^{-1}] \to 1$; similarly for $B_n^*(\theta)^{-1}l_n'(\theta_0)$. Hence $(2.6) \to 1$ and the criterion of Theorem 1 is fulfilled.

Theorem 2 gives a breakdown of the single condition of Theorem 1 into more manageable bits. The reduction depends critically upon (i), that $Q_n(\theta)H^T(\theta_0, \theta)$ should be bounded uniformly in both $n$ and $\theta$, and this can fail as in Example 1 below. However, in many situations the non-zero eigenvalues of $B_n$ will all have the same order of magnitude, say $O(b_n)$, and then the choice $C_n = b_n I$ will ensure that $B_n^*(\theta)b_n^{-1}$ is bounded uniformly in $n$ (Lemma 5 below), and (i) follows; the i.i.d. case is covered since there $b_n = n$. In (ii) $E |B_n^*(\theta)^{-1}l_n'(\theta_0)|^2 = \text{trace} [B_n^*(\theta)^{-1}B_n B_n^* (\theta)^{-1}]$ and straightforward choices for $C_n$ and $M$ will often suffice to show that this tends to zero, so that $B_n^*(\theta)^{-1}l_n'(\theta_0) \underset{p}{\to} 0$. In (iii) $E[R_n (\theta_0)] = 0$ so stability of $R_n(\theta_0)$, in the sense that $|R_n(\theta_0)|$ is small, and a certain continuity of $R_n(\theta)$ are called for.

In the nonsingular case $B_n^*(\theta) = B_n$ and (ii) reduces to trace $B_n^{-1} \to 0$, which would normally hold in practice. Also, (iii) becomes a condition requiring $B_n^{-1}l_n''(\theta_0, \theta) + I_k$ to be "small in probability" near $\theta_0$; such a property is discussed in Crowder [4].

The following two theorems concern the asymptotic distribution of the constrained m.l.e. We use the notations $H_0 = H(\theta_0, \theta_0)$, $Q_{n0} = Q_n(\theta_0)$, $B_{n0}^* = B_n^*(\theta_0)$, and $\underset{d}{\sim}$ to mean "is asymptotically distributed as".

THEOREM 3. *Suppose that*
( i ) *there is a consistent solution $(\dot{\theta}_n, \dot{\lambda}_n)$ of the likelihood equations,*
( ii ) $l_n'(\theta_0) \underset{d}{\sim} N(0, B_n)$,
(iii) $R_n(\dot{\theta}_n) \underset{p}{\to} 0$, $Q_n(\dot{\theta}_n) - Q_{n0} \underset{p}{\to} 0$, $|Q_{n0}| \leq K < \infty \ \forall n$.

*Then* $\begin{pmatrix} \dot{\theta}_n \\ \dot{\lambda}_n \end{pmatrix} \underset{d}{\sim} N\left( \begin{pmatrix} \theta_0 \\ 0 \end{pmatrix}, \begin{pmatrix} V_{n1} & 0 \\ 0 & V_{n2} \end{pmatrix} \right)$ *where* $V_{n1} = (I_k - Q_{n0}H_0^T)B_{n0}^{*-1}$, $V_{n2} = Q_{n0} \cdot B_n Q_{n0}^T$.

PROOF. Let $\mu = H^T(\theta_0, \theta)B_n^*(\theta)^{-1}H(\theta, \theta)\lambda$ and $D_n(\theta) = \begin{pmatrix} I_k - R_n(\theta) & -Q_n(\theta) \\ H^T(\theta_0, \theta) & 0 \end{pmatrix}$. Then the likelihood equations (2.2), (2.3) can be written, after multiplying the latter by $B_n^*(\theta)^{-1}$, as

$$(2.7) \qquad D_n(\theta)\begin{pmatrix}\theta-\theta_0\\ \mu\end{pmatrix}=\begin{pmatrix}B_n^*(\theta)^{-1}l_n'(\theta_0)\\ 0\end{pmatrix}.$$

Let $D_{n0}=\begin{pmatrix}I_k & -Q_{n0}\\ H_0^T & 0\end{pmatrix}$, so $D_{n0}^{-1}=\begin{pmatrix}I_k-Q_{n0}H_0^T & Q_{n0}\\ -H_0^T & I_r\end{pmatrix}$ and

$$D_{n0}^{-1}D_n(\theta)=\begin{pmatrix}I_k-(I_k-Q_{n0}H_0^T)R_n(\theta)+Q_{n0}(H(\theta_0,\theta)-H_0)^T\\ H_0^T R_n(\theta)+(H(\theta_0,\theta)-H_0)^T\end{pmatrix}$$
$$\begin{pmatrix}-(Q_n(\theta)-Q_{n0})-Q_{n0}(H(\theta_0,\theta)-H_0)^T Q_n(\theta)\\ I_r-(H(\theta_0,\theta)-H_0)^T Q_n(\theta)\end{pmatrix}.$$

Under (i) and (iii) $D_{n0}^{-1}D_n(\dot\theta_n)\underset{p}{\to}I_{k+r}$. Hence, from (2.7),

$$\begin{pmatrix}\dot\theta_n-\theta_0\\ \dot\mu_n\end{pmatrix}\underset{d}{\sim}D_{n0}^{-1}\begin{pmatrix}B_{n0}^{*-1}l_n'(\theta_0)\\ 0\end{pmatrix},$$

using a matrix version of a theorem of Cramér ([3], § 20.6). By (ii) the asymptotic distribution of the right-hand side is $N(0, V_n)$, where

$$V_n=D_{n0}^{-1}\begin{pmatrix}B_{n0}^{*-1}B_n B_{n0}^{*-1} & 0\\ 0 & 0\end{pmatrix}(D_{n0}^{-1})^T.$$

The forms given above for $V_{n1}$, $V_{n2}$ follow after some matrix manipulation.

Condition (i) of Theorem 3 seems natural, and the behaviour in (ii) is discussed in Crowder [4] based on the work of Brown [2] and Scott [7]. Condition (iii) is designed to support $D_{n0}^{-1}D_n(\dot\theta_n)\underset{p}{\to}I_{k+r}$ in the proof; alternatives are possible.

The variance matrices $V_{n1}$, $V_{n2}$ are, in fact, independent of the particular choice of $M$, $C_n$; this must be so since $\dot\theta_n$, $\dot\lambda_n$ are defined as solutions of (2.1), (2.2) which do not involve $M$, $C_n$, and the assertion also follows from Lemma 6, Section 4. Since $V_{n1}H_0=0$ variation orthogonal to the constraint surface is suppressed, i.e., $\mathrm{Var}\,(a^T\dot\theta_n)\sim 0$ whenever $a=H_0 b$ for some $b$. Theorem 4 elaborates this point.

Let $\hat\theta_n$ be the "unconstrained" m.l.e., satisfying the identifiability constraints but not the rest. Then $\hat\theta_n$ will be consistent for $\theta_0$ but, as the next theorem shows, less efficient than $\dot\theta_n$. We use the notation $Q_{n0}'=B_{n0}^{*-1}H_{10}(H_{10}^T B_{n0}^{*-1}H_{10})^{-1}$, $H_{10}=H_1(\theta_0,\theta_0)$, $H_{20}=H_2(\theta_0,\theta_0)$.

THEOREM 4. *Let $a^T\theta_0$ be a linear parametric function to be estimated, where $a$ is $k\times 1$. Then* $\mathrm{Var}\,(a^T\hat\theta_n)\geqq\mathrm{Var}\,(a^T\dot\theta_n)$ *with equality iff* $H_{20}^T(I_k-Q_{n0}'H_{10}^T)B_{n0}^{*-1}a=0$.

PROOF. Under the conditions assumed for consistency and asymptotic normality of $\dot\theta_n$, $\hat\theta_n$ will also be so, with asymptotic variance $V_{n1}'$ $=(I_k-Q_{n0}'H_{10}^T)B_{n0}^{*-1}$, from Theorem 3. Now

$$\mathrm{Var}\,(a^T\hat{\theta}_n) - \mathrm{Var}\,(a^T\dot{\theta}_n) \sim a^T(Q_{n0}H_0^T - Q'_{n0}H_{10}^T)B_{n0}^{*-1}a$$
$$= a^T B_{n0}^{*-1}\{H_0(H_0^T B_{n0}^{*-1}H_0)^{-1}H_0^T - H_{10}(H_{10}^T B_{n0}^{*-1}H_{10})^{-1}H_{10}^T\}B_{n0}^{*-1}a\ .$$

Also,

$$(H_0^T B_{n0}^{*-1}H_0)^{-1} = \begin{pmatrix} H_{10}^T B_{n0}^{*-1}H_{10} & H_{10}^T B_{n0}^{*-1}H_{20} \\ H_{20}^T B_{n0}^{*-1}H_{10} & H_{20}^T B_{n0}^{*-1}H_{20} \end{pmatrix}$$
$$= \begin{pmatrix} F^{-1} + F^{-1}H_{10}^T B_{n0}^{*-1}H_{20}E^{-1}H_{20}^T B_{n0}^{*-1}H_{10}F^{-1} & -F^{-1}H_{10}^T B_{n0}^{*-1}H_{20}E^{-1} \\ -E^{-1}H_{20}^T B_{n0}^{*-1}H_{10}F^{-1} & E^{-1} \end{pmatrix}$$

where

$$E = H_{20}^T(B_{n0}^{*-1} - B_{n0}^{*-1}H_{10}F^{-1}H_{10}^T B_{n0}^{*-1})H_{20}\ , \qquad F = H_{10}^T B_{n0}^{*-1}H_{10}\ .$$

Hence

$$\mathrm{Var}\,(a^T\hat{\theta}_n) - \mathrm{Var}\,(a^T\dot{\theta}_n) \sim \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}^T (H_0^T B_{n0}^{*-1}H_0)^{-1}\begin{pmatrix} b_1 \\ b_2 \end{pmatrix} - b_1^T(H_{10}^T B_{n0}^{*-1}H_{10})^{-1}b_1$$
$$= |E^{-1/2}(H_{20}^T B_{n0}^{*-1}H_{10}F^{-1}b_1 - b_2|^2 \geqq 0\ ,$$

where

$$b_1 = H_{10}^T B_{n0}^{*-1}a\ , \qquad b_2 = H_{20}^T B_{n0}^{*-1}a\ .$$

The difference is zero iff $H_{20}^T B_{n0}^{*-1}H_{10}F^{-1}b_1 = b_2$, i.e., iff $H_{20}^T P_n B_{n0}^{*-1}a = 0$ where $P_n = I_k - B_{n0}^{*-1}H_{10}F^{-1}H_{10}^T = I_k - Q'_{n0}H_{10}^T$.

The condition for equality of variances in Theorem 4 may be interpreted as follows. The matrix $P_n = (I_k - Q'_{n0}H_{10}^T)$ is idempotent, non-symmetric and satisfies $H_{10}^T P_n = 0$. It therefore represents non-orthogonal projection into the subspace orthogonal to the columns of $H_{10}$, i.e., the $H_1$-constraint surface. The condition is that when $B_{n0}^{*-1}a$ is so projected, the resulting vector also lies within the $H_2$-constraint surface; since $P_n B_{n0}^{*-1} = V'_{n1}$ is independent of $M$ and $C_n$, the set of vectors $a$ satisfying the condition is likewise independent of $M$ and $C_n$. In particular, this will always hold when there are only identifiablility constraints, i.e., when $H_2$ is null. In the identified case, when there are no identifiability constraints, $H_1$ is null, $B_{n0}^* = B_n$, $P_n = I_k$ and the condition reduces to $H_{20}^T B_n^{-1}a = 0$, i.e., $B_n^{-1}a$ lies within the constraint surface.

## 3. Examples

*Example* 1. We give a very simple, non-pathological case in which condition (i) of Theorem 2 fails, but where it is possible to verify the criterion in Theorem 1 directly. It thus shows that Theorem 1 is more

sensitive than Theorem 2, and that straightforward adaptation of the conditions for the i.i.d. case is not sufficient.

Suppose that $y_t$ $(t=1,\cdots,n)$ are independent, with distributions $N(\theta_1+\theta_2 t^{-\beta}, 1)$, where $0<\beta<1/2$ and $\beta$ is known, and that the constraint is $\theta_1=0$. Thus $k=2$, $r=1$, $s=0$, $B_n^*=B_n$,

$$-l_n''(\theta_0, \theta) = B_n = \begin{pmatrix} n & \sum t^{-\beta} \\ \sum t^{-\beta} & \sum t^{-2\beta} \end{pmatrix}, \quad H(\theta_0, \theta) = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad Q_n(\theta) = \begin{pmatrix} 1 \\ -r_n \end{pmatrix},$$

where $r_n = \sum t^{-\beta}/\sum t^{-2\beta} = O(n^\beta)$; also $R_n(\theta) = 0$ so $w_n = \begin{pmatrix} 0 & 0 \\ r_n & 1 \end{pmatrix} B_n^{-1} l_n'(\theta_0)$.

Theorem 2 fails because $Q_n(\theta) H^T(\theta_0, \dot\theta)$ has an element $O(n^\beta)$. However, $B_n^{-1/2} l_n'(\theta_0)$ is constant in mean square so Theorem 1 will work if $\begin{pmatrix} 0 & 0 \\ r_n & 1 \end{pmatrix}$ $\cdot B_n^{-1/2} \to 0$, which can be verified. In fact, the constrained m.l.e. $\hat\theta_{2n} = \sum y_t t^{-\beta}/\sum t^{-2\beta}$ is unbiased with variance $(\sum t^{-2\beta})^{-1} \to 0$.

*Example 2. Inequality constraints.* Suppose the model has log-likelihood $l_n(\phi)$ and is identified in the parameters $\phi$, i.e., $B_{\phi\phi} = \mathrm{E}[-\partial^2 l_n/\partial\phi^2]$ is nonsingular. General inequality constraints can be written as $h_i(\phi) \geqq 0$ $(i=1,\cdots,r)$ and can be converted to equality constraints by introducing $r$ parameters $\psi$ and writing $h_i(\phi) = \psi_i^2$. Thus $\theta = \begin{pmatrix} \phi \\ \psi \end{pmatrix}$ and $H = \begin{pmatrix} H_\phi \\ -2\Psi \end{pmatrix}$ where $(H_\phi)_{ij} = \partial h_j/\partial\phi_i$ and $\Psi = \mathrm{diag}(\psi_1,\cdots,\psi_r)$. All the constraints are required for identifiability so $s=r$ and $H_\lambda$ is null. We have $l_n'(\theta) = \begin{pmatrix} l_\phi \\ 0 \end{pmatrix}$, $l_n''(\theta) = \begin{pmatrix} l_{\phi\phi} & 0 \\ 0 & 0 \end{pmatrix}$ and $B_n = \begin{pmatrix} B_{\phi\phi} & 0 \\ 0 & 0 \end{pmatrix}$ where $(l_\phi)_i = \partial l_n/\partial\phi_i$ and $(l_{\phi\phi})_{ij} = \partial^2 l_n/\partial\phi_i\partial\phi_j$. For convenience we will use the notation $H_{00} = H(\theta_0, \theta_0)$, $H_{01} = H(\theta_0, \theta)$, $H_{11} = H(\theta, \theta)$, with similar definitions for $H_\phi$ and $\Psi$. Now $B_n^*(\theta) = B_n + H_{01} C_n H_{01}^T$, so

$$B_n^*(\theta) = \begin{pmatrix} B_{\phi\phi} + H_{\phi01} C_n H_{\phi01}^T & -2H_{\phi01} C_n \Psi_{01}^T \\ -2\Psi_{01} C_n H_{\phi01}^T & 4\Psi_{01} C_n \Psi_{01}^T \end{pmatrix},$$

$$B_n^*(\theta)^{-1} = \begin{pmatrix} B_{\phi\phi}^{-1} & (1/2) B_{\phi\phi}^{-1} H_{\phi01} \Psi_{01}^{-1} \\ (1/2)\Psi_{01}^{-1} H_{\phi01}^T B_{\phi\phi}^{-1} & (1/4)\Psi_{01}^{-1}(C_n^{-1} + H_{\phi01}^T B_{\phi\phi}^{-1} H_{\phi01})\Psi_{01}^{-1} \end{pmatrix}.$$

The conditions of Theorem 2 may be checked as follows:
(i) $B_n^*(\theta)^{-1} H(\theta, \theta)$

$$= \begin{pmatrix} B_{\phi\phi}^{-1}[H_{\phi11} - H_{\phi01}\Psi_{01}^{-1}\Psi_{11}] \\ (1/2)\Psi_{01}^{-1} H_{\phi01}^T B_{\phi\phi}^{-1}[H_{\phi11} - H_{\phi01}\Psi_{01}^{-1}\Psi_{11}] - (1/2)\Psi_{01}^{-1} C_n^{-1}\Psi_{01}^{-1}\Psi_{11} \end{pmatrix}. \text{ Hence}$$

$H^T(\theta_0, \theta) B_n^*(\theta)^{-1} H(\theta, \theta)$

$$= C_n^{-1}\Psi_{01}^{-1}\Psi_{11} \text{ and } Q_n(\theta) = \begin{pmatrix} B_{\phi\phi}^{-1} D C_n \\ (1/2)\Psi_{01}^{-1} H_{\phi01}^T B_{\phi\phi}^{-1} D C_n - (1/2)\Psi_{01}^{-1} \end{pmatrix},$$

where $D = H_{\phi11}\Psi_{11}^{-1}\Psi_{01} - H_{\phi01}$. Uniform boundedness of $Q_n(\theta) H^T(\theta_0, \theta)$ will

then follow from that of $B_{\phi\phi}^{-1}DC_n$, and this will obtain in the "usual case" where the eigenvalues of $B_{\phi\phi}$ all have the same order of magnitude; if they are all $O(b_n)$ the choice $C_n = b_n I$ suffices.

(ii)  Trace $[B_n^*(\theta)^{-1}B_n B_n^*(\theta)^{-1}] = \text{trace } (B_{\phi\phi}^{-1} + (1/4)\Psi_{01}^{-1}H_{\phi01}^T B_{\rho\phi}^{-1}H_{\phi01}\Psi_{01}^T) \to 0$ iff $B_{\phi\phi}^{-1} \to 0$, as would be expected in practice.

(iii)  $R_n(\theta) = B_n^*(\theta)^{-1}\{B_n + l_n''(\theta_0, \theta)\} = \begin{pmatrix} I_{k-r} \\ (1/2)\Psi_{01}^{-1}H_{\rho01}^T \end{pmatrix}[I_{k-r} + B_{\phi\phi}^{-1}l_{\phi\phi}(\phi_0, \phi)].$

Conditions under which $B_{\rho\phi}^{-1}l_{\phi\phi}(\phi_0, \phi) + I_{k-r}$ is "small in probability", for an identified model, are discussed in Crowder [4].

Asymptotic normality of $\dot{\theta}_n$ obtains if the conditions of Theorem 3 hold, and these have been essentially covered above. The asymptotic variance of $\dot{\theta}_n$ is $V_{n1} = (I_k - Q_{n0}H_0^T)B_{n0}^{*-1}$

$$= \begin{pmatrix} B_{\phi\phi}^{-1} & (1/2)B_{\phi\phi}^{-1}H_0\Psi_0^{-1} \\ (1/2)\Psi_0^{-1}H_{\phi0}^T B_{\phi\phi}^{-1} & (1/4)\Psi_0^{-1}H_{\rho0}^T B_{\phi\phi}^{-1}H_{\phi0}\Psi_0^{-1} \end{pmatrix}.$$

In particular, the asymptotic variance of $\dot{\phi}_n$ is $B_{\phi\phi}^{-1}$, the same as that of the unconstrained m.l.e. $\hat{\phi}_n$, there being no gain in efficiency since $H_2$ is null, as noted in Section 2. The reason for this is that $\dot{\phi}_n$ and $\hat{\phi}_n$ are asymptotically equivalent, as a consequence of the assumption that $\theta_0$ is an interior point of the constraint space. For the boundary case a different analysis is required, see Moran [6].

*Example 3.  A class of constrained regression problems.* Hudson [5] discusses polynomial regression where the regression function is constrained be non-positive, non-negative, non-increasing, non-decreasing, convex, or concave over a specified interval. Suppose $E(y|x) = \beta(x)$ for $x \in (a, b)$, where $\beta(x)$ is a polynomial in $x$, then the constraint is of the form $\beta(x) \leqq 0$, $\beta(x) \geqq 0$, $\beta'(x) \leqq 0$, $\beta'(x) \geqq 0$, $\beta''(x) \leqq 0$, or $\beta''(x) \geqq 0$ on $(a, b)$. Hudson remarks that such a constraint generates an infinite set of linear inequalities on the regression coefficients, which makes conventional statistical inference difficult. His paper concentrates on the computational aspect of fitting the constrained regression by least-squares. We show now that such a problem can be accomodated to some extent within the present framework, and thus derive conventional (asymptotic) statistical inference.

The constrained polynomial function ($\beta(x)$, or $\beta'(x)$, or $\beta''(x)$), say of degree $p$, is expressible as $\gamma(x) = \gamma_0(x - \alpha_0) \prod (x^2 - 2\alpha_i x + \gamma_i)$; the product of quadratic factors runs over $i = 1, \cdots, [p/2]$, and the linear factor $(x - \alpha_0)$ is absent if $p$ is even. The constraint is equivalent to (1) $\gamma(x_0) \geqq 0$ (or $\leqq 0$) for some $x_0 \in (a, b)$, and (2) $\gamma(x)$ has no real roots in $(a, b)$. Condition (2) may be broken down as (2a) $\alpha_0 \notin (a, b)$ (if $p$ odd), and (for each $i$) *either* (2b) $\alpha_i^2 < \gamma_i$ (complex roots) *or* (2c) $\alpha_i^2 \geqq \gamma_i$ and $\alpha_i \pm (\alpha_i^2 - \gamma_i)^{1/2} \notin (a, b)$. Condition (2a) can be expressed as $(\alpha_0 - a)(\alpha_0 - b) - \psi_0^2 = 0$,

introducing a parameter $\phi_0$, and "either (2b) and (2c)" similarly as

$$(4.1) \quad (\alpha_i^2 - \gamma_i) \min \{(\alpha_i^2 - \gamma_i), [\alpha_i + (\alpha_i^2 - \gamma_i)_+^{1/2} - a][\alpha_i + (\alpha_i^2 - \gamma_i)_+^{1/2} - b],$$
$$[\alpha_i - (\alpha_i^2 - \gamma_i)_+^{1/2} - a][\alpha_i - (\alpha_i^2 - \gamma_i)_+^{1/2} - b]\} - \phi_i^2 = 0 \;,$$

where $(z)_+$ denotes $\max(0, z)$. These constraint functions will be continuously differentiable in the neighbourhood of a true parameter which does not lie on an implied boundary; for example, $\min(z^2, z^3)$ has a continuous derivative everywhere except at $z=1$. The regression function $\beta(x)$ is obtainable from $\gamma(x)$, by integration if necessary. Thus, as soon as an error distribution is specified for the regression, the problem comes under the general outline of Example 2.

For illustration consider fitting a cubic curve which is monotone increasing on $(a, b)$. Corresponding to Hudson's least-squares approach we will take the usual Normal, homoscedastic model in which the observations $y_t$ $(t=1,\cdots,n)$ are independent $N(\beta(x_t), \sigma^2)$ variates. Now $\beta'(x) = \gamma(x) = \gamma_0(x^2 - 2\alpha_1 x + \gamma_1)$, so $\beta(x) = \alpha_0 + \gamma_0(x^3/3 - \alpha_1 x^2 + \gamma_1 x)$. Constraint (1) may be written as $h_1(\alpha_1, \gamma_0, \gamma_1, \phi_0) \equiv \beta'(x_0) - \phi_0^2 = 0$, where $x_0 \in (a, b)$ is specified. Constraint (2) is $h_2(\alpha_1, \gamma_1, \phi_1) = 0$, where $h_2$ has the form (4.1). The situation is thus covered by Example 2 with $\boldsymbol{\phi} = (\alpha_0, \alpha_1, \gamma_0, \gamma_1, \sigma)^T$, $\boldsymbol{\psi} = (\phi_0, \phi_1)^T$. We will assume that the conditions for consistency and asymptotic Normality of $\hat{\boldsymbol{\theta}}_n$ hold, these being the standard ones for nonlinear regression. The asymptotic variance of $\hat{\boldsymbol{\phi}}_n$, the constrained estimator, is $\boldsymbol{B}_{\phi\phi}^{-1}$, the same as that of the unconstrained m.l.e.

## 4. Identifiability, and the information and constraint matrices

Some particulars are listed in this section connecting parameter identifiability, the information matrix, and the constraints. Although some of this material is familiar in non-rigorous terms, it does not seem readily available elsewhere in suitable form; it is needed to support the work in Section 2.

We will say that there is local non-identifiability at $\boldsymbol{\theta}_0$, in the direction of vector $\boldsymbol{u}$, when $\boldsymbol{u}^T l_n'(\boldsymbol{\theta}_0) = 0$ a.s. In this case the likelihood $l_n(\boldsymbol{\theta})$ has zero derivative along $\boldsymbol{u}$ at $\boldsymbol{\theta}_0$ for all possible data, so $\boldsymbol{u}$ is tangent at $\boldsymbol{\theta}_0$ to a contour of a.s. constant likelihood. Since $E|\boldsymbol{u}^T l_n'(\boldsymbol{\theta}_0)|^2 = \boldsymbol{u}^T \boldsymbol{B}_n \boldsymbol{u}$, $\boldsymbol{u}^T l_n'(\boldsymbol{\theta}_0) = 0$ a.s. iff $\boldsymbol{B}_n \boldsymbol{u} = 0$, i.e., $\boldsymbol{u} \in U_n$, the null space of $\boldsymbol{B}_n$. Every $\boldsymbol{\theta}$ has representation $\boldsymbol{\theta}_u + (\boldsymbol{\theta} - \boldsymbol{\theta}_u)$, where $\boldsymbol{\theta}_u \in U_n$ and $\boldsymbol{\theta} - \boldsymbol{\theta}_u \perp U_n$, and identification is achieved by selecting a particular $\boldsymbol{\theta}_u$.

The constraints, $h_i(\boldsymbol{\theta}) = 0$ $(i=1,\cdots,r)$, are linearized in (2.2) as $\boldsymbol{H}^T(\boldsymbol{\theta}_0, \boldsymbol{\theta})(\boldsymbol{\theta} - \boldsymbol{\theta}_0) = 0$. In Lemma 1 $\boldsymbol{H}_1$ represents the leading $k \times s$ submatrix of the partition $\boldsymbol{H}(\boldsymbol{\theta}_0, \boldsymbol{\theta}) = [\boldsymbol{H}_1(\boldsymbol{\theta}_0, \boldsymbol{\theta}), \boldsymbol{H}_2(\boldsymbol{\theta}_0, \boldsymbol{\theta})]$, $0 \leqq s \leqq r$.

LEMMA 1. *If constraints $h_i(\boldsymbol{\theta}) = 0$ $(i=1,\cdots,s)$ are sufficient for*

*local identification at $\boldsymbol{\theta}_0$ then rank $(\boldsymbol{B}_n, \boldsymbol{H}_1) = k$.*

PROOF. If rank $(\boldsymbol{B}_n, \boldsymbol{H}_1) < k$ there exist non-zero $\boldsymbol{u} \in U_n$ s.t. $\boldsymbol{H}_1^T \boldsymbol{u} = 0$. In that case the constraint $\boldsymbol{H}_1^T(\boldsymbol{\theta} - \boldsymbol{\theta}_0) = 0$ satisfied by a particular $\boldsymbol{\theta}$ would also be satisfied by $\boldsymbol{\theta} + \boldsymbol{u}$, and identification would not be achieved.

Note that Lemma 1 does not show that $[\boldsymbol{B}_n, \boldsymbol{H}_1(\boldsymbol{\theta}_0, \boldsymbol{\theta}_0)]$ has full rank. Only in the case of linear constraints will $\boldsymbol{H}_1(\boldsymbol{\theta}_0, \boldsymbol{\theta})$ be independent of $\boldsymbol{\theta}$.

The following lemmas are stated in general terms but with notation corresponding closely to their application here. Lemmas 2, 3 and 4 serve Lemma 5 which shows that the matrix $\boldsymbol{B}_n^*$ is nonsingular and gives a bound for its minimum eigenvalue.

LEMMA 2. *Suppose $U$ is a vector subspace of $R^k$ spanned by the columns of $A$ $(k \times s)$ of rank $s$, $0 < s \leq k$. Let $a = \inf \boldsymbol{u}^T A A^T \boldsymbol{u}$, inf over $\boldsymbol{u} \in U$ with $|\boldsymbol{u}| = 1$. Then $a > 0$.*

PROOF. Since $\boldsymbol{u} \in U$, $\boldsymbol{u} = A\boldsymbol{v}$ for some $\boldsymbol{v}$ and $1 = |\boldsymbol{u}|^2 = |A\boldsymbol{v}|^2 = \boldsymbol{v}^T A^T A \boldsymbol{v} = |(A^T A)^{1/2} \boldsymbol{v}|^2$.

$$a = \inf \boldsymbol{v}^T A^T (A A^T) A \boldsymbol{v} = \inf \{(A^T A)^{1/2} \boldsymbol{v}\}^T (A^T A) \{(A^T A)^{1/2} \boldsymbol{v}\} ,$$

inf over $\boldsymbol{v}$ s.t. $|(A^T A)^{1/2} \boldsymbol{v}| = 1$. Thus $a$ is the smallest eigenvalue of $A^T A$, which is positive definite since $A$ has full rank.

LEMMA 3. *Suppose $B$ $(k \times k)$ is positive semidefinite with null space $U$ of dimension $s$, $0 \leq s < k$. Let $b = \inf \boldsymbol{u}^T B \boldsymbol{u}$, inf over $\boldsymbol{u} \perp U$ with $|\boldsymbol{u}| = 1$. Then $b > 0$.*

PROOF. The eigenvectors $e_1, \cdots, e_k$ of $B$ span $R^k$ and correspond to eigenvalues $\lambda_1 = \cdots = \lambda_s = 0$, $0 < \lambda_{s+1} \leq \cdots \leq \lambda_k$. Let $E_2$ be $k \times (k-s)$ with columns $e_{s+1}, \cdots, e_k$, and $\Lambda_2 = \text{diag}(\lambda_{s+1}, \cdots, \lambda_k)$. Then $\boldsymbol{u} = E_2 \boldsymbol{v}$ for some $\boldsymbol{v}$, and $1 = |\boldsymbol{u}| = |E_2 \boldsymbol{v}| = |\boldsymbol{v}|$ since $E_2^T E_2 = I_{k-s}$. Hence $b = \inf_{|\boldsymbol{v}|=1} (E_2 \boldsymbol{v})^T \cdot B(E_2 \boldsymbol{v}) = \inf_{|\boldsymbol{v}|=1} \boldsymbol{v}^T \Lambda_2 \boldsymbol{v} = \lambda_{s+1} > 0$.

LEMMA 4. *Let $B$ be positive semidefinite with null space $U$ and let the columns of $B$ and $H_1$ together span $R^k$. Then $H_1$ is expressible as $A_1 + BA_2$ where $BA_1 = 0$ and the columns of $A_1$ span $U$.*

PROOF. Each column of $\boldsymbol{H}_1$ can be expressed uniquely as $\boldsymbol{a}_1 + \boldsymbol{a}_2$ where $\boldsymbol{a}_1 \in U$ and $\boldsymbol{a}_2 \perp U$, i.e., there exist matrices $A_1$ (unique) and $A_2$ (non-unique) such that $\boldsymbol{H}_1 = A_1 + BA_2$, and $BA_1 = 0$.

For arbitrary $\boldsymbol{z} \in R^k$ we have $\boldsymbol{z} = B\boldsymbol{v}_1 + \boldsymbol{H}_1 \boldsymbol{v}_2$ for some $\boldsymbol{v}_1$, $\boldsymbol{v}_2$. Thus $\boldsymbol{z} = B\boldsymbol{v}_3 + A_1\boldsymbol{v}_2$, where $\boldsymbol{v}_3 = \boldsymbol{v}_1 + A_2\boldsymbol{v}_2$, so the columns of $B$ and $A_1$ together span $R^k$. If $\boldsymbol{z} \in U$ then $0 = B\boldsymbol{z} = B^2\boldsymbol{v}_3$, so $\boldsymbol{v}_3^T B^2 \boldsymbol{v}_3 = 0$ and thus $B\boldsymbol{v}_3 = 0$. It

follows that $z = A_1 v_2$, so the columns of $A_1$ span $U$.

LEMMA 5.   *Let $B$, $U$, $b$ and $H_1$ be defined as in Lemmas 2 and 3. Let $B^* = B + GCG^T$ where $C$ is positive definite with minimum eigenvalue $c$, and $G$ is partitioned as $(H_1, G_2)$ in which $G_2$ is arbitrary. Then $\inf z^T B^* z \geq \min \{b\varepsilon^2, ca(1 - \varepsilon^2)\}$, inf over $|z| = 1$, for some $a > 0$ and $\varepsilon \in (0, 1)$.*

PROOF.   Let $z = z_1 + z_2$ where $z_1 \in U$, $z_2 \perp U$ and $|z| = 1$. Then $z^T B^* z = z_2^T B z_2 + z^T GCG^T z \geq b|z_2|^2 + c|G^T z|^2$ (Lemma 2). Choose $\varepsilon \in (0, 1)$ s.t. $\varepsilon^2 \leq a(a + 16|H_1|^2)^{-1}$; note that, since $|H_1|^2 = |A_1|^2 + |BA_2|^2 \geq |BA_2|^2$ (Lemma 3), $\varepsilon^2 \leq a(a + 16|BA_2|^2)^{-1}$ which implies $2\varepsilon|A_2^T B| \leq \{a(1 - \varepsilon^2)\}^{1/2}/2$. If $|z_2| \geq \varepsilon$ then $z^T B^* z \geq b\varepsilon^2$. Otherwise $|z_1|^2 = 1 - |z_2|^2 \geq 1 - \varepsilon^2$ and

$$|G^T z|^2 = |H_1^T z|^2 + |G_2^T z|^2 \geq |(A_1 + BA_2)^T z|^2 \quad \text{(Lemma 3)}$$
$$= |A_1^T z_1 + A_2^T B z_2|^2 \geq (|A_1^T z_1| - |A_2^T B z_2|)^2 \geq |A_1^T z_1|^2 - 2|A_1^T z_1||A_2^T B z_2| .$$

But, when $|z_2| < \varepsilon$, $2|A_2^T B z_2| \leq 2\varepsilon|A_2^T B| \leq \{a(1 - \varepsilon^2)\}^{1/2}/2$ from above, and $|A_1^T z_1|^2 \geq a|z_1|^2$ (Lemmas 4 and 2) $\geq a(1 - \varepsilon^2)$, so $|G^T z|^2 \geq a(1 - \varepsilon^2)/2$.

LEMMA 6.   *Let $B^* = B + HKH^T$, and suppose that $B^*$ and $H^T B^{*-1} H$ are non-singular. Then $Q = B^{*-1} H (H^T B^{*-1} H)^{-1}$ and $V = (I - QH^T) B^{*-1}$ are both independent of $K$.*

PROOF.   Suppose $K$ varies with a parameter $x$, but $B$ and $H$ do not. Then $\partial Q/\partial x = -B^{*-1} H (\partial K/\partial x) H^T B^{*-1} H (H^T B^{*-1} H)^{-1} + B^{*-1} H (H^T \cdot B^{*-1} H)^{-1} H^T B^{*-1} H (\partial K/\partial x) H^T B^{*-1} H (H^T B^{*-1} H) = 0$ and $\partial V/\partial x = (I - QH^T) \cdot (-B^{*-1} H(\partial K/\partial x) H^T B^{*-1}) = -V H (\partial K/\partial x) H^T B^{*-1} = 0$ since $VH = 0$.

UNIVERSITY OF SURREY

REFERENCES

[1] Aitchison, J. and Silvey, S. D. (1958). Maximum likelihood estimation of parameters subject to restraints, *Ann. Math. Statist.*, 29, 813-828.

[2] Brown, B. M. (1971). Martingale central limit theorems, *Ann. Math. Statist.*, 42, 59-66.

[3] Cramér, H. (1946). *Mathematical Methods of Statistcs*, Princeton University Press.

[4] Crowder, M. J. (1975). Maximum likelihood estimation for dependent observations, *J. R. Statist. Soc.*, B, 38, 45-53.

[5] Hudson, D. J. (1969). Least-squares fitting of a polynomial constrained to be either non-negative, non-decreasing or convex, *J. R. Statist. Soc.*, B, 31, 113-118.

[6] Moran, P. A. P. (1971). Maximum likelihood estimation in non-standard conditions, *Proc. Camb. Phil. Soc.*, 70, 441-450.

[7] Scott, D. J. (1973). Central limit theorems for martingales and for processes with stationary increments using a Skorokhod representation approach, *Adv. Appl. Prob.*, 5, 119-137.

[8] Silvey, S. D. (1959). The Lagrangian multiplier test, *Ann. Math. Statist.*, 30, 389-407.

[9] Wald, A. (1949). Note on the consistency of the maximum likelihood estimate, *Ann. Math. Statist.*, 20, 595-601.