# INFERENCE IN THE INDETERMINATE PARAMETERS PROBLEM

Marco Barnabani

## 1. INTRODUCTION

Let $f(x,\theta)$ $\theta \in \Theta \subseteq R^k$ be a density function continuous on $\Theta$, defining the distribution corresponding to the parameter $\theta$ in a neighbourhood of a particular point, $\theta_0$, say in $U_\delta = \{\theta; \|\theta - \theta_0\| \leq \delta\}$ where $\|.\|$ is the square norm and $\theta_0$ is the true, though unknown, parameter value; $(x_1, x_2, ..., x_n, ....)$ is a given sequence of independent observations; $\log L(\theta) = \sum_{i=1}^n \log f(x_i, \theta)$ is the log-likelihood function defined on $\Theta$ and $B(\theta_0)$ is the (Fisher) information matrix in an observation.

Assume $\theta$ to be partitioned into two subvectors, $\theta' = [\psi' \ \gamma']$ with $\psi$ of order $m$ and $\gamma$ of order $q = k - m$. We face an indeterminacy problem when there exist two disjoint and exhaustive subsets of $\psi$, $\{\psi_j, j \in J\}$, $\{\psi_t, t \in T\}$ say, such that the null hypothesis $H_0 : \psi_j = \psi_{j0}$ for all $j \in J$ makes the likelihood independent of $\gamma$ (see Cheng and Traylor (1995) for a definition of indeterminacy based on a general transformation $\phi = \phi(\theta)$). A common case is when $\psi_j = \psi_{j0}$ makes $\gamma$ indeterminate. In applications the complementary subset $\{\psi_t, t \in T\}$ can be the null set. In this case $\{\psi_j, j \in J\}$ coincides with $\psi$ and the null hypothesis involves the whole subvector $\psi$.

A well-known example of indeterminacy is the simple mixture model with probability density function

$$f(x,\theta) = (2\pi)^{-1/2}[(1-\psi)\exp(-x^2/2) + \psi \exp(-(x-\gamma)^2/2)], \qquad 0 \leq \psi \leq 1.$$

Setting either $\psi = 0$ or $\gamma = 0$ eliminates the other from the expression for $f(x,\theta)$. Examples of indeterminacy abound in non linear regression models.

Let $\eta = \sum_{i=1}^r \psi_i x_i + \psi \exp(\sum_{i=1}^q \gamma_i z_i)$. Then, $\psi = 0$ eliminates $\gamma_i, i = 1, ..., q$ from the model.

Consequences of indeterminacy are

a) The score is a vector with a first component of order $m$ (the first derivative of the log-likelihood with respect to $\psi$) which depends on the parameter $\gamma$ and can depend on $\{\psi_t, t \in T\}$; a second component of order $q$ (the first derivative of the log-likelihood with respect to $\gamma$) which is zero.

b) The expected information matrix in an observation is singular and block diagonal with all submatrices zeroes and the north-west block matrix of order $m \times m$ which depends on the parameter $\gamma$ and on $\{\psi_t, t \in T\}$. That is, when $\psi_j = \psi_{j0}$ the expected information matrix assumes the following form

$$B(\psi_{j0}, \psi_t, \gamma) = \begin{bmatrix} B_{\psi\psi}(\psi_{j0}, \psi_t, \gamma) & 0 \\ {\scriptstyle m \times m} & {\scriptstyle m \times q} \\ 0 & 0 \\ {\scriptstyle q \times m} & {\scriptstyle q \times q} \end{bmatrix}$$

that shows both a singularity and a local orthogonality between $\psi$ and $\gamma$.

c) Let $\psi_0 = [\psi_{j0} \; \psi_{t0}]$ be the "true" parameter of $\psi$. Then, in the indeterminate parameters problem the submatrix $B_{\psi\psi}(\psi_0, \gamma)$ is non-singular for any $\gamma$.

d) The Hessian of the log-likelihood is not singular and a solution to the log-likelihood equation can be computed.

Given the above features of the likelihood, the score and the information matrix, we cannot use $\tilde{\theta}_n$ the joint estimation by maximum likelihood of both $\psi$ and $\gamma$, for inferential purposes. As known, in this case the standard results, such as the asymptotic chi-squared distribution of the Wald test statistic or of the likelihood ratio statistic, are generally not true and the correct results depend very much on the precise problem being investigated.

In the "regular" case if the value of $\gamma$ were known and the hypothesis is $\psi = \psi_0$ then, under the usual regularity conditions, the asymptotic distribution of the maximum likelihood estimator $\tilde{\psi}_n$ of $\psi$ is well known to be normal with mean vector $\psi_0$ and variance-covariance matrix $B_{\psi\psi}^{-1}(\psi_0, \gamma)$. Moreover, because $\gamma$ is assumed known, the Wald test $W = n(\tilde{\psi}_n - \psi_0)' B_{\psi\psi}(\psi_0, \gamma)(\tilde{\psi}_n - \psi_0)$ is distributed asymptotically as a central $\chi^2(m)$. Durbin (1970) called naive a Wald type test based on an estimator of $\psi$ (with $\gamma$ estimated consistently) that has the same asymptotic distribution as the maximum likelihood estimator $\tilde{\psi}_n$ assuming the true value of $\gamma$ known.

In his paper Durbin argues that the maximum likelihood estimator of $\psi$ assuming $\gamma$ equal to the solution of the (constrained) equation

$$\frac{\partial}{\partial \gamma} \log L(\psi_0, \gamma) = 0 \tag{1}$$

produces a naive test if the maximum likelihood estimators of $\psi$ and $\gamma$ in the full model are asymptotically uncorrelated.

We note that this condition holds for the indeterminate parameters problem (consequence (b) above), nevertheless, in this case, because of the disappearance of the parameter $\gamma$ from the likelihood function it is not possible to solve equation (1), to calculate the maximum likelihood estimator of $\psi$ and to derive its asymptotic properties. Therefore, when we face an indeterminacy problem, the Durbin's approach based on a constrained estimator of $\gamma$ is unfeasible.

The aim of this paper is to look for an estimator of the parameters of interest, $\psi$ (treating the parameter $\gamma$ somehow), so that a Wald-type test statistic can be used for testing $H_0 : \psi_j = \psi_{j0}$ for all $j \in J$. We'll require that this estimator has the same asymptotic distribution as that of $\tilde{\psi}_n$ in the case of nonsingularity of the information matrix. We continue to call naive such an estimator and naive test the corresponding Wald type test.

As said above, in the indeterminate parameters problem the information matrix is block diagonal and positive semi-definite. Then, to tackle the indeterminacy we must face up to these characteristics of $B(\theta_0)$. The paper is organized as follows. In Section 2 we briefly review some existing results on the singularity of the information matrix and in a work of Silvey (1959) we found a possible approach to tackle the indeterminacy problem. After we had briefly recalled (Section 3) the properties of the maximum likelihood estimator in the regular case, in Section 4.1 we deal with the genesis of a naive maximum likelihood estimator and in Section 4.2 we detect its properties and its applicability to the indeterminate parameters problem. Finally, in Section 5 we show a Monte Carlo simulation applied to two nonlinear statistical models detecting the performance of the proposed estimator in small samples.

## 2. PREVIOUS WORKS ON THE SINGULARITY OF THE INFORMATION MATRIX

Perhaps, the author who first tackled the problem of the singularity of $B(\theta_0)$ was Silvey (1959). He recognized that the singularity of the information matrix is the main symptom of the lack of identifiability (a necessary but not sufficient condition for the non-identification problem) and he proposed a solution in this field. Silvey's approach is based on a modification of the information matrix adding an appropriate matrix to $B(\theta)$ obtained by imposing some restrictions on the parameters of the model so that the restricted parameters are identified and the modified matrix is positive definite. Poskitt and Tremayne (1981) have pointed out that the inverse of this matrix is in fact a generalized inverse of the information matrix. El-Helbawy and Hassan (1994) further generalized Silvey's results. Silvey's approach is very simple and elegant but its applicability is limited to the non-identification problem. In particular it is not applicable when the sin-

gularity of $B(\theta)$ is caused by one or more nuisance parameters vanishing under the null hypothesis.

In finite mixture models such as in the typical example of the previous section, a likelihood-based approach does not produce a satisfactory solution (Hartigan, 1985) and some authors suggest following other procedures (for example Wald's approach to testing) in alleviating problems caused by the singularity of $B(\theta)$ (Kay, 1995, discussion of the paper by Cheng and Traylor). Examples concerning hypothesis tests involving parameters not identifiable under the null hypothesis abound in nonlinear regression models (Seber and Wild, 1989) and several *ad hoc* solutions have been proposed. Cheng and Traylor (1995) introduced the "intermediate model" between the models where parameters are missing and where they are present. This approach is based on suitable reparameterizations and its success depends on how well the reparameterization positions the "intermediate model" between the two extremes. This procedure seems to be very difficult to apply when the number of vanishing parameters is relatively high.

Davies (1977, 1987) proposed an interesting approach to the problem of hypothesis testing when a nuisance parameter is present only under alternative. Given a suitable test statistic he suggested treating it as a function of the underidentified nuisance parameters and basing the test upon the maximum of this function. The asymptotic distribution of this maximum is not standard but Davies provided an upper bound for the significance level of his procedure. Though elegant, "Davies' method is quite elaborate to implement in practice and difficult to generalize" (Cheng and Traylor, 1995) particularly when several nuisance parameters vanish under the null hypothesis. Moreover, "there is no analytically tractable solution to Davies's maximization problem" (Godfrey, 1990, p. 90).

Segmented regression is another subject where singularity of the information matrix can occur. For example in the two phases linear regression, the null hypothesis of one single segment creates difficulties with the usual asymptotic chisquare theory for the likelihood ratio test for one phase against two. In this subject several *ad hoc* solutions have been proposed (Smith, 1989).

Rotnitzky *et al.*, (2000) provided an asymptotic distribution of the maximum likelihood estimator and of the likelihood ratio test statistic when the model is identified and the information matrix has rank one less than full. This approach is based on a suitable reparameterization of the model and was motivated by models with selection-bias but it seems quite complex and difficult to apply to models where the rank of $B(\theta)$ is arbitrary.

In the above brief survey, the solutions proposed are generally based on suitable reparameterizations of the model so that to remove the causes of singularity. As a consequence of this approach the solutions proposed are often difficult to generalize because they usually depend on the particular issue being investigated.

From a thorough analysis of the above works the mathematical aspect of singularity emerges. It affects the asymptotic approximating quadratic model of the log-likelihood function which can have a whole linear sub-space of maxima in a neighborhood of the "true" parameter. In that case we can say that we are faced

by (asymptotic) unstable parameters Ross (1990), in the sense that in a neighbourhood of the true parameter the asymptotic log-likelihood function cannot be approximated by a concave quadratic form using the second-order term in the Taylor series expansion about $\theta_0$. Therefore, a possible solution to the problem of singularity could be passed through a modification of the curvature of this quadratic model.

In our opinion, the author who first tackled the problem of singularity following this approach was Silvey (1959). As said, he proposed to replace $B(\theta)$ by $B(\theta_0) + A$ where $A$ is an appropriate matrix obtained introducing some restrictions on $\theta$. As we pointed out, this approach is very simple and gives an elegant solution to the problem, but it is of limited applicability. Nevertheless, we think that Silvey's idea could be generalized replacing $B(\theta_0)$ by $B(\theta_0) + \lambda I$ with $\lambda > 0$ (strictly positive) and $I$ the identity matrix of appropriate dimension. We show that this modified matrix is compatible with the definition of a penalized log-likelihood function, and inferences on the non-vanishing parameters can be based on the maximizing point of this function. Under usual regularity conditions, the estimator so obtained is consistent and asymptotically normally distributed with a variance-covariance matrix approximated by the Moore-Penrose pseudoinverse of the information matrix, which always exists and is unique (Rao and Mitra, 1971). In an indeterminacy problem this result allows us to construct a naive test useful for inferential purposes.

## 3. THE REGULAR CASE

We assume the following conditions which are straightforward generalizations of Cramer's conditions (Aitchison and Silvey, 1958).

$\mathfrak{F}1$) $\Theta$ is a compact subset of the Euclidian k-space and $\theta_0$ is an interior point.

$\mathfrak{F}2$) For every $\theta \in \Theta$, $Q(\theta) = E_0[\log f(x, \theta)]$ that is, the expected value of $\log f(x, \theta)$ taken with respect to a density function characterized by the parameter vector $\theta_0$, exists.

$\mathfrak{F}3$) For every $\theta \in U_\delta$ (and for almost all $x \in R$) first, second and third order derivatives with respect to $\theta$ of $\log f(x, \theta)$ exist and are bounded by functions independent of $\theta$ whose expected values are finite.

$\mathfrak{F}4$) The information matrix in an observation is positive definite (local identifiability condition).

In the regular case the classical proof of the consistency of a solution of the likelihood equations, $D \log L(\theta) = 0$, is based on the analysis in $U_\delta$ of the behaviour of the maximizing point of the quadratic model obtained from a Taylor series expansion of , $n^{-1} \log L(\theta)$ about $\theta_0$

$$\frac{1}{n}\log L(\theta) = \frac{1}{n}\log L(\theta_0) + \frac{1}{n}D'\log L(\theta_0)h + \frac{1}{2n}h'D^2\log L(\theta_0)h + R \qquad (2)$$

where $h = \theta - \theta_0$, $D = [\partial/\partial\,\theta_i]$ $i = 1,...,k$ is the column vector of a differential operator; $D^2 = [\partial^2/\partial\,\theta_i\partial\,\theta_j]$ $i,j = 1,...,k$ is the matrix of second derivatives. By imposing the first order necessary conditions for a maximum to the log-likelihood function or by expanding the likelihood equations about $\theta_0$ after re-scaling by $n^{-1}$, we have:

$$\frac{1}{n}D\log L(\theta_0) + \frac{1}{n}D^2\log L(\theta_0)h + \frac{1}{2}V(x,\theta) = 0 \qquad (3)$$

where $V(x,\theta)$ is a vector whose $i$-th component may be expressed in the form $n^{-1}(\theta - \theta_0)'\Delta_i(\theta^*)(\theta - \theta_0)$, $\Delta_i(\theta^*)$ being a matrix whose $(j,m)$ element is $\sum_{t=1}^{n}(\partial^3/\partial\theta_i\partial\theta_j\partial\theta_m)\log f(x_t;\theta^*)$, $j,m = 1,...,k$ bounded in $U_\delta$ and $\theta^*$ a point such that $\|\theta^* - \theta_0\| < \|\theta - \theta_0\|$. Conditions $\mathfrak{F}1-\mathfrak{F}3$ ensure that $n^{-1}D\log L(\theta_0)$ converges in probability to $0 \in R^k$; $n^{-1}D^2\log L(\theta_0)$ converges in probability to $-B(\theta_0)$, and the elements of $n^{-1}\Delta_i(\theta^*)$ are bounded for $\theta \in U_\delta$. We have the following Lemma

*Lemma 1* (Aitchison and Silvey, 1958). Subject to the conditions $\mathfrak{F}1-\mathfrak{F}4$, for large enough $n$, and $\delta$ sufficiently small, the equation (3) has a (unique) solution $\tilde{h} = \tilde{\theta}_n - \theta_0$ such that $\tilde{h}'\tilde{h} \le \delta^2$ if and only if $\tilde{h}$ satisfies a certain equation of the form

$$-B(\theta_0)h + m(x,\theta)\delta^2 = 0 \qquad (4)$$

where $m(x,\theta)$ is a continuous function on $U_\delta$ and $\|m(x,\theta)\|$ is bounded in $U_\delta$ by a positive number $\tau$, say.

Because of condition $\mathfrak{F}4$ the latent roots $\mu_1 \le \mu_2 \le ... \le \mu_k$ of the information matrix are all positive. Using an equivalent of Brower's fixed point theorem as in Aitchison and Silvey (1958), $\delta < \mu_1/\tau$ is a sufficient condition for equation (4) to have a unique solution $\tilde{h}$ such that $\tilde{h}'\tilde{h} \le \delta^2$.

Taking the probability limit of both sides of (2) and using the above assumptions, we have

$$p\lim\frac{1}{n}\log L(\theta) = Q(\theta) = Q(\theta_0) - \frac{1}{2}h'B(\theta_0)h + h'm(x,\theta)\delta^2$$

Equation (4) may be seen as the first order necessary conditions for the unconstrained maximum of the quadratic model $Q(\theta)$. Then, the crucial point of the consistency of a solution to the likelihood equations is that for $\delta$ sufficiently small (actually for $\delta < \mu_1/\tau$), $Q(\theta)$ has a unique maximizing point in $U_\delta$.

As to the asymptotic distribution of the maximum likelihood estimator we have

$$\left(\frac{1}{n}D^2 \log L(\theta_0) + \tilde{h}' \mathrm{R}^*\right) n^{1/2} \tilde{h} = -\frac{1}{\sqrt{n}} D \log L(\theta_0) \tag{5}$$

where $\mathrm{R}^*$ is a vector whose $i-th$ component may be expressed as $(2n)^{-1} \Delta_i(\theta^*)$ and $\theta^*$ a point such that $\left\| \theta^* - \theta_0 \right\| < \left\| \tilde{\theta}_n - \theta_0 \right\|$. In the regular case, $p\lim n^{-1} D^2 \log L(\theta_0) = -B(\theta_0)$, $n^{-1/2} D \log L(\theta_0)$ has the limiting normal distribution $N(0, B(\theta_0))$ by the central limit theorem and because of the consistency of the estimator, $\tilde{h}' \mathrm{R}^* = o_P(1)$. On compiling these results, we see that $n^{1/2} \tilde{h}$ tends in distribution to a random vector $B^{-1}(\theta_0)\eta$ where $\eta \sim N(0, B(\theta_0))$; and we conclude that $n^{1/2} \tilde{h}$ has the limiting normal distribution $N(0, B^{-1}(\theta_0))$.

We point out that in the regular case a solution of the likelihood equation has the same limiting distribution as the (unfeasible) linearized estimator

$$S_n = \theta_0 - \left(\frac{1}{n}D^2 \log L(\theta_0)\right)^{-1} \frac{1}{n} D \log L(\theta_0)$$

obtained by maximizing the quadratic model to $n^{-1} \log L(\theta_0)$ given by (2) with approximation error of order $o(\left\| \theta - \theta_0 \right\|)$ in $U_\delta$. As known, $S_n$ is the basis of several numerical procedures used to obtain a maximum likelihood estimator.

## 4. THE INDETERMINATE PARAMETERS PROBLEM

### 4.1. *An unfeasible estimator*

Suppose that the conditions $\mathfrak{F}1 - \mathfrak{F}3$ are satisfied. Then, when $\gamma$ vanishes the information matrix is singular and the asymptotic approximation, $Q(\theta)$, will not have a unique maximizing point in a neighbourhood of $\theta_0$ but a whole (linear) sub-space of maxima. The demands that $Q(\theta)$ should have a maximum in $U_\delta$ and that $B(\theta_0)$ should be positive definite are, clearly, related. In fact, if $B(\theta_0)$ is singular, nothing guarantees the existence of a unique solution $\tilde{h}$ which maximizes $Q(\theta)$ and such that $\tilde{h}' \tilde{h} \leq \delta^2$ for any $\delta$ sufficiently small.

Following Silvey (1959) a way to tackle the problem of the singularity of $B(\theta_0)$ is to modify the information matrix constraining the function $Q(\theta)$ in $U_\delta$. Using the Lagrange multiplier method, we could proceed to maximize $Q(\theta)$ subject to the constraint $\|\theta - \theta_0\| \leq \delta$. As known, a solution to this constrained problem, $\hat{h} = \hat{\theta}_\lambda^{(n)} - \theta_0$ say, must satisfy the following equation (Dennis and Schnabel, 1983, p. 131)

$$
\begin{aligned}
&-(B(\theta_0) + \lambda I)\, h + m(x,\theta)\, \delta^2 \equiv -A_\lambda(\theta_0)\, h + m(x,\theta)\, \delta^2 = 0 \\
&\Rightarrow \hat{h} = A_\lambda^{-1}(\theta_0)\, m(x,\theta)\, \delta^2
\end{aligned}
\tag{6}
$$

where $I$ is the identity matrix of an appropriate dimension and $\lambda > 0$ (strictly positive) a scalar determined so that $\left\| \hat{\theta}_\lambda^{(n)} - \theta_0 \right\| = \delta$. That is, the constrained maximum of $Q(\theta)$ occurs on the boundary of the region $\|\theta - \theta_0\| \leq \delta$ fixing appropriately $\lambda$.

If we compare (6) with that obtained in the regular case given by (4) we can observe that the information matrix is now modified by adding a scalar diagonal matrix giving rise to a "new" matrix $A_\lambda(\theta_0)$ which is positive definite. What about the meaning and the interpretation of $\hat{\theta}_\lambda^{(n)}$ fixing $\lambda$ arbitrarily?. The following results are well known in numerical analysis (Goldfeld *et al.*, 1966).

a) Given $\lambda$, $\hat{\theta}_\lambda^{(n)}$ is the maximizing point of the function $P(\theta) = Q(\theta) - (\lambda/2)\|\theta - \theta_0\|^2$ obtained by penalizing the asymptotic approximation $Q(\theta)$ with a quadratic penalty term. Because $A_\lambda(\theta_0)$ is positive definite, $P(\theta)$ has a global maximum at $\hat{\theta}_\lambda^{(n)}$.

b) From a) $P(\hat{\theta}_\lambda^{(n)}) = Q(\hat{\theta}_\lambda^{(n)}) - (\lambda/2)\left\| \hat{\theta}_\lambda^{(n)} - \theta_0 \right\|^2 \geq Q(\theta) - (\lambda/2)\|\theta - \theta_0\|^2$ and $Q(\hat{\theta}_\lambda^{(n)}) \geq Q(\theta)$ for all $\theta$ such that $\|\theta - \theta_0\| = \left\| \hat{\theta}_\lambda^{(n)} - \theta_0 \right\| = \delta_\lambda$. That is, if we define a region consisting of all $\theta$ such that $\|\theta - \theta_0\| \leq \delta_\lambda$ then the maximum of $Q(\theta)$ occurs on the boundary of this region.

c) $0 \leq \delta_\lambda = \left\| (B(\theta_0) + \lambda I)^{-1} m(x,\theta) \right\| \delta^2 \leq \delta^2 \tau/\lambda$. Then, if $\delta \leq \lambda/\tau$, $\delta_\lambda \leq \delta$ that is, $\hat{\theta}_\lambda^{(n)}$ is in $U_\delta$.

The above remarks suggest a way to use $A_\lambda(\theta_0)$. Define the following (penalized) log-likelihood function

$$
P_n(\theta) = \log L(\theta) - \frac{\lambda}{2}\|\theta - \theta_0\|
\tag{7}
$$

and let $\hat{\theta}_\lambda^{(n)}$ be a solution of the (penalized) likelihood equations

$$\frac{1}{n} D \log L(\theta_0) + \left(\frac{1}{n} D^2 \log L(\theta_0) - \lambda I\right) h + \frac{1}{2} V(x, \theta) = 0 \tag{8}$$

We can state the following Lemma equivalent to Lemma 1 for the "regular" case.

*Lemma 2* Given $\lambda > 0$, under the conditions $\mathfrak{F}1 - \mathfrak{F}3$, for large enough $n$, and $\delta$ sufficiently small, the equation (8) has a (unique) solution, $\hat{h} = \hat{\theta}_\lambda^{(n)} - \theta_0$ such that $\hat{h}' \hat{h} \le \delta^2$ if and only if $\hat{h}$ satisfies a certain equation of the form

$$-(B(\theta_0) + \lambda I)h + m(x, \theta)\delta^2 = 0 \tag{9}$$

where $m(x, \theta)$ is a continuous function on $U_\delta$ and $\|m(x, \theta)\|$ is bounded in $U_\delta$ by a positive number $\tau$, say.

Because $B(\theta_0) + \lambda I$ is positive definite, the system (9) has a unique solution in a neighborhood of $\theta_0$ if $\delta$ is sufficiently small. Indeed, it is sufficient $\delta < \lambda / \tau$.

The asymptotic distribution of $\hat{\theta}_\lambda^{(n)}$ is immediate following the the same line of reasoning as in the "regular" case. We have

$$\left(\frac{1}{n} D^2 \log L(\theta_0) - \lambda I + \hat{h}' R^o\right) n^{1/2} \hat{h} = -n^{-1/2} D \log L(\theta_0)$$

where $R^o$ is a vector calculated at some point in $U_\delta$ and bounded in $U_\delta$.

Under above conditions, $p \lim [n^{-1} D^2 \log L(\theta_0)] = -B(\theta_0)$, $n^{-1/2} D \log L(\theta_0)$ has the limiting normal distribution $N(0, B(\theta_0))$ and, because of the consistency of the estimator, $\hat{h}' R^o = o_P(1)$. Then, $n^{1/2} \hat{h}$ tends in distribution to a random vector $(B(\theta_0) + \lambda I)^{-1} \eta$ where $\eta \sim N(0, B(\theta_0))$ and we conclude that $n^{1/2} \hat{h}$ has the limiting normal distribution $N(0, A_\lambda^{-1}(\theta_0) B(\theta_0) A_\lambda^{-1}(\theta_0))$.

### 4.2. *The naive maximum likelihood estimator*

The definition and the use of the (penalized) log-likelihood function, $P_n(\theta)$, given in the previous section, leads to the following observations.

i) $P_n(\theta)$ can be interpreted as a penalty function where the penalty term is expressed in quadratic form. In a "non-regular" theory, the approach based on a modified log-likelihood function is certainly not new. The logarithmic barrier

function has been used in recent times to overcome the boundary problem and the non-identifiability in mixture models (Chen *et al.*, (2001)).

*ii)* $P_n(\theta)$ can be motivated by a Bayesian procedure or by incorporating a stochastic constraint. In the Bayesian motivation, let $\theta$ have the prior density proportional to $\exp[(-\lambda/2)\|\theta - \theta_0\|^2]$ so that $\exp[P_n(\theta)]$ is proportional to the posterior density. Alternatively, we can think of equation (7) as a constrained log-likelihood where the constraint is of the form $\theta = \theta_0 + v, \quad v \sim (0, \lambda^{-1}I)$ where $I$ is the identity matrix of an appropriate dimension. The stochastic constraint is introduced into the log-likelihood function through the penalty function approach.

*iii)* The maximization of $P_n(\theta)$ is not a feasible procedure because, given $\lambda$, the procedure depends on the unknown "true" parameter $\theta_0$ and the problem on fixing $\lambda$ arises.

The observation *iii)* is closely bound up with the aim of our paper and it can be solved if we can answer to the following question. Given the (unfeasible) estimator $\hat{\theta}_\lambda^{(n)}$ how can we construct a naive test?. In other words, when $\hat{\psi}_\lambda^{(n)}$, the first component of $\hat{\theta}_\lambda^{(n)}$, could have (at least approximately) the same asymptotic distribution as the maximum likelihood estimator $\hat{\psi}_\lambda^{(n)}$, given $\lambda$?. To solve this problem we propose to take $\lambda$ very small, formally, $\lambda \to 0$.

Then, what about the asymptotic properties of a solution of the following (penalized) likelihood equations?

$$\lim_{\lambda \to 0}\left[\frac{1}{n}D\log L(\theta_0) + \left(\frac{1}{n}D^2\log L(\theta_0) - \lambda I\right)h + \frac{1}{2}V(x,\theta) = 0\right] \qquad (10)$$

We can state the following Theorem

*Theorem* Under the conditions $\mathfrak{F}1 - \mathfrak{F}3$, for large enough $n$, and $\delta$ sufficiently small, the equation (10) has a (unique) solution $\hat{h}_0 = \hat{\theta}_{\lambda 0}^{(n)} - \theta_0$ in a neighborhood of the true parameter $\theta_0$. Moreover,

$$n^{1/2}(\hat{\theta}_{\lambda 0}^{(n)} - \theta_0) \sim N(0, B^+(\theta_0)) \qquad (11)$$

where $B^+(\theta_0)$ is the Moore-Penrose pseudoinverse of $B(\theta_0)$.

*Proof.* Consistency is immediate invoking Lemma 2. The asymptotic distribution of the estimator emerges following the same line of reasoning as that at the end of the previous Section. Therefore, under the conditions $\mathfrak{F}1 - \mathfrak{F}3$ we can state that $n^{1/2}\hat{h}_0$ tends in distribution to $\lim_{\lambda \to 0}(B(\theta_0) + \lambda I)^{-1}\eta$ where

$\eta \sim N(0, B(\theta_0))$. Then, $n^{1/2}\hat{b}_0$ has the limiting normal distribution $N(0, B^+(\theta_0))$ where $B^+(\theta_0) = \lim_{\lambda \to 0}(A_\lambda(\theta_0)^{-1}B(\theta_0)A_\lambda(\theta_0)^{-1})$. It is immediate to show that $B^+(\theta_0)$ is the Moore-Penrose pseudoinverse of $B(\theta_0)$ which always exists and is unique (Albert, 1972).

For the indeterminate parameters problem given the particular form assumed by the information matrix, the pseudoinverse of $B(\psi_0, \gamma_0)$ is given by

$$B^+(\psi_0, \gamma_0) = \begin{bmatrix} B_{\psi\psi}^{-1}(\psi_0, \gamma_0) & 0 \\ 0 & 0 \end{bmatrix}$$

then, $W_{\lambda 0} = n(\hat{\psi}_{\lambda 0}^{(n)} - \psi_0)B_{\psi\psi}(\psi_0, \gamma)(\hat{\psi}_{\lambda 0}^{(n)} - \psi_0)$ is distributed as a central $\chi^2(m)$.

Let $B_{\psi\psi}^{-1}(\psi_0, \gamma)$ be partitioned in four blocks, $B^{11}$, $B^{12}$, $B^{21}$, $B^{22}$ and call $\psi_j$ and $\psi_t$ respectively the first and the second (block) component of the vector $\psi$. Then, we can test a subset of parameters $H_0 : \psi_j = \psi_{j0}$ through the statistic $n(\hat{\psi}_j - \psi_{j0})(B^{11})^{-1}(\hat{\psi}_j - \psi_{j0})$ which is distributed as $\chi^2(rank(B^{11}))$. It is immediate to note that $\hat{\theta}_{\lambda 0}^{(n)} - \tilde{\theta}_n = o(\lambda^{-\varepsilon})$ and $W_{\lambda 0} - W = o(\lambda^{-\varepsilon})$ with $\varepsilon > 0$.

## 5. SOME EXAMPLES

Applications of the naive test, $W_{\lambda 0}$, are closely associated with the possibility of obtaining a solution of the naive maximum likelihood estimator through equation (10). With respect to this problem, we first note that for any $\lambda$, the estimator $\hat{\theta}_\lambda^{(n)}$ has the same limiting distribution as the (unfeasible) linearized estimator

$$T_n = \theta_0 - \left(\frac{1}{n}D^2 \log L(\theta_0) - \lambda I\right)^{-1} \frac{1}{n} D \log L(\theta_0) \tag{12}$$

in the sense that $n^{1/2}(\hat{\theta}_\lambda^{(n)} - \theta_0) = n^{1/2}(T_n - \theta_0) + o_p(1)$. We underline that in the indeterminacy problem $T_n$ plays the same role as the (unfeasible) linearized estimator $S_n$ given for the regular case. Then, we can use (12) to obtain a solution to equation (10) through an iterative algorithm fixing in advance a sequence of $\lambda$ converging to zero. More specifically in the subsequent examples we computed the estimate following these steps:

*i)* Fix a sequence $\{\lambda_i\}$, typically $\{1, 10^{-1}, 10^{-2}, ...\}$, choose a starting point, $\theta^{(s)}$ and set $i = 1$.

*ii*) Check the termination condition. When a sufficiently small value of $\lambda_i$ has been reached the algorithm terminates.

*iii*) Compute an analytical Hessian matrix, $J(\theta^{(s)})$, and the matrix $A_\lambda = J(\theta^{(s)}) + \lambda_i I$.

*iv*) Find iteratively a solution to (12), call it $\theta^{(F)}$.

*v*) Set $\theta^{(s)} = \theta^{(F)}$, set i=i+1 and return to (*ii*).

This algorithm works quite well in the examples discussed in this paper.

*Example 1* (Gallant, 1987): Let $Y_1, Y_2,..., Y_n$ be a sequence of independent normal random variables with (known) variance $\sigma^2$ and expectations given by

$$E(Y_i) = \psi_1 x_{i1} + \psi_2 x_{i2} + \psi_3 \exp\left(\sum_{i=1}^{q} \gamma_i z_{ij}\right)$$

The inputs correspond to a one way "treatment-control" design that uses experimental variables that affect the response exponentially. Suppose we want to test the hypothesis $H_0 : \psi_3 = 0$. Then, under $H_0$, $\log L(\psi, \gamma) \propto 2^{-1} \sum_i (y_i - v_i)^2$, $v_i = \psi_1 x_{i1} + \psi_2 x_{i2}$, is independent on the $q$ nuisance parameters $\gamma_j$, $j = 1,...,q$ but depends on two parameters, $\psi_1$ and $\psi_2$ to be estimated. The elements of the score vector are

$$\frac{\partial}{\partial \psi_j} \log L(\psi_3, \psi_2, \psi_1, \gamma)\bigg|_{\psi_3=0} = \sigma^{-2} \sum_i v_i x_{ij} \qquad j = 1, 2$$

$$\frac{\partial}{\partial \psi_3} \log L(\psi_3, \psi_2, \psi_1, \gamma)\bigg|_{\psi_3=0} = \sigma^{-2} \sum_i v_i a_i$$

$$\frac{\partial}{\partial \gamma_j} \log L(\psi_3, \psi_2, \psi_1, \gamma)\bigg|_{\psi_3=0} = 0 \qquad j=1,...,q$$

where $a_i = \exp\left(\sum_j \gamma_j z_{ij}\right)$.

The information matrix on $n$ observations is given by

$$B_n(\psi_3 = 0, \psi_1, \psi_2, \gamma) = \sigma^{-2} \begin{bmatrix} \sum_i x_{i1}^2 & \sum_i x_{i1} x_{i2} & \sum_i x_{i1} a_i & 0 \\ \sum_i x_{i1} x_{i2} & \sum_i x_{i2}^2 & \sum_i x_{i2} a_i & 0 \\ \sum_i x_{i1} a_i & \sum_i x_{i2} a_i & \sum_i \gamma_i^2 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

which shows both a singularity and a local orthogonality between $\gamma$ and $\psi$.

For simulation purposes we construct independent variables following Gallant (1987, p. 19). The first two coordinates consist of the replication of a fixed set of design points determined by the design structure

$$(x_{i1}, x_{i2}) = \begin{cases} (1,1) & \textit{if } i \textit{ is odd} \\ (0,1) & \textit{if } i \textit{ is even} \end{cases}$$

As to the $q$ variables $z_{ij}$ we limited these to $q = 2$ and generated $z_{ij}$, $j = 1, 2$ by random selections from the uniform distribution in the interval $[0,10]$. Results are based on 5000 replications of samples of different sizes with $\psi_1 = -0.05$, $\psi_2 = 1$, $\psi_3 = 0$ and $\sigma^2 = 0.001$. The model is very sensitive to the choice of the functional form of the distributions of $z_{ij}$, which must be positive everywhere on some known interval. Moreover, the initial point for the iterative process is crucial to be successful in the simulation. Therefore, a particular care with these aspects is required (Gallant, 1987). The naive test is given by $W_{\lambda 0} = b^{33}(\hat{\gamma})(\hat{\psi}_3^2) \sim \chi^2(1)$ where $b^{33}(\hat{\gamma})$ is the inverse of the third element of the principal diagonal of the pseudoinverse of $B_n(\psi_3 = 0, \psi_1, \psi_2, \gamma)$. Proportion of rejections of $H_0$ for different sample sizes are shown in Tab. 1.

TABLE 1

*Proportion of rejections of $H_0 : \psi_3 = 0$*

| Hypothesis | Sample size | | | |
|---|---|---|---|---|
| | $n = 30$ | $n = 50$ | $n = 70$ | $n = 80$ |
| $\psi_3 = 0$ | 0.57 | 0.28 | 0.105 | 0.057 |

The table shows that the proportion of rejections reaches the 0.05-significance level when the sample size is about 80.

*Example 2*: (Davies, 1987). Let $Y_1, ...., Y_n$ be a sequence of independent normal random variables with a unit variance and expectations given by

$$E(Y_i) = \begin{cases} a + bx_i & \text{if } x_i < \gamma \\ a + bx_i + c(x_i - \gamma) & \text{if } x_i \geq \gamma \end{cases}$$

where $x_i$ denotes the time and $\gamma$ the unknown time, at which the change in a slope occurs. We want to test the null hypothesis $H_0 : c = 0$ against the alternative that $c \neq 0$. We use simulation to investigate how rapidly the finite-sample performance of the test statistic based on the naive maximum likelihood estima-

tor approaches its asymptotic limit. For simulation purposes we construct an X matrix which has one in the first column, time such that $\sum_i x_i = 0$ in the second column, zero if $x_i < \gamma$ or $(x_i - \gamma)$ if $x_i \geq \gamma$ in the third. Then, we generated samples of different sizes starting from $n = 20$ using the following model $y_i = 1 + 3x_i + c(x_i - 1) + u_i$, $u_i \sim N(0,1)$, giving several values to the parameter c. Under $H_0$, one immediately notes that when the null hypothesis is true $\gamma$ vanishes from the model and the expected information matrix becomes singular

$$B_n(c = 0, a, b, \gamma) = \begin{bmatrix} n & \sum_i x_i & \sum_{2i}(x_i - \gamma) & 0 \\ \sum_i x_i & \sum_i x_i^2 & \sum_{2i} x_i(x_i - \gamma) & 0 \\ \sum_{2i}(x_i - \gamma) & \sum_{2i} x_i(x_i - \gamma) & \sum_{2i}(x_i - \gamma)^2 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$\sum_{2i}$ denotes the summation over $x_i \geq \gamma$.

In small samples, the application of the naive test to the two-phase model leads to define the test statistic, $W_{\lambda 0} = b^{33}(\hat{\gamma})\hat{c}^2 \sim \chi^2(1)$ where $b^{33}(\hat{\gamma})$ is the inverse of the third elementt of the principal diagonal of the pseudoinverse of $B_n(c = 0, a, b, \gamma)$.

Proportion of rejections of a null hypothesis for some value of c and different sample sizes are shown in Table 2. Results are based on 1000 simulation runs at a 5% level of confidence.

TABLE 2

*Proportion of rejections of $H_0 : c = 0$*

| Parameter | Sample size | | | |
|---|---|---|---|---|
| c | $n = 20$ | $n = 30$ | $n = 40$ | $n = 50$ |
| 0.0 | 0.134 | 0.124 | 0.053 | 0.037 |
| 0.1 | 0.142 | 0.144 | 0.154 | 0.145 |
| 0.2 | 0.265 | 0.33 | 0.387 | 0.773 |
| 0.3 | 0.42 | 0.64 | 0.942 | 1 |
| 0.4 | 0.651 | 0.85 | 1 | 1 |

The table shows that there are differences in the performance of the test when we move from samples of size 20 to 50. In particular, under the null hypothesis $H_0 : c = 0$ the proportion of rejections reaches the 0.05-significance level with a 95% confidence interval [36,64] when the sample size is 40. Moreover, when data are generated with $c = 0.1$ (we also tried with different values of $0 \leq c \leq 0.1$) the proportion of rejections is nearly constant at about 14-16 per cent. We have an increase of this percentage when $n$ is raised from 50 to 100 as shown in Table 3.

TABLE 3

*Proportion of rejections of $H_0 : c = 0.1$*

| Parameter | Sample size | | | | |
|---|---|---|---|---|---|
| c | $n = 60$ | $n = 70$ | $n = 80$ | $n = 90$ | $n = 100$ |
| 0.1 | 0.174 | 0.221 | 0.412 | 0.584 | 0.645 |

Because the two-phase model is taken from Davies (1987), a brief comment may be appropriate. Our remarks concern the approach used rather than the results obtained. The test based on the naive maximum likelihood estimator proposed in this paper may be considered "standard" because asymptotically the test statistic has a central chi-square distribution. Moreover, it is relatively simple to apply as it emerges from the above application. Davies' approach, though elegant, is quite elaborate to implement in practice and it is difficult to generalize when more than one parameter vanishes under the null hypothesis. In models more complex than those described in this paper, the asymptotic distribution of the test statistic constructed following Davies' method is unknown. Approximated distributions using simulation techniques are necessary and tabulation of critical values is impossible. Recent works that follow Davies' approach are Andrews and Ploberger (1994) and Hansen (1996).

6. CONCLUSIONS

In this paper we proposed a way to make inference in the indeterminate parameter problem. The approach is based on the definition of a modified (penalized) log-likelihood function letting a penalty parameter going to zero. The maximizing point of this function has attractive statistical properties. It is consistent and asymptotically normally distributed with variance-covariance matrix approximated by the Moore-Penrose pseudoinverse of the information matrix. These properties allow one to construct a Wald-type test statistic with a "standard" distribution both under the null and alternative hypotheses. This test is relatively simply to apply to the indeterminacy problem. The performance in small samples of the proposed test statistic is detected on two nonlinear regression models.

*Dipartimento di Statistica "G. Parenti"*      MARCO BARNABANI
*Università di Firenze*

REFERENCES

J. AITCHISON, S. D. SILVEY (1958), *Maximum-likelihood estimation of parameters subject to restraints.* "The Annals of Mathematical Statistics", 29, 813-828.

A. ALBERT (1972), *Regression and the Moore-Penrose pseudoinverse*, Academic Press, New York.

D. W. K ANDREWS, W. PLOBERGER (1994), *Optimal tests when a nuisance parameter is present only under the alternative.* "Econometrica", 62, 1383-1414.

H. CHEN, J. CHEN, J. D. KALBFLEISCH (2001), *A modified likelihood ratio test for homogeneity in finite mixture models.* "Journal of Royal Statistical Society", B, 63, 1, 19-29.

R. C. H. CHENG, L. TRAYLOR (1995), *Non-regular maximum likelihood problems.* "Journal of Royal Statistical Society", 57, 1, 3-44.

R. B. DAVIES (1977), *Hypothesis testing when a nuisance parameter is present only under alternative.* "Biometrika", 64, 247-254.

R. B. DAVIES (1987), *Hypothesis testing when a nuisance parameter is present only under the alternative.* "Biometrika", 74, 33-43.

J. E. DENNIS, R. E. SCHNABEL (1983), *Numerical methods for unconstrained optimization and nonlinear equations*, Prentice-Hall Inc., New Jersey.

J. DURBIN (1970), *Testing for serial correlation in least squares regression when some of the regressors are lagged dependent variables.* "Econometrica", 38, 410-421.

A. T. EL-HELBAWY, T. HASSAN (1994), *On the Wald, Lagrangian multiplier and the likelihood ratio tests when the information matrix is singular.* "Journal of The Italian Statistical Society", 1, 51-60.

R. FLETCHER (1980), *Practical methods of optimization.* Vol. 1, 2, Wiley, New York.

R. A. GALLANT (1987), *Nonlinear statistical models*, Wiley, New York.

L. G. GODFREY (1990), *Misspecification tests in econometrics.* Cambridge University Press, Cambridge.

M. GOLDFELD, R. E. QUANDT, H. F. TROTTER (1966), *Maximization by quadratic hill-climbing.* "Econometrica", 34, 541-551.

B. E. HANSEN (1996), *Inference when a nuisance parameter is not identified under the null hypothesis.* "Econometrica", 64, 2, 413-30.

J. A. HARTIGAN (1985), *A failure of likelihood asymptotics for normal mixtures.* Proc. Berkeley Symp. In Honor of J. Neyman and J. Kiefer, (eds L. LeCam and R.A.Olshen), vol. II, 807-810, New York, Wadsworth.

E. L. LEHMAN (1991), *Theory of point estimation.* Wadsworth, Inc, Belmont, Ca.

D. S. POSKITT, A. R. TREMAYNE (1981), *An approach to testing linear time series models.* "The Annals of Statistics", 9, 974-86.

C. R. RAO, S. K. MITRA (1971), *Generalized inverse of matrices and its applications.* New York: Wiley.

G. I. S. ROSS (1990), *Nonlinear estimation.* Springer, New York.

A. ROTNITZKY, D. R. COX, M. BOTTAI, J. ROBINS (2000), *Likelihood-based inference with singular information matrix.* "Bernoulli", 6(2), 243-284.

G. A. F. SEBER, C. J. WILD (1989), *Nonlinear regression.* Wiley, New York.

S. D. SILVEY (1959), *The Lagrangian multiplier test.* "The Annals of Mathematical Statistics", 30, 389-407.

R. L. SMITH (1989), *A survey of non-regular problems.* Proceedings of the International Statistical Institute, 47th Session, Paris, 353-372.

RIASSUNTO

*Inferenza nel problema dei parametri indeterminati*

Il problema dei parametri indeterminati nasce quando ci sono due insiemi di parametri, $\psi$ e $\gamma$, tali che l'ipotesi nulla $H_0 : \psi = \psi_0$ rende la verosimiglianza indipendente da $\gamma$. Una conseguenza di questa situazione è la singolarità della matrice di informazione. Per questo tipo di problema i risultati asintorici standard legati allo stimatore di massima verosimiglianza non sono generalmente validi. Nel lavoro si propone uno stimatore per il pa-

rametro di interesse, $\psi$, tale che, asintoticamente sia possibile sottoporre a verifica l'ipotesi $H_0$ utilizzando la statistica test di Wald. Un tale stimatore è ottenuto tramite la massimizzazione di una funzione di log-verosimiglianza penalizzata. Si mostra che una soluzione dell'equazione di verosimiglianza (penalizzata) è consistente, e distribuita asintoticamente in modo normale con matrice di varianze-covarianze approssimata dalla pseudoinversa di Moore-Penrose della matrice di informazione. Queste proprietà consentono di costruire una statistica test di Wald utile per scopi inferenziali.

<div align="center">SUMMARY</div>

*Inference in the indeterminate parameters problem*

We face an indeterminate parameters problem when there are two sets of parameters, $\psi$ and $\gamma$, say, such that the null hypothesis $H_0 : \psi = \psi_0$ makes the likelihood independent of $\gamma$. A consequence of indeterminacy is the singularity of the information matrix. For this problem the standard results, such as the asymptotic chi-squared distribution of the Wald test statistic, are generally false. In the paper we propose an estimator of the parameters of interest, $\psi$, so that a Wald-type test statistic can be used for testing $H_0$. Such an estimator is obtained through the maximization of a modified (penalized) log-likelihood function. We show that a solution to the (penalized) likelihood equation is consistent and asymptotically normally distributed with variance-covariance matrix approximated by the Moore-Penrose pseudoinverse of the information matrix. These properties allow one to construct a Wald-type test statistic useful for inferential purposes.