# Dipartimento di Statistica
## "Giuseppe Parenti"

# Maximum likelihood estimator and singularity of the information matrix

Marco Barnabani

## Università degli Studi di Firenze

# Maximum likelihood estimator and singularity of the information matrix

Marco Barnabani[1]
Dipartimento di Statistica "G. Parenti"
Università di Firenze
barnaban@ds.unifi.it

**Abstract:** When the model is identified but the information matrix is singular, the classic asymptotic properties of the maximum likelihood estimator are not clear and an inferential procedure based on it is not viable. In the paper a solution of a loglikelihood equation appropriately penalized is shown to be consistent and asymptotically normal distributed with variance-covariance matrix approximated by the Moore-Penrose pseudoinverse of the information matrix. These properties allow one to get a quadratic function based on a standard Chi-square distribution for hypothesis testing. A simulation applied to a simplified Engle's model is presented to support the theoretical results.

**Keywords:** Singular Information Matrix, Moore-Penrose pseudoinverse, Maximum likelihood estimator, penalized loglikelihood equation.

## 1  Introduction

Let $f(t; \theta)$ $\theta' = [\alpha' \ \beta'] \in \Theta \subseteq \mathbb{R}^k$, $\alpha \in \Theta_r \subseteq \mathbb{R}^r$, $\beta \in \Theta_{k-r} \subseteq \mathbb{R}^{k-r}$, $t \in \mathbb{R}$ be a density function continuous on $\Theta$ defining the distribution corresponding to the parameter $\theta$ in a neighbourhood of a true unknown parameter value, $\theta_0' = [\alpha_0' \ \beta_0']$. In this paper we tackle the problem of the asymptotic properties of maximum likelihood estimator when the information matrix is singular and the model is identified. We propose an estimator which allows one to do inference on the whole set of parameters, $\theta_0$ or on the parameter of interest $\alpha_0$, say.

Statistical literature on the singularity of the information matrix is large (see Rotnitzky *et al.* (2000) and the associated bibliografy). Models more relating to this paper concern hypothesis tests involving parameters not identifiable under the null hypothesis. Consider, for example, the following simplification of Engle's (1984) model

$$y/x \sim N(\alpha x^\beta, \sigma^2 = 1), \qquad x > 0, \qquad H_0 : \alpha = 0 \tag{1}$$

where $x$ is non-stochastic. In the model the parameter $\beta$ is estimable only when the null hypothesis is false. Under $H_0$ the hessian matrix is non-singular while the (expected)

---

[1]Viale Morgagni, 59 - 50134 Firenze.

information matrix in an observation, $B(\theta)$, given by

$$B(\alpha = 0\,,\beta) = \begin{bmatrix} x^{2\beta} & 0 \\ 0 & 0 \end{bmatrix}$$

is nonnegative definite. In this model an inference on the parameter of interest is possible only if we were able to handle $\beta$ somehow. In small samples the loglikelihood function of both $\alpha$ and $\beta$ can be maximized under the null and the alternative hypothesis, but because of the singularity of the information matrix the asymptotic properties of the joint estimator is not clear.

Models of this type abound in nonlinear regression where several *ad hoc* solutions have been suggested. For example, Cheng and Traylor (1995) proposed an "intermediate model" between the model where parameters are missing and where they are present. The solution proposed is based on suitable reparameterizations and the 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 parameter, 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. It has been observed (Cheng and Traylor, 1995) that, though elegant, "Davies' method is quite elaborate to implement in practice and difficult to generalize".

In general, most of the solutions proposed in the statistical literature are based on suitable reparameterizations of the particular model analyzed so that to remove the causes of singularity and to obtain (stable) asymptotic estimates. As a consequences of this approach the solutions are often difficult to generalize because they usually depend on the particular issue being investigated.

Perhaps the author who first suggested a solution to the singularity of the information matrix susceptible of generalization was Silvey (1959). Within the non-identification problem he proposed to replace the information matrix by $B(\theta_0) + F$ where $F$ is an appropriate matrix obtained imposing some restrictions on the parameters of the model so that the restricted parameters are identified and the "new" matrix is positive definite. More precisely, he suggested to set $F = H'_r H_r$ where $H_r$ is the jacobian of r *ad hoc* constraints imposed on $\theta$. In his work Silvey showed that statistical tests (Wald or Score) based on the inverse $B(\theta_0)^- = (B(\theta_0) + H'_r H_r)^{-1}$ are "standard" in the sense that under the null hypothesis they are asymptotically central chi-square distributions. Silvey's approach is very simple and elegant but, is not applicable when the singularity of $B(\theta_0)$ cannot be removed by constraining some parameters because the singularity of the matrix is caused, for example, by one or more nuisance parameters vanishing under the null hypothesis.

Several authors (Poskitt and Tremayne, 1981) have pointed out that $B(\theta_0)^-$ is a generalized inverse of $B(\theta_0)$. Then, a first step towards a generalization of the above ap-

proach could be based on the search of an estimator and consequently on the choice of an appropriate matrix $F$ such that a "standard" test based on a generalized inverse of the information matrix is possible. Unfortunately, this approach is unfeasible because of the non-uniqueness of $B(\theta_0)^-$ which causes some difficulties in finding a test invariant to the choice of this matrix. To overcome the invariancy problem we propose to replace $B(\theta_0)^-$ by the Moore-Penrose pseudoinverse $B(\theta_0)^+$ which always exists and is unique. Of course, in this case the main problem is to find an estimator computed somehow which is compatible (at least asymptotically) with this matrix. The search of this estimator is the goal of the paper.

The work is organized as follows. In Section 2 we review the asymptotic properties of maximum likelihood estimator in the regular case both from an analytical and geometrical point of view. In this section we repropose well known results which are preliminary for subsequent sections. In Section 3 we analyze the consequences of the singularity of the information matrix on the asymptotic properties of maximum likelihood estimator. We show that in a neighborhood of the true parameter still exists a solution to the likelihood equations but this solution is no more unique. Nothing we can say about the asymptotic distribution of the estimator. Section 4 is devoted to describe how to pick up one of the solutions which exist near the true parameter. We show that such an estimate can be chosen, following Silvey's idea, replacing $B(\theta_0)$ by $B(\theta_0) + \lambda I$, $\lambda > 0$ a scalar and $I$ an identity matrix of appropriate dimension, letting $\lambda \to 0$. In this way a solution near the maximum likelihood estimate can be found. We prove that this estimator is consistent and asymptotically normally distributed with variance-covariance matrix approximated by the Moore-Penrose pseudoinverse. These properties allow one to construct a Wald-type test statistic with a "standard" distribution both under the null and the alternative hypotheses. Finally, in Section 4 a numerical solution is given and a simulation of the properties of the estimator of the Engle's model is analyzed.

## 2 The Regular Case

The theory is said to be regular if, in a neighborhood of the true parameter $\theta_0$, "the log-likelihood function is closely approximated, in probability, by a concave quadratic function whose maximum point converges in some efficient sense to the true parameter value as the sample size increases. Conditions ensuring this are called regular conditions" (Cheng and Traylor, 1995).

Let $U_\delta = \{\theta; \|\theta - \theta_0\| \leq \delta\}$ be a neighborhood of $\theta_0$ where $\|.\|$ is the square norm; $x = (x_1, x_2, ..., x_n, ....)$ a given sequence of independent observations on $X$ and $logL(\theta) = \sum_{i=1}^{n} log f(x_i; \theta)$ the log-likelihood function defined on $\Theta$.

We assume the following 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 $\theta \in \Theta$, $z(\theta) = E_0[log f(t, \theta)]$ that is, the expected value of $log f(t; \theta)$ 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 $t \in \mathbb{R}$ ) first and second order derivatives with respect to $\theta$ of $log f(t; \theta)$ exist, are continuous functions of $\theta$ and are bounded by functions independent of $\theta$ whose expected values are finite. $\mathfrak{F}4-$ For every $\theta \in U_\delta$ and for $i, j, m = 1, \cdots k$, $|(\partial^3/\partial\theta_i\partial\theta_j\partial\theta_m)log f(t, \theta)| < G(t)$ where $E_0[G(t)] = M(t)$. $\mathfrak{F}5-$ For every $\theta \in U_\delta$ the information matrix in an observation, is positive definite with latent roots $\mu_1 \leq \mu_2 \leq \cdots \leq \mu_k$.

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

$$\frac{1}{n}log L(\theta) = \frac{1}{n}log L(\theta_0) + \frac{1}{n}D'log L(\theta_0)h + \frac{1}{2n}h'D^2 log L(\theta_0)h + \frac{1}{6}h'V(x;\theta^*) \quad (2)$$

where $h = \theta - \theta_0$; $D = [\partial/\partial\theta_i]$ $i = 1, \cdots, k$ is the column vector of a differential operator; $D^2 = [\partial^2/\partial\theta_i\partial\theta_j]$ $i, j = 1, \cdots, k$ is the matrix of second derivatives; $V(x;\theta^*)$ 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 $(\partial^3/\partial\theta_i\partial\theta_j\partial\theta_m)\sum_{t=1}^n log f(x_t, \theta^*)$ and $\theta^*$ a point such that $\| \theta^* - \theta_0 \| < \| \theta - \theta_0 \|$.

By imposing the first order necessary conditions for a maximum to the function $(2)$, or by expanding the likelihood equations about $\theta_0$ after rescaling by $n^{-1}$, we have:

$$\frac{1}{n}Dlog L(\theta_0) + \frac{1}{n}D^2 log L(\theta_0)h + \frac{1}{2}V(x;\theta^*) = 0 \quad (3)$$

Conditions $\mathfrak{F}1 - \mathfrak{F}4$ ensure that, for large enough $n$, $n^{-1}log L(\theta_0)$ is near $z(\theta_0)$, $\| n^{-1}Dlog L(\theta_0) \|$ is small, $-n^{-1}D^2 log L(\theta_0)$ is near a certain positive definite matrix $B(\theta_0)$ and $(\partial^3/\partial\theta_i\partial\theta_j\partial\theta_m)log f(x_t, \theta^*)$ is bounded in $U_\delta$. As $n$ goes to infinity, the $(j, m)$ element of $n^{-1}\Delta_i(\theta^*)$ converges in probability to its expected value that exists and does not depend on $\theta$. Therefore, $V(x;\theta^*)$ converges in probability to a function, $m(x)$, continuous on $U_\delta$ and such that $\| m(x) \|$ is bounded in $U_\delta$ by a positive number $\tau$, say. Then, for large $n$, $n^{-1}log L(\theta)$ can be approximated by the following quadratic model,

$$Q(\theta) \equiv z(\theta_0) - \frac{1}{2}h'B(\theta_0)h + h'm(x)\delta^2 \quad (4)$$

Moreover we have the following result

**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 likelihood equations have a solution $\widetilde{h} = \widetilde{\theta}_n - \theta_0 \in U_\delta$ if (and only if) it satisfies a certain equation of the form*

$$-B(\theta_0)h + m(x)\delta^2 = 0 \quad (5)$$

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

*Proof.* It is a straightforward generalization of Cramér's proof. See Aitchison and Silvey (1958) for details of the proof. $\square$

The fact that $B(\theta_0)$ is positive definite (condition $\mathfrak{F}5$) allows one to state that, if $\delta$ is less than a certain value, a solution to the system of equation (5) exists, is unique and belongs to $U_\delta$. Because $\delta$ can be chosen arbitrarily small, this is sufficient to show the statistical consistency of a solution to the likelihood equations. Indeed, we have the following Lemma

**Lemma 2.** *if $B(\theta_0)$ is positive definite and $\delta < \mu_1/\tau$ where $\mu_1 > 0$ is the minimum eigenvalue of $B(\theta_0)$, then $\widetilde{h}$ is the unique solution of equation (5) belonging to $U_\delta$*

*Proof.* If $B(\theta_0)$ is positive definite then the latent roots are all positive and the equation (5) has a unique solution given by $\widetilde{h} = B^{-1}(\theta_0)m(x)\delta^2$.
Because $\| B^{-1}(\theta_0)\, m(x) \| \leq \| m(x) \| \, max\, (\mu_i^{-1}), i = 1, \cdots, k$, then

$$\| \widetilde{h} \| = \| B^{-1}(\theta_0)\, m(x) \| \, \delta^2 \leq \| m(x) \| \, \mu_1^{-1}\, \delta^2 \leq \delta^2 \frac{\tau}{\mu_1}$$

If $\tau \mu_1^{-1} < \delta^{-1}$ which implies $\delta < \mu_1/\tau$ then, $\| \widetilde{h} \| < \delta$ and $\widetilde{h}$ belongs to $U_\delta$. $\square$
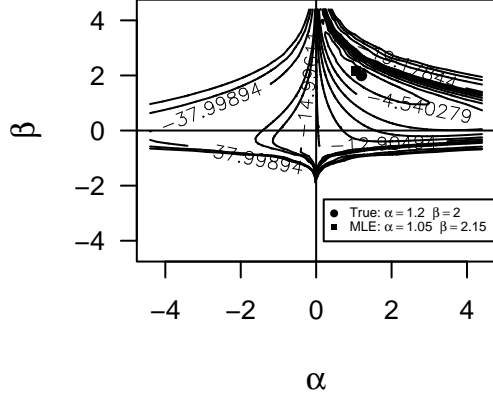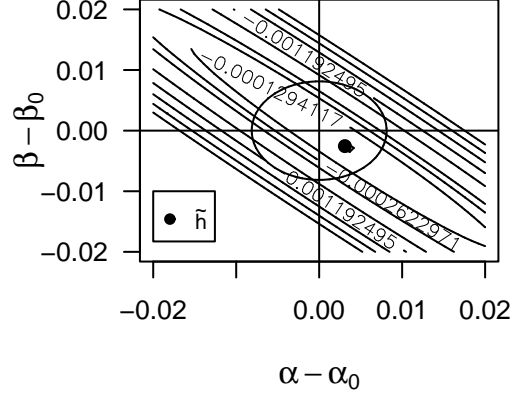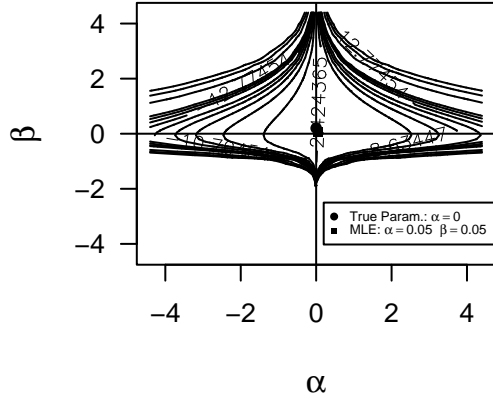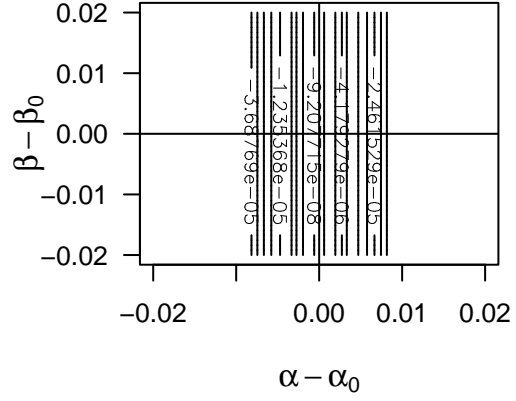
Therefore, under conditions $\mathfrak{F}1 - \mathfrak{F}5$, for large enough n, and $\delta < \mu_1/\tau$ there exists a (unique) consistent solution to the likelihood equations. Moreover, by a straightforward generalization of Huzurbazar's results (1948), we can show that $\widetilde{h}$ maximizes the log-likelihood funtion.

It is interesting to look at the consistency from a geometrical point of view. In ensuring the consistency of a solution to the likelihood equations it is important that $z(\theta)$ should have a maximum turning value at $\theta_0$ and that for $\delta$ sufficiently small $Q(\theta)$ has a unique maximizing point in $U_\delta$.

$Fig.\,1$ shows the simulated contour lines of $z(\alpha, \beta)$ for the model (1) with true parameter $\alpha_0 = 1.2$ and $\beta_0 = 2$, sample size 1000 and $x$ generated by random selection from a uniform distribution in the interval $(0, 1)$ held fixed on replications of samples. In the graph the true parameter is marked with a black point while the maximum likelihood estimate (MLE) is marked with a black square. $z(\alpha, \beta) - z(1.2, 2) < 0$ in a neighbourhood of $\theta_0 = (1.2, 2)$ with a maximum of the function $z(\alpha, \beta)$ which occurs in that point and is equal to $z(1.2, 2) = log(1/\sqrt{2\pi}) - (1/2) = -1.4189$.

Equation (5) may be seen as the first order necessary conditions for the unconstrained maximum of the quadratic approximation $Q(\theta)$. Therefore, by $Lemma\,2$, the positive definiteness of $B(\theta_0)$ and $\delta < \mu_1/\tau$ ensure the existence of such maximizing point, unique in $U_\delta$.

From a geometric point of view, the locus of points $Q(\theta) - z(\theta_0)$ is a quadric surface. Let define the set

5

Fig. 1: z(α,β)



Fig.2: Q(α,β)−z(α₀,β₀)

True: α = 1.2  β = 2
MLE: α = 1.05  β = 2.15



Fig. 3: z(α,β)



Fig.4: Q(α,β)−z(α₀,β₀)

True Param.: α = 0
MLE: α = 0.05  β = 0.05

$$C_\theta \equiv \{\theta\,;\; Q(\theta) - z(\theta_0) < 0\} = \left\{\theta\,;\; -\frac{1}{2}h'B(\theta_0)h + h'm(x)\delta^2 < 0\right\} \qquad (6)$$

then, in the regular case, we have the following *Lemma*

**Lemma 3.** *For any δ, $C_\theta$ is non empty and bounded.*

*Proof.* Because of the positive definiteness of $B(\theta_0)$ the quadric surface may be written as

$$Q(\theta) - z(\theta_0) = -\frac{1}{2}(h - \widetilde{h})'B(\theta_0)(h - \widetilde{h}) + \frac{1}{2}d < 0$$

6

where $\widetilde{h} = B^{-1}(\theta_0)m(x)\delta^2$ and $d = \delta^2 m'(x)B^{-1}(\theta_0)m(x)\delta^2 > 0$. By the spectral decomposition $B(\theta_0) = P'\Lambda P$ where $\Lambda = diag(\mu_1, \cdots, \mu_k)$ and P is an orthogonal matrix, the inequality $Q(\theta) - z(\theta_0) < 0$ may be reexpressed as

$$\sum_{i=1}^{k} \left( \frac{z_i}{\beta_i} \right)^2 > d$$

where $z = P(h - \widetilde{h})$ and $\beta_i = \sqrt{1/\mu_i}$. The transformation $z = P(h - \widetilde{h})$ represents a translation followed by a rotation, so $C_\theta$ is equivalent to the set

$$C_z \equiv \left\{ z; \ \sum_{i=1}^{k} \left( \frac{z_i}{\beta_i} \right)^2 > d \right\}$$

Because $d$ is greater than zero, the above inequality describes a non empty area inside an ellipsoid with center $\widetilde{h}$. Therefore, this area is a bounded set. $\qquad\square$

To guarantee the statistical consistency of a solution to the likelihood equations it is necessary that the center of the ellipsoid is in $U_\delta$. This condition is ensured if $\delta$ is taken sufficiently small. This may be shown in the following way. For every $\theta$ such that $\| \theta - \theta_0 \| = \delta$ we have

$$-\frac{1}{2}h'B(\theta_0)h + h'm(x)\delta^2 \leq -\frac{1}{2}\mu_1 \| h \|^2 + \| h \| \tau\delta^2 = -\frac{1}{2}\mu_1\delta^2 + \tau\delta^3 = \delta^2 \left( -\frac{1}{2}\mu_1 + \tau\delta \right)$$

which is less than zero if $\delta < (1/2)\mu_1/\tau$. $Q(\theta) - z(\theta_0) < 0$ implies $-h'B(\theta_0)h + h'm(x)\delta^2 < 0$ and the inequality is still valid if we divide both terms by $\delta$. Therefore, for every $\theta$ such that $\|\theta - \theta_0\| = \delta$ and $\delta$ arbitrarily small (in this case it is sufficient $\delta < \mu_1/\tau$) we have

$$-\frac{1}{\delta}h'B(\theta_0)h + \frac{1}{\delta}h'm(x)\delta^2 = \frac{1}{\delta}h'\big( -B(\theta_0)h + m(x)\delta^2 \big) = \eta'g(\eta) < 0 \qquad (7)$$

where $\eta = h/\delta$ with $\| \eta \| = 1$ and $g(\eta) = -B(\theta_0)(\theta - \theta_0) + m(x)\delta^2$ is a continuous function mapping $\mathbb{R}^k$ into itself. To get on we require the following Lemma

**Lemma 4.** *(Aitchison and Silvey, 1958) If g is a continuous function mapping $\mathbb{R}^k$ into itself with the property that, for every $\theta$ such that $\|\theta\| = 1$, $\theta'g(\theta) < 0$, then there exists a point $\widehat{\theta}$ such that $\|\widehat{\theta}\| < 1$ and $g(\widehat{\theta}) = 0$ .*

*Proof.* This result is equivalent to Brouwer's fixed point Theorem. See Aitchison and Silvey (1958) for a complete proof. $\qquad\square$

Therefore, by $Lemma$ 4, in our case we can say that there exists a point $\widetilde{\eta} = (\widetilde{\theta} - \theta_0)/\delta$ such that $\| \widetilde{\eta} \| < 1$ that is, $\| \widetilde{\theta} - \theta_0 \| < \delta$ and $g(\widetilde{\eta}) = -B(\theta_0)(\widetilde{\theta} - \theta_0) + m(x)\delta^2 = 0$. Moreover, because $B(\theta_0)$ is positive definite, this point is unique.

$Fig.$ 2 shows the simulated behavior of the quadratic form $Q(\alpha, \beta) - z(\alpha, \beta)$, with a sample size 1000. Given $\delta$ the center of the ellipsoid is in $U_\delta$ and because $\delta$ may be chosen arbitrarily small this point is as close as we like to the center of $U_\delta$ ensuring in this way the consistency of the estimator.

As to the asymptotic distribution of the maximum likelihood estimator, taking the probability limit of equation (3) after replacing $h$ by $\widetilde{h}$, we have

$$plim\left(\frac{1}{n}D^2 logL(\theta_0) + \frac{R^*}{2n}\right)\sqrt{n}\,\widetilde{h} = -\eta \tag{8}$$

where $\eta \sim N(0, B(\theta_0))$ is the asymptotic distribution of the score scaled by $n^{-1/2}$ and $R^*$ is a matrix whose $i^{th}$ component may be expressed as $\widetilde{h}'\Delta_i(\theta^*)$ and $\theta^*$ a point such that $\| \theta^* - \theta_0 \| < \| \theta - \theta_0 \|$.

Under above conditions $plim\,(n^{-1})D^2 logL(\theta_0) = -B(\theta_0)$. Moreover, because of the consistency of the estimator, $plim\,R^*/2n = o_p(1)$ so that $plim\,n^{1/2}\,\widetilde{h} = B^{-1}(\theta_0)\eta$ and asymptotically $n^{1/2}\,\widetilde{h} \sim N(0, B^{-1}(\theta_0))$.

## 3  Singular Information Matrix

As known, the whole problem of maximum likelihood estimation is closely bound up with the behavior of the function $z(\theta)$ which should have a unique maximum at $\theta_0$ (local asymptotic identifiability condition). The demands that $z(\theta)$ is a maximum at $\theta_0$ and that the information matrix in an observation, $B(\theta_0)$, is positive definite are related. In fact, under regularity conditions on $f(t; \theta)$, a Taylor series expansion of $z(\theta)$ about $\theta_0$, yields

$$z(\theta) - z(\theta_0) = -\frac{1}{2}h'B(\theta^*)h, \qquad \|\theta^* - \theta_0\| < \|\theta - \theta_0\|$$

so that $B(\theta^*)$ is positive definite if $z(\theta) - z(\theta_0) < 0$ in a neighbourhood of $\theta_0$. If one assumes that the rank of $B(\theta)$ does not change in an open neighborhood of $\theta_0$ (the Rothenberg's regularity condition of $B(\theta)$ in $\theta_0$), then one can conclude that $B(\theta_0)$ is positive definite. Moreover, if $B(\theta)$ is regular in a neighborhood of $\theta_0$, the positive definiteness of $B(\theta_0)$ implies local identifiability of $\theta_0$.

The singularity of $B(\theta_0)$, only by itself, does not necessarily imply the local unidentifiability of $\theta_0$. This fact can be understood from a Taylor series expansion of $z(\theta)$ near $\theta_0$,

$$z(\theta) - z(\theta_0) = -\frac{1}{2}h'B(\theta_0)h + O(\|\theta - \theta_0\|^3)$$

the higher order terms can ensure that $z(\theta) - z(\theta_0) < 0$ for every $\theta \neq \theta_0$ in a neighbourhood of $\theta_0$, even though the quadratic form in the above expression be null.

Moreover, in some statistical applications $\theta_0$ is identified but $B(\theta)$ does not satisfy the Rothenberg's regularity condition in $\theta_0$. It may happen that $B(\theta^*)$ is of full rank and positive definite for some $\theta^*$ in a neighborhood of $\theta_0$ while $B(\theta_0)$ is of lower rank.

Sometimes, we could be interested in doing inference on the parameter of interest $\alpha_0$, say. In this case it may happen that even if $B(\theta_0)$, $\theta_0' = [\alpha_0' \ \beta_0']$, is not positive definite and $z(\theta_0)$ is not a maximum turning value of $z(\theta)$, it may still be the case that setting the (nuisance) parameter $\beta$ equal to some constant, $\beta^\circ$, the information matrix $B(\alpha_0, \beta^\circ)$ is positive definite and if Rothenberg's condition is satisfied in the point $(\alpha_0, \beta^\circ)$, $z(\alpha, \beta^\circ) - z(\alpha_0, \beta^\circ) < 0$ in a neighbourhood of $\alpha_0$. Often this situation occurs for any $\beta$ as in the so called "indeterminate parameter problem" (Cheng and Traylor, 1995) where the information matrix is usually block diagonal with the northwest submatrix $B_{11}(\alpha_0, \beta)$ positive definite satisfying the Rothenberg's condition in $\alpha_0$, $\forall \beta$ and the southeast submatrix $B_{22}(\alpha_0, \beta) = 0$. In this case $z(\alpha, \beta) - z(\alpha_0, \beta) < 0$ and $\alpha_0$ is identified for any $\beta$.

$Fig.\,3$ shows the simulated contour lines of $z(\alpha \ \beta)$ for the model (1) with a sample size 1000. As said above the information matrix in an observation is now singular and the graph shows a whole set of maxima of the function $z(\alpha\,, \beta)$ in the point $\theta_0 = (0, \beta)$. In this case $z(\alpha, \beta) = log(1/\sqrt{2\pi}) - (1/2)E_0(y - \alpha x^\beta)^2$. If $\alpha = 0$, $E_0(y - \alpha x^\beta)^2 = \sigma^2 = 1$ for any $\beta$ while if $\alpha \neq 0$, $E_0(y - \alpha x^\beta)^2 = 1 + T^2$ where T is a constant. Therefore, $z(\alpha, \beta) - z(0, \beta) = -(1/2)T^2 < 0$ in a sufficiently small neighborhood of $\alpha = 0$ and $B_{11}(0, \beta) = x^{2\beta}$ is positive (definite) for any $\beta$.

Then, when the information matrix is singular but the parameter $\theta_0' = [\alpha_0' \ \beta']$ is locally ($\beta = \beta_0$) or partially ($\forall \beta$) identified, we can ask how to do inference on the whole set of parameter $\theta$ or on the parameter of interest, $\alpha$. In this regard the starting point is the analysis of the asymptotic properties of maximum likelihood estimator when the information matrix is singular.

Let begin with the statistical consistency. We can observe that $Lemma\,1$ is still valid because the asymptotic result given by equation (5) does not involve the assumption on the singularity of the information matrix. The problem rises with $Lemma\,2$. More precisely, the problem concerns the existence of a unique solution in $U_\delta$ that satisfies equation (5). This system can be algebraical consistent or not, in both cases we can write an (approximate) solution in the following form

$$\widetilde{h} = \widetilde{h}_1 + \widetilde{h}_2 = B^+(\theta_0)m(x)\delta^2 + \left[I - B^+(\theta_0)B(\theta_0)\right]u \qquad (9)$$

for some $u$ where $B^+(\theta_0)$ is the Moore-Penrose pseudoinverse of $B(\theta_0)$ and $\widetilde{h}_2 = [I - B^+(\theta_0)B(\theta_0)]u$ is the projection of $u$ on the $kernel$ of $B(\theta_0)$. It is well known that if the system (5) is algebraical consistent then $(9)$ is a solution to the system, otherwise it is a solution which minimize $\|m(x)\delta^2 - B(\theta_0)h\|^2$.

Of course there is no guarantee that $\widetilde{h}$ is unique in $U_\delta$ unless we could say something

on the norm of the arbitrary vector $u$ and on the information matrix. In this regard from the following inequalities

$$\|B^+(\theta_0)m(x)\delta^2\| < \mu_{min}^{-1}\tau\delta^2 \quad and \quad \left\|\left[I - B^+(\theta_0)B(\theta_0)\right]u\right\| \le \|u\|$$

where $\mu_{min}$ is the minimum eigenvalue non zero of $B(\theta_0)$, we have

$$\|\widetilde{h}\| < \mu_{min}^{-1}\tau\delta^2 + \|u\|$$

Let $\|u\| = \xi\,\delta$ be the norm of $u$ with $\xi > 0$. Then, $\|\widetilde{h}\| < \delta$ if $\mu_{min}^{-1}\tau\delta + \xi < 1$ that is, if $\xi < 1 - \mu_{min}^{-1}\tau\delta$ which is valid if $\delta < \mu_{min}/\tau$. Therefore, if $\delta$ is sufficiently small and $\|u\| < (1 - \mu_{min}^{-1}\tau\delta)\delta$, we can define a neighborhood of $\theta_0$ where we have a solution of equation (5). However, in this neighborhood $\widetilde{h}$ is not unique because the $kernel$ of $B(\theta_0)$ does not consist only of the zero vector, that is, $[I - B^+(\theta_0)B(\theta_0)]u$ does not vanish for all $u$ in $U_\delta$. Then, when the information matrix is singular, a solution to the likelihood equation is not statistically consistent.

To detect the asymptotic distribution of the maximum likelihood estimator we refer to $(8)$. Taking the probability limit of the expressions on the left-hand side of $(8)$, problems rise with $R^*/2n$ which is now a quantity $O_p(1)$ because the estimator is no more consistent. Then, we have

$$plim\left(\frac{1}{n}D^2logL(\theta_0) + \frac{R^*}{2n}\right)\sqrt{n}\,\widetilde{h} = [-B(\theta_0) + F]\,plim\left(\sqrt{n}\,\widetilde{h}\right) = -\eta$$

where the symbols are the same as in $(8)$. From above equality we observe that if the information matrix is singular nothing we know about the invertibility of the matrix $[-B(\theta_0) + F]$ and we can not derive the asymptotic distribution of $n^{1/2}\,\widetilde{h}$.

From a geometric point of view, the set $C_\theta$ defined by $(6)$ is still non empty but unlike the regular case, now it is unbounded. In fact, when the information matrix is singular the inequality $-\frac{1}{2}h'B(\theta_0)h + h'm(x)\delta^2 < 0$ may be reexpressed as (Shilov, 1977, p.288)

$$-\frac{1}{2}\sum_{i=1}^{r}\eta_i\,z_i^2 + \sum_{i=r+1}^{k}z_i\,\gamma_i + d < 0$$

where, $\eta_i,\ i = 1,\cdots,r$ are the r eigenvalues of $B(\theta_0)$ greater than zero, $\gamma = \delta^2 Pm(x)$, $d = (1/2)\sum_{i=1}^{r}(\gamma_i^2/\eta_i) > 0$ and

$$z = Ph + u \qquad with \qquad u_i = \begin{cases} -\gamma_i/\eta_i & \text{if } \eta_i > 0 \\ 0 & \text{if } \eta_i = 0 \end{cases}$$

In the new space given by the transformation $z = Ph + u$ the quadric surface may take many forms according to the number of non-zero eigenvalues (Shilov, 1977, p. 295). However, this set will always be non empty and unbounded.

For example, with one eigenvalue zero, in $\mathbb{R}^2$ we have a pair of parallel lines, in $\mathbb{R}^k$ with one eigenvalue zero, the surface is generated by translating the ellipsoid described by the remaining $k-1$ eigenvalues in the $(k-1)$-dimensional space along a perpendicular to $\mathbb{R}^{k-1}$. In a three-dimensional space we must translate the ellipses in a two dimensional space along a third axis giving elliptic cylinders. $Fig.\,4$ shows the contour lines of the quadratic form for the model $(1)$.

# 4    A Solution to the Singularity of the Information Matrix

As said above, in the identification problem Silvey (1959) proposed to replace the singular information matrix $B(\theta_0)$ by $B(\theta_0) + F$ where $F$ is an appropriate matrix obtained imposing some restrictions on the parameters of the model so that the restricted parameters are identified and the new matrix is positive definite. To generalize Silvey's approach we suggest to modify the information matrix adding an arbitrary positive constant $\lambda^2$ to the diagonal element of $B(\theta_0)$ producing $A_\lambda(\theta_0) = B(\theta_0) + \lambda^2\, I$ where $I$ is an identity matrix of appropriate dimension. To investigate the consequences of this transformation we replace $B(\theta_0)$ by $A_\lambda(\theta_0)$ wherever it appears in the regular theory.

## 4.1    An Unfeasible Solution

By construction, $A_\lambda(\theta_0)$ is positive definite with eigenvalues given by $\mu_i + \lambda^2$, $i = 1, \cdots, k$, $\mu_i \geq 0$ and $\lambda > 0$ arbitrarily chosen.

Consider first what happens to the quadratic approximation $(4)$. Adding and subtracting the quantity $\frac{1}{2}\lambda^2 \parallel \theta - \theta_0 \parallel^2$ to $(2)$, taking the probability limit of both sides and using conditions $\mathfrak{F}1 - \mathfrak{F}4$, we have that for large $n$, $n^{-1}logL(\theta) - \frac{1}{2}\lambda^2 \parallel \theta - \theta_0 \parallel^2$ can be approximated by the following quadratic model,

$$P(\theta, \lambda) \equiv Q(\theta) - \frac{1}{2}\lambda^2 \parallel \theta - \theta_0 \parallel^2 \tag{10}$$

$P(\theta, \lambda)$ may be seen as a penalty function given by $Q(\theta)$ "penalized" by a quadratic term, $\parallel \theta - \theta_0 \parallel^2$, with a penalty parameter $\lambda^2$. If we maximize $(10)$, by imposing the first order necessary conditions we get

$$-(B(\theta_0) + \lambda^2\, I)h + m(x)\delta^2 = -A_\lambda(\theta_0)h + m(x)\delta^2 = 0 \tag{11}$$

$A_\lambda(\theta_0)$ is positive definite for any $\lambda > 0$, and $(11)$ is an algebraic $consistent$ system of equations with a unique solution given by

$$\widehat{h}_\lambda = \left(\widehat{\theta}_\lambda - \theta_0\right) = \left(B(\theta_0) + \lambda^2 I\right)^{-1} m(x)\delta^2$$

Because $\| A_\lambda^{-1}(\theta_0)m(x) \| \leq \lambda^{-2}\tau$, we have $\| \widehat{h}_\lambda \| \leq \lambda^{-2}\tau\delta^2$. If $\delta < \lambda^2\tau^{-1}$ then $\| \widehat{h}_\lambda \|$ is in $U_\delta$. Therefore, given $\lambda > 0$ there always exists a $\delta$ sufficiently small such that $P(\theta, \lambda)$ has a unique maximizing point in a neighborhood of $\theta_0$.

In this case $P(\theta, \lambda)$ plays the same role as $Q(\theta)$ for the regular case and equation (11) may be seen as an asymptotic result of a Taylor series expansion about $\theta_0$ of what we call "penalized" likelihood equations. That is, if we maximize the following "penalized" likelihood function

$$\frac{1}{n}logL(\theta) - \frac{1}{2}\lambda^2 \| \theta - \theta_0 \|^2 \qquad \lambda > 0$$

then, by imposing the first order necessary conditions, we get the "penalized" likelihood equations given by

$$\frac{1}{n}DlogL(\theta) - \lambda^2(\theta - \theta_0) = 0 \tag{12}$$

that now plays the same role as the likelihood equations for the regular case. Then, we can restate $Lemma$ 1 as follows

**Theorem 1.** *Subject to the conditions $\mathfrak{F}1 - \mathfrak{F}4$, for large enough n and $\delta$ sufficiently small, the "penalized" likelihood equations have a solution $\widehat{h}_\lambda = \widehat{\theta}_\lambda - \theta_0 \in U_\delta$ if (and only if) it satisfies a certain equation of the form given by $(11)$ where $\lambda > 0$, $m(x)$ is a continuous function on $U_\delta$ and $\| m(x) \|$ is bounded in $U_\delta$ by a positive number $\tau$, say.*

*Sketch of the Proof.* A Taylor series expansion of $(12)$ about $\theta_0$ gives

$$\frac{1}{n}DlogL(\theta_0) + \left(\frac{1}{n}D^2logL(\theta_0) - \lambda^2I\right) h + \frac{1}{2}V(x;\theta^*) = 0 \tag{13}$$

Then, under conditions $\mathfrak{F}1 - \mathfrak{F}4$, equation (11) is obtained following the same lines of reasoning as in the regular case. $\qquad\square$

Then, above arguments allow one to state that a solution to the "penalized" likelihood equations, $\widehat{h}_\lambda$ is statistically consistent for any $\lambda > 0$. Moreover, following the same line of reasoning as in the regular case, it is immediate to show that asymptotically $n^{1/2}\widehat{h}_\lambda \sim N(0, V)$ where $V = A_\lambda^{-1}(\theta_0)B(\theta_0)A_\lambda^{-1}(\theta_0)$ is singular with $Rank(V) = r$.

From a geometric point of view we have the same situation as in the regular case with $B(\theta_0)$ replaced by $A_\lambda(\theta_0)$. In fact, the locus of points $P(\theta, \lambda) - z(\theta_0)$ is a quadric surface and the set defined as in (6) is equivalent to an area inside an ellipsoid with center in $U_\delta$. In fact,

$$\left\{\theta\,;\; -\frac{1}{2}h'A_\lambda(\theta_0)h + h'm(x)\delta^2 < 0\right\} \Leftrightarrow \left\{z\,;\; \sum_{i=1}^{r}\left(\frac{z_i}{\gamma_i}\right)^2 + \sum_{i=r+1}^{k}\left(\frac{z_i}{\sqrt{(1/\lambda^2)}}\right)^2 > d\right\}$$

where $\gamma_i = \sqrt{1/(\mu_i + \lambda^2)}$, $\lambda > 0$ and $\mu_i > 0$. This set is non empty and bounded.

As it emerges looking at the "penalized" likelihood equations, the main problem connected to the estimator proposed is its feasibility because given $\lambda$ the search of a solution to (12) depends on the unknown true parameter. In the paper the problem is solved fixing appropriately the magnitude of $\lambda$ so that the knowledge of $\theta_0$ is unnecessary.

## 4.2 A Feasible Solution

Our assumption is to take $\lambda$ small enough, formally $\lambda \to 0$. In this case we must investigate the consequences of this assumption on the asymptotic properties of a solution to the "penalized" likelihood equations given by

$$\lim_{\lambda \to 0} \left[ \frac{1}{n} DlogL(\theta) - \lambda^2(\theta - \theta_0) = 0 \right] \tag{14}$$

The main result is the following Theorem

**Theorem 2.** *Let $Rank(B(\theta_0)) = r < k$. Subject to the conditions $\mathfrak{F}1 - \mathfrak{F}4$ for large enough $n$ and $\delta$ sufficiently small, equations (14) have a (unique) solution, $\lim_{\lambda \to 0} \widehat{h}_\lambda = \widehat{h}_{\lambda 0}$ in $U_\delta$ if (and only if) it satisfies a certain equation of the form*

$$\lim_{\lambda \to 0} \left[ -(B(\theta_0) + \lambda^2 I)h + m(x)\delta^2 = 0 \right] \tag{15}$$

*Moreover,*

$$\lim_{\lambda \to 0} \sqrt{n} \left( \widehat{\theta}_\lambda - \theta_0 \right) = \sqrt{n} \left( \widehat{\theta}_{\lambda 0} - \theta_0 \right) = \sqrt{n} \ \widehat{h}_{\lambda 0} \sim N \left( 0, B^+(\theta_0) \right)$$

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

$$W_0 = n \ \widehat{h}'_{\lambda 0} \ B(\theta_0) \ \widehat{h}_{\lambda 0} \sim \chi^2(r)$$

*Proof.* The if and only if part of the theorem is immediate following the "regular" case. We show that as $\lambda \to 0$, (15) has a unique solution in $U_\delta$. In previous section we have seen that for any $\lambda$, $\widehat{h}_\lambda$ is unique in $U_\delta$ if $\delta < \lambda^2 \tau^{-1}$ or if $\delta = \xi \lambda^2 \tau^{-1}$ with $0 < \xi < 1$. Therefore, we can reexpress $\widehat{h}_\lambda$ as

$$\widehat{h}_\lambda = \left( B(\theta_0) + \lambda^2 I \right)^{-1} m(x)\xi\lambda^2\tau^{-1}\delta = \left( B(\theta_0) + \lambda^2 I \right)^{-1} \lambda^2 s(x)\delta$$

where $s(x) = \xi m(x)\tau^{-1}$ with $\|s(x)\| < 1$. But

$$\left( B(\theta_0) + \lambda^2 I \right)^{-1} \lambda^2 = I - \left( B(\theta_0) + \lambda^2 I \right)^{-1} B(\theta_0)$$

13

and (Albert, 1971, p.19)

$$\lim_{\lambda \to 0} \left(B(\theta_0) + \lambda^2 I\right)^{-1} \lambda^2 = I - \lim_{\lambda \to 0} \left(B(\theta_0) + \lambda^2 I\right)^{-1} B(\theta_0) = I - P_B$$

where $P_B = B^+(\theta_0)B(\theta_0)$ and $I - P_B$ is a projector on the *kernel* of $B(\theta_0)$. Therefore,

$$\lim_{\lambda \to 0} \widehat{h}_\lambda = \widehat{h}_{\lambda 0} = \left(I - P_B\right) s(x)\delta$$

that is, $\widehat{h}_{\lambda 0}$ is the projection of the vector $s(x)\delta$ on the *kernel* of the information matrix with $\delta$ taken arbitrarily small, formally $\delta \to 0$. Moreover,

$$\|\widehat{h}_{\lambda 0}\| = \| \left(I - P_B\right) s(x)\delta\| < \|s(x)\|\delta < \delta \qquad \delta \to 0$$

which proves the first part of the Theorem.

As to the asymptotic distribution, we apply the probability limit to (13) after replacing $h$ by $\widehat{h}_\lambda$, letting $\lambda \to 0$ and following the same lines of reasoning as in the regular case. Then, we have that $\lim_{\lambda \to 0} \sqrt{n} \left(\widehat{\theta}_\lambda - \theta_0\right)$ tends in distribution to a random vector $\lim_{\lambda \to 0} \left(B(\theta_0) + \lambda^2 I\right)^{-1} \eta$ where $\eta \sim N\left(0, B(\theta_0)\right)$. Therefore, asymptotically

$$\lim_{\lambda \to 0} \sqrt{n} \left(\widehat{\theta}_\lambda - \theta_0\right) \sim N\left(0, \lim_{\lambda \to 0} \left(B(\theta_0) + \lambda^2 I\right)^{-1} B(\theta_0) \left(B(\theta_0) + \lambda^2 I\right)^{-1}\right)$$

It is immediate to show that (Albert, 1972)

$$\lim_{\lambda \to 0} \left(B(\theta_0) + \lambda^2 I\right)^{-1} B(\theta_0) \left(B(\theta_0) + \lambda^2 I\right)^{-1} = B^+(\theta_0)$$

where $B^+(\theta_0)$ always exists and is unique.

Finally the last part of the Theorem. By the properties of $B^+(\theta_0)$, the matrix $B(\theta_0)B^+(\theta_0)$ is idempotent, then

$$Rank\left(B(\theta_0)B^+(\theta_0)\right) = tr\left(B(\theta_0)B^+(\theta_0)\right) = tr\left(P\Lambda P'P\Lambda^+P'\right) = tr\left(P\Lambda\Lambda^+P'\right)$$

with $\Lambda^+ = diag\left(\mu_1^+, \mu_2^+, \cdots, \mu_k^+\right)$ where

$$\mu_j^+ = \left\{ \begin{array}{cc} \mu_j^{-1} & \text{if } \mu_j > 0 \\ 0 & \text{if } \mu_j = 0 \end{array} \right.$$

Therefore, $Rank\left(B(\theta_0)B^+(\theta_0)\right) = r = Rank\left(B(\theta_0)\right)$. Moreover,

$$B^+(\theta_0)B(\theta_0)B^+(\theta_0)B(\theta_0)B^+(\theta_0) = B^+(\theta_0)B(\theta_0)B^+(\theta_0)$$

then by a Theorem on the quadratic forms (Searle, 1971, p. 69), the chi-square distribution follows. □

We conclude this section discussing the relationship between $\widehat{h}_{\lambda 0}$ and a solution to the likelihood equations given by (9) in the case of a singular information matrix. If we replace the arbitrary vector $u$ in (9) by $s(x)\delta$, we have

$$\widetilde{h} = B^+(\theta_0)m(x)\delta^2 + \left[I - B^+(\theta_0)B(\theta_0)\right]s(x)\delta \qquad (16)$$

that is,

$$\widetilde{h} - \left[I - B^+(\theta_0)B(\theta_0)\right]s(x)\delta = \widetilde{h} - \widehat{h}_{\lambda 0} = B^+(\theta_0)m(x)\delta^2$$

but $\widetilde{h} - \widehat{h}_{\lambda 0} = \widetilde{\theta}_n - \widehat{\theta}_{\lambda 0}$ therefore, in $U_\delta$ we have

$$\|\widetilde{\theta}_n - \widehat{\theta}_{\lambda 0}\| = \|B^+(\theta_0)m(x)\|\delta^2 < \delta$$

with $\delta$ arbitrarily small. Letting $\delta \to 0$, by (16) we get

$$\widetilde{h} = \left[I - B^+(\theta_0)B(\theta_0)\right]s(x)\delta + o(\delta)$$

and a unique solution in $U_\delta$ is found in the $kernel$ of $B(\theta_0)$.

# 5   Numerical solution of penalized likelihood equations

A first order approximation to (14) about $\theta_0$ gives

$$\lim_{\lambda \to 0}\left[\frac{1}{n}DlogL(\theta_0) - \left(-\frac{1}{n}D^2logL(\theta_0) + \lambda I\right)(\theta - \theta_0) = 0\right]$$

that is,

$$\lim_{\lambda \to 0}\left[\theta - \theta_0 = \left(-\frac{1}{n}D^2logL(\theta_0) + \lambda I\right)^{-1}\frac{1}{n}DlogL(\theta_0)\right]$$

then, we propose the following algorithm

(i) Fix a decreasing sequence $\{\lambda_i\}$, typically $\{1, 10^{-1}, 10^{-2}, \cdots\}$ and choose a starting point $\theta^{(r)}$.

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

(iii) Find iteratively a solution to

$$\theta^{(r+1)} = \theta^{(r)} + \left(-\frac{1}{n}D^2logL(\theta^{(r)}) + \lambda_i I\right)^{-1}\frac{1}{n}DlogL(\theta^{(r)})$$

call $\theta^{(F)}$ such solution.

15

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

An estimate of the information matrix $B(\theta_0)$ can be computed replacing $\theta_0$ by $\widehat{\theta}_{\lambda 0}$.

A simulation applied to the Engle's model is presented to support the theoretical results. Fig. 3(a) shows the simulated distribution of an estimate of $\alpha$ obtained as a solution to the penalized loglikelihood equation from 100 generated random samples of size 1000. This estimate is compared with an underlying normal distribution. In Fig. 3(b) the cumulative distribution of an estimate of $W_0$, $\widehat{W}_0$ is compared with a $\chi^2(1)$ distribution. From Figure 3 it emerges the good fits of the simulated distributions.
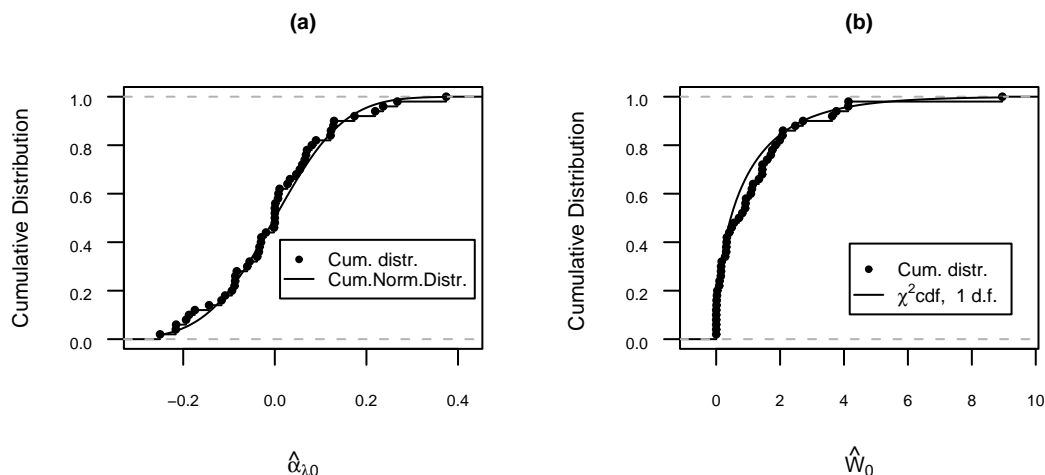


Figure 3: Simulated cumulative distribution functions of $\widehat{\alpha}_{\lambda 0}$ (graphic $(a)$) and of $W_0$ (graphic $(b)$) for the Engel's model. $H_0 : \alpha = 0$, sample size 1000, 100 replications.

# 6  Conclusions

In this paper we proposed a way to solve the singularity of the information matrix. The approach is based on the definition of a penalized loglikelihood function letting the penalty parameter going to zero. In this way we get a solution in a neighborhood of the maximum likelihood estimate with attractive statistical properties. More precisely, the estimator 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.

16

# References

Aitchison, J. and Silvey, S. D. (1958) Maximum-likelihood estimation of parameters subject to restraints. *The Annals of Mathematical Statistics*, **29**, 813–828.

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

Cheng, R. C. H. and Traylor, L. (1995) Non-regular maximum likelihood problems. *Journal of Royal Statistical Society, B*, **57**, 3–44.

Davies, R. B. (1977) Hypothesis testing when a nuisance parameter is present only under the alternative. *Biometrika*, **64**, 247–254.

— (1987) Hypothesis testing when a nuisance parameter is present only under the alternative. *Biometrika*, **74**, 33–43.

Engle, R. F. (1984) Wald, likelihood ratio and lagrange multiplier tests in econometrics. In *Handbook of Econometrics* (eds. Z. Griliches and M. Intriligator), vol. 2, 775–826. Amsterdam: North Holland.

Huzurbazar, V. S. (1948) The likelihood equation, consistency and the maxima of the likelihood function. *Annals of Eugenics*, **14**.

Poskitt, D. S. and Tremayne, A. R. (1981) An approach to testing linear time series models. *The Annals of Statistics*, **9**, 974–86.

Rotnitzky, A., Cox, D. R., Bottai, M. and Robins, J. (2000) Likelihood-based inference with singular information matrix. *Bernoulli*, **6**, 243–284.

Searle, S. R. (1971) *Linear Models*. New York: Wiley.

Shilov, G. E. (1977) *Linear Algebra*. New York: Dover.

Silvey, S. D. (1959) The lagrange multiplier test. *The Annals of Mathematical Statistics*, **30**, 389–407.