

Is this the least squares estimate?

By EUGENE DEMIDENKO

Dartmouth Medical School and Dartmouth College
7927 Rubin Building, One Medical Center Drive, Lebanon, NH 03756, USA
e-mail: eugene.demidenko@dartmouth.edu

To be published in Biometrika

Summary

It is shown that the sum of squares can have several local minima with a positive probability for any intrinsically nonlinear regression with infinite tails. Therefore, the availability of global criteria is crucial. The concept of the local convexity level for sum of squares in nonlinear regression model is introduced. A general formula for the local convexity level of the sum of squares is derived and link to the curvature of nonlinear model is established. It is shown that the local convexity level is equal to the minimum of the squared radius of the full curvature of the expectation surface of nonlinear regression. Two general global criteria are formulated. The calculation of the local convexity level and construction of global criteria are illustrated by four types of nonlinear models: polylinear, power, linear M-regression, and exponential regression models. The suggested global criteria work well for real life data.

Some key words: nonlinear regression; uniqueness; convexity; unimodality; robust regression; curvature.

1. Introduction

Only in rare cases estimation procedure admits a closed form solution, the least squares estimate in linear regression model is such a fortunate example. Then, after iterations on minimization of the sum of squares in nonlinear regression are done and we look at the final parameter value, a reasonable question arises: "Is this the least squares estimate?", i.e. did we find the global minimum? Amazingly, until today it is extremely difficult to answer this simple to formulate question. Moreover, as is shown in this paper, it is quite possible for a sum of squares to have at least two local minima.

The global optimization is under constant interest, particularly a special journal has been established recently, Journal of Global Optimization. Several books have been published on the topic. The objective of the global optimization is to find the global minimum of

a continuous function (possibly under some restrictions). Two approaches are available: deterministic and stochastic. In deterministic approach the parameter set is divided into a finite number of small subsets where it is known there is a unique local minimum, e.g. Horst and Tuy (1996), Floudas and Pardalos (1992), Hansen (1992). The drawbacks of this approach is that it requires a compact parameter set and the number of the divided subsets may be very large. In stochastic global search the global minimum is found with a certain probability, Mockus (1989), Zhigljavsky (1991). In several books merits of different global optimization techniques are compared, Horst and Pardalos (1995), Kearfort (1996), Pinter (1996). The lack of well working global criteria could be explained by the fact that all general criteria, i.e. applied to a general class of optimization functions cannot work because they are not specific, Horst and Tuy (1996). The success in construction of a global criterion is determined by the class of nonconvex functions it covers. The sum of squares for nonlinear regression is a very interesting and broad class of multiextremal functions yet allowing construction of working global criteria.

We argue that in practical optimization, particularly in optimization driven by statistical problems we can rarely expect that the minimized function has many local minima. Thus there is just a possibility of multiextremality. Consequently, instead of doing a very expensive global search we need just criteria for global minimum, i.e. conditions under which a found local minimum is in fact the global one. To work out such criteria is the goal of the present paper. In particular, links between the full curvature of the nonlinear model and local convexity of the sum of squares are established: the minimum of the squared radius of the full curvature of the expectation surface is equal to the level of the local convexity of the sum of squares. Based on the introduced concepts we are able to formulate global criteria, i.e., to propose conditions under which a found local minimizer is the global one.

Now let us review some computational basics of nonlinear regression estimation and give necessary definitions. Let y_i be the i th observation and $f(x_i; \beta)$ the regression function where x_i is a fixed vector, $\beta \in \mathbb{R}^m$ is a unknown vector parameter. We assume that f is a twice continuously differentiable function and E is a convex m_j dimensional parameter set. Since $f(x_i; \beta)$ are nonrandom we can simplify the notation letting $f_i(\beta) = f(x_i; \beta)$. Then, the standard nonlinear regression model takes the form:

$$y_i = f_i(\beta) + \varepsilon_i; \quad i = 1; \dots; n \quad (1.1)$$

where it is assumed $E(\varepsilon_i) = 0$; $\varepsilon_1; \dots; \varepsilon_n$ are independent and identically distributed with a positive probability function, and $n \geq m$. The set of points $\{f_1(\beta); \dots; f_n(\beta)\}$ in \mathbb{R}^n is called expectation surface (e.g. Bates and Watts 1988, Seber and Wild 1989). The object of our

study is the Sum of Squares (SS), as a function of parameter given data,

$$S(\theta) = \sum_{i=1}^n (y_i - f_i(\theta))^2 \quad (1.2)$$

Our concern is the global minimum of $S(\theta)$ on E : The point of the minimum is called the Least Squares Estimate (LSE). Geometrically, the LSE corresponds to the nearest point on the expectation surface from the data vector $y = (y_1, \dots, y_n) \in \mathbb{R}^n$. It is assumed that the nonlinear regression is regular (e.g., Gallant 1987, Pazman 1993), i.e. the derivative matrix has full rank, $\text{rank}(G(\theta)) = m$ for all $\theta \in E$; where $G(\theta)$ is the $n \times m$ matrix of first derivatives, $G(\theta) = [\partial f_1 / \partial \theta_1, \dots, \partial f_m / \partial \theta_m]$; $f = (f_1, \dots, f_m)^T$: For linear model it means that the design matrix has full rank. In most cases when the derivative matrix has not full rank the nonlinear regression model is not identified, i.e., for different θ_1 and θ_2 we have $f(\theta_1) = f(\theta_2)$:

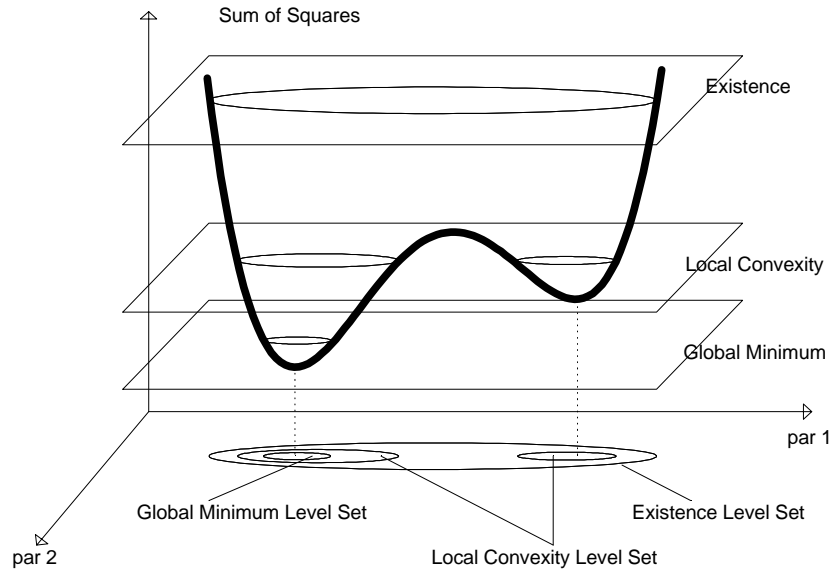


Fig. 1. Three level sets of the Sum of Squares: the Existence LS is compact; the Local Unimodality LS consists of two components corresponding to two local minima with only one local minimizer on each component. The Local Convexity LS is convex, and the SS is a convex function on this set.

In this paper we investigate the properties of the sum of squares in terms of its convexity and possibility to have several local minima, refer Fig.1 for geometrical illustration. The

concept of the Level (or Lebegue) Set (LS),

$$L(S_{\alpha}) = \{\theta \in \mathbb{R}^m: S(\theta) \leq S_{\alpha}\}; \quad (1.3)$$

given level S_{α} is crucial to our analysis. Our general strategy is to express properties of the SS through its levels. For instance, the Existence Level (EL), as the minimum of the SS on the border of the parameter set, was defined in Demidenko (1981, 1989, 1996), Nakamura (1984). If the parameter set \mathbb{E} coincides with the entire space the EL is defined as

$$\bar{S}_E = \lim_{r \rightarrow 1} \inf_{\|\theta\| \leq r} S(\theta);$$

Then the existence criteria works simply: if there is an initial parameter θ_0 such that the value of the SS is less than EL, \bar{S}_E then the LSE exists. It is easy to prove that the level set $L(S(\theta_0))$ is compact, implying that any minimization algorithm produces a sequence of parameters with at least one limit point, Demidenko (1996). In the present paper we introduce another level set, the Local Convexity Level (LCL) set, as the level set where the SS is locally convex, i.e., the Hessian is positive definite. Based on this level we can construct ultimate global criteria. Again, the criterion works easy: if the value of the found local minimum of the SS is less than the calculated threshold value, then this minimum is the global one, see Fig.1.

The outline of the paper is as follows. We start with the warning result that for an intrinsically nonlinear regression with infinite tails the SS quite likely may have at least two local minima, i.e. to be multimodal. In the third section the LCL is introduced and simple global criteria based on general optimization theory are formulated. The following sections contain several examples of the LCL calculation and global criteria construction: polylinear and power regression, M-regression, exponential regression.

2. How likely can one expect multimodality of the Sum of Squares?

The possibility to have several local extrema is studied here for the Sum of Squares in nonlinear regression model. How likely the Sum of Squares (SS) may have two or more local minima? In particular, do there exist classes of nonlinear regressions with unimodal SS for any data y ? As it is shown in this section, for any intrinsic nonlinear regression with infinite tails, the probability that the SS has at least two local minima is positive. Consequently, there are no nonlinear regressions with unimodal SS for any data. Therefore, the problem of determining whether a found local minimizer is indeed the global one, i.e., the LSE, is crucial.

In this section it is assumed that the parameter set coincides with the entire space \mathbb{R}^m : The following result on multimodality will be proved for regressions with infinite tails: $\|\theta\| \rightarrow \infty$

implies $\|f^{(\beta)}\| \leq 1$. We say that a multivariate regression is intrinsically linear if its regression function can be written as $f^{(\beta)} = g_1^{(\beta)}x_1 + \dots + g_m^{(\beta)}x_m + z$ where $x_1, \dots, x_m, z \in \mathbb{R}^n$; the system of vectors $\{x_1, \dots, x_m\}$ has full rank, and $g_1^{(\beta)}, \dots, g_m^{(\beta)}$ are twice continuously differentiable scalar functions such that the map $g: \mathbb{R}^m \rightarrow \mathbb{R}^m$ is injective (Pazman 1993). A characteristic feature of intrinsically linear regression is that its expectation surface lies in a linear m -dimensional manifold. We call regression intrinsically nonlinear if it is not intrinsically linear. Now we formulate the main result on multimodality of the SS for nonlinear regression.

Theorem 2.1. The probability, that the Sum of Squares for any three times differentiable intrinsically nonlinear regression with infinite tails and $n \geq m(m+3)/2$ has at least two local minima, is positive.

Proof of this Theorem is found in the Appendix.

On contrary, as was shown by Pazman (1984), the probability that several local minimizers give the same value for the global minimum of the SS, is zero.

Example. The concept of possible multimodality of the SS is illustrated on the one-parameter quadratic regression model of the form $y_i = \beta x_{i1} + \beta^2 x_{i2} + \varepsilon_i$ where $-1 < \beta < 1$; $\sum x_{i2}^2 > 0$. It is easy to show that this regression has infinite tails. We also assume that $\{x_{i1}\}$ and $\{x_{i2}\}$ are not collinear which implies that this regression is intrinsically nonlinear. Therefore, the above theorem is applicable and the probability to have at least two local minimum of the SS is positive. We aim to determine regions in the outcome space where the SS is not convex, multiextremal, has at least two local minima, etc. For simplicity I consider the regression on the plane with orthogonal covariates, i.e. $n = 2$; and $x_{11} = 1; x_{12} = 0$; and $x_{21} = 0; x_{22} = 1$: Then $f_1(\beta) = \beta$; $f_2(\beta) = \beta^2$ and the expectation curve becomes a parabola in \mathbb{R}^2 (Fig. 2). We have

$$S(\beta) = (y_1 - \beta)^2 + (y_2 - \beta^2)^2; \quad \frac{1}{4} \frac{dS}{d\beta} = \beta^3 + \beta \left(\frac{1}{2} - y_2 \right) - \frac{1}{2} y_1; \quad \frac{1}{4} \frac{d^2S}{d\beta^2} = 3\beta^2 + \left(\frac{1}{2} - y_2 \right);$$

Therefore, the region of the convexity of the SS is below the line $y_2 = 1/2$: This means that if $y_2 < 1/2$ for the data vector (y_1, y_2) then $d^2S/d\beta^2 > 0$ for all $\beta \in \mathbb{R}^1$: If $y_2 > 1/2$; there exists at least one β_0 such that $d^2S/d\beta^2(\beta = \beta_0) = d^2S/d\beta^2 < 0$: Now we find the region of multimodality. Since $dS/d\beta = 0$ is a cubic equation for β we can apply the condition for the existence of three different real roots (e.g., Abramowitz and Stegun p.17, 1972) which leads to the inequality $y_2 > 3^{-3/4} \sqrt{y_1^2 + 1/2}$; $-1 < y_1 < 1$: Therefore, if (y_1, y_2) satisfies the preceding inequality, the SS will have two local minima.

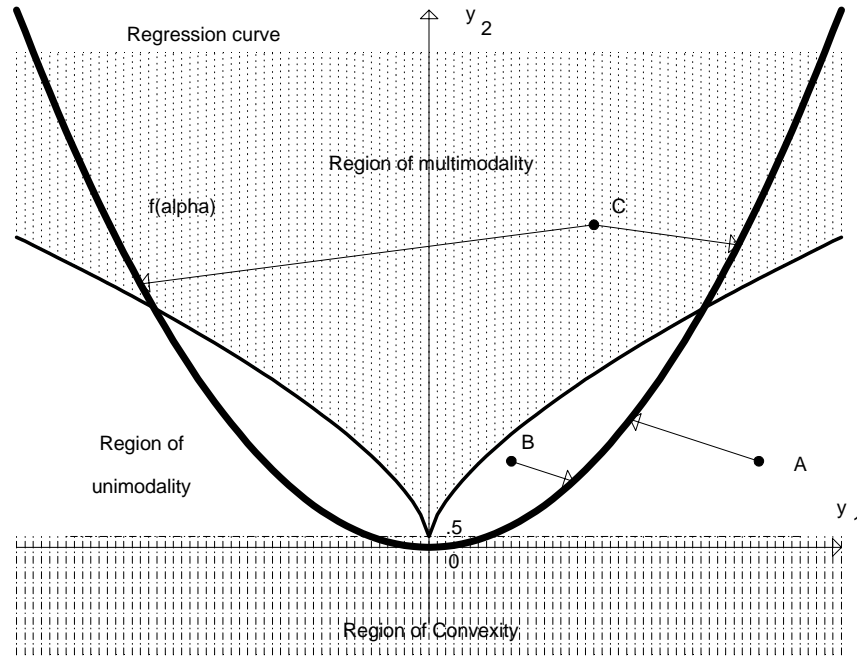


Fig. 2. The regions of the local convexity and multimodality of the SS for the quadratic one-parameter regression. Bold line { regression expectation curve $y_2 = y_1^2$ (parabola). For points A and B there is only one local minimum of the SS, for point C two local minimizers exist, the true least squares estimate corresponds to the right arrow because the distance to the curve is smaller. The region of the local (global) convexity is below the line $y_2 = .5$:

Three regions are shown in Fig. 2. Here, for data vectors A and B the according SS has one local minimum although it is not a convex function (arrows point to the nearest point on the expectation curve, i.e., the LSE). On contrary, for the observation point C the normal equation $dS=d^{\otimes} = 0$ has two solutions, and the LSE corresponds to the right, shortest arrow. For given $\frac{3}{4}^2$ and the true parameter \otimes we can calculate

$$\text{prfSS has two local minima} = \frac{1}{\frac{1}{4}\frac{3}{4}^2} \int_0^1 \int_{\frac{3}{4}^2 y_1^2 + 0.5}^{\frac{3}{4}^2 y_1^2} e^{i \cdot 0.5(y_2 i - \otimes^2) = \frac{3}{4}^2} dy_2 e^{i \cdot 0.5(y_1 i - \otimes) = \frac{3}{4}^2} dy_1.$$

Let us take for example $\frac{3}{4}^2 = 1$: Then, for $\otimes = 0$ the probability of multimodality is quite small, 0:056. However, for moderate $\otimes = 1$ it is 0:15 and for large $\otimes = 3$ the probability that the SS for the quadratic model has two local minima is 0:999:

3. Local convexity

The concept of convexity is fundamental in the theory of optimization, e.g., Rockafellar (1970), Ortega and Rheinboldt (1970), Webster (1994), Chong and Zak (1996). For the

reader's convenience I formulate essential concepts of the optimization theory here, in terms of an arbitrary function.

It is said that a set $Q \subseteq \mathbb{R}^m$ is convex if along with any pair of points it contains the whole segment. Let $F(u)$ be an arbitrary continuous function of vector argument $u \in Q \subseteq \mathbb{R}^m$ where Q is a convex m -dimensional set, or function domain. We say that $F(u)$ is convex on Q if for all $u_1, u_2 \in Q; u_1 \neq u_2; 0 < \lambda < 1$ we have $F(\lambda u_1 + (1 - \lambda)u_2) \leq \lambda F(u_1) + (1 - \lambda)F(u_2)$: It is said that $F(u)$ is strictly convex if the last inequality is strict. We say that u_* is the global minimizer, or the point of the global minimum on Q if $F(u_*) < F(u)$ for all $u \in Q; u \neq u_*$: As follows from this definition, the global minimizer is unique, if it exists. It is said that u_* is a local minimizer, or a point of local minimum if there is a neighborhood V at u_* such that $F(u_*) < F(u)$ for all $u \in Q \cap V; u \neq u_*$: We say that $F(u)$ is unimodal if it has only one local minimum, otherwise we call it multimodal. The central theorem of the optimization theory is that a strictly convex function on a convex set has at most one local minimizer which is the global one. We say that F is locally strictly convex on a set if it is locally strictly convex at each point of the set. It is easy to show that a locally strictly convex function on a convex set is convex. There is a well known sufficient criterion for the local convexity: $F(u)$ is strictly locally convex at u_* if its Hessian matrix at u_* is positive definite (p.d.). These facts are enough to formulate the following

General criterion for the global minimum, I (GCGMI). Let $F(u)$ be a function on $Q \subseteq \mathbb{R}^m$ which is locally strictly convex on a convex subset $L \subseteq Q$. Let us define $\underline{F}_{LC} = \inf_{u \in L} F(u)$: If $u_* \in Q$ is a local minimizer and $F(u_*) < \underline{F}_{LC}$ then u_* is the global minimizer of F on Q :

Proof is short yet important for a better understanding of how this criterion works. It is easy to see that $u_* \in L$ because otherwise $F(u_*) > \underline{F}_{LC}$: Now let u be any point from Q : If $u \in L$ and $u \neq u_*$ we have $F(u) > F(u_*)$ because F is locally convex on a convex set. If $u \notin L$ then $F(u_*) < \underline{F}_{LC} \leq F(u)$: Therefore, u_* is the global minimizer on Q .

We have to mention that this criterion, as all others considered later, is sufficient, not necessary. In other words, its failure does not imply that the minimized function is multimodal.

Now we introduce an important concept, the local convexity level. We know that at the minimum the Hessian is at least nonnegative definite. Clearly, lower the value of the function is more likely the Hessian is positive definite. Therefore, we can define a region where the Hessian is positive definite as a collection of parameter points which give the value of the function less than a certain value (level). It explains the following definition.

Definition 3.1. A value F_{LC} is called Local Convexity (LC) level for function $F(u)$ if

its Hessian is positive definite on the level set $L(F_{LC}) = \{x : F(x) < F_{LC}\}$:

We can speak of the Upper LC level as the maximum of all LC levels. The concept of the LC level provides the following global criterion.

General criterion for the global minimum, II (GCGMII). If the level set $L(F_{LC})$ is convex and a local minimizer u_* belongs to it, then u_* is the global minimizer.

Proof is again easy but worth to present. Since $F(u)$ is a convex function on a convex set $F(S_{LC})$ we obtain $F(u_*) < F(u)$ for all $u \in L(S_{LC})$; $u \notin u_*$: If $u \in L(F_{LC})$; by definition $F(u) < F_{LC} < F(u_*)$: Thus, u_* is the global minimizer.

In the rest of this section we apply this theory to the Sum of Squares (1.2). Notice, the SS is locally strictly convex at θ if matrix

$$\frac{1}{2} \frac{\partial^2 S}{\partial \theta^2} = G^T(\theta) G(\theta) + \sum_{i=1}^n (y_i - f_i(\theta)) H_i(\theta) \quad (3.1)$$

is p.d. where $G(\theta) = \partial f / \partial \theta$ is the $n \times m$ matrix of the first derivatives and $H_i(\theta) = \partial^2 f_i / \partial \theta^2$ is the $m \times m$ symmetric matrix of the second derivatives of f_i , $i = 1, \dots, n$:

We want to express the positive definiteness of (3.1) in terms of the SS-value. For this purpose let us multiply (3.1) by p^T and p where p is an arbitrary nonzero $m \times 1$ vector (T denotes vector/matrix transpose, hereafter). Applying the Cauchy inequality to the second term, $p^T \left(\sum_{i=1}^n (y_i - f_i) \right) p$ we obtain

$$p^T \frac{1}{2} \frac{\partial^2 S}{\partial \theta^2} p \geq p^T G^T(\theta) G(\theta) p + \frac{p^T \left(\sum_{i=1}^n (y_i - f_i) \right)^2}{\sum_{i=1}^n (p^T H_i(\theta) p)^2} \quad (3.2)$$

Define the radius of the full curvature as a function of the parameter vector θ and the $m \times 1$ direction vector p

$$R_F(\theta; p) = \frac{p^T G^T(\theta) G(\theta) p}{\sum_{i=1}^n (p^T H_i(\theta) p)^2} \quad (3.3)$$

following Bates and Watts (1988). If the derivative matrix G has full rank the numerator is positive, thus we can put $R_F = +1$ if the denominator is zero. Thus, based on the inequality (3.2) we can determine the LC level in nonlinear regression as follows.

Local Convexity Level for SS. The Upper Local Convexity Level (ULCL) for the SS of nonlinear regression model (1.1) is equal to the minimum of the squared radius of the full curvature of the expectation surface, or

$$\bar{S}_{LC} = \min_{\theta} \min_p R_F^2(\theta; p) \quad (3.4)$$

The Local Convexity Level (LCL) is any number less or equal \bar{S}_{LC} ; and will be denoted S_{LC} :

Shortly, the ULCL is the minimum of the squared radius of the full curvature. The word "upper" in the last definition can be justified by the following comment. If for some θ_0 and p_0 we have $R_F^2(\theta_0; p_0) > \bar{S}_{LC}$ then there are observations $y_1; \dots; y_n$ with the $SS = \sum_{i=1}^n (y_i - f_i(\theta_0))^2 = R_F^2(\theta_0; p_0)$ and $p_0^T (\frac{\partial^2 S(\theta_0)}{\partial \theta^2}) p_0 < 0$; i.e. the Hessian is not p.d. at $\theta = \theta_0$. It follows from the fact that the Cauchy inequality in (3.2) turns into equality if vectors $(y_1 - f_1; \dots; y_n - f_n)$ and $(p_0^T H_1 p_0; \dots; p_0^T H_n p_0)$ are collinear.

Now we will consider the problem of calculating $\min_p R_F^2(\theta; p)$ given θ (therefore θ is omitted). As follows from (3.3) this optimization problem is equivalent to

$$\lambda' = \lambda'(p; G, H_i; i = 1; \dots; n) = \frac{(\sum_{i=1}^n p_i^T G_i G_i p_i)^2}{(\sum_{i=1}^n p_i^T H_i p_i)^2} \rightarrow \min_p \quad (3.5)$$

We show that this problem can be reduced to a series of eigenvalue problems; called repeated eigenvalue problem. Indeed, let T be the Cholesky decomposition factor of $G^T G$; i.e., T is a triangle nonsingular matrix such that $T^T T = G^T G$. Denote $q = Tp$; then $p^T H_i p = q^T Q_i q$ where $Q_i = (T^T)^{-1} H_i T^{-1}$ and λ'_{\min} is equal to reciprocal of

$$\lambda'_{\max} = \max_{q^T q = 1} \sum_{i=1}^n (q^T Q_i q)^2 \quad (3.6)$$

To solve (3.6) we denote $K(q) = \sum_{i=1}^n Q_i q q^T Q_i$; then $\sum_{i=1}^n (q^T Q_i q)^2 = q^T K(q) q$. Let q_0 be any initial vector of unit length. We solve the eigenvalue problem $q^T K(q_0) q \rightarrow \max$ and obtain the solution q_1 . Then matrix K is recalculated, i.e., we solve $q^T K(q_1) q$ and repeat this process until convergence. A good starting value for q would be the eigenvector correspondent to the maximum eigenvalue of matrix $\sum_{i=1}^n Q_i^2$. The sequence $\{q_s\}$ generated by this process has at least one limit point because it belongs to a compact set, the sphere of the unit radius. Moreover, the sequence $\{q_s^T K(q_s) q_s\}$ is increasing because

$$q_s^T K(q_s) q_s \leq q_{s+1}^T K(q_s) q_{s+1} = q_s^T K(q_{s+1}) q_s \leq q_{s+1}^T K(q_{s+1}) q_{s+1}; \quad s = 0; 1; \dots$$

The minimum of (3.5) is denoted as $\lambda'_{\min} = 1/\lambda'_{\max}$. Since $q^T q = 1$ it implies $(q^T Q_i q)^2 \leq p^T H_i^2 p$; and $\lambda'_{\min} \geq \lambda'_{\min}(G^T G (\sum_{i=1}^n H_i^2)^{-1} G^T G)$. Obviously, for linear regression $\sum_{i=1}^n H_i^2(\theta) = 0$ and $\bar{S}_{LC} = 1$:

It is important to understand that it is not required calculating the value \bar{S}_{LC} precisely using some optimization algorithm (otherwise all above does not make sense because it might be even more difficult task than the SS minimization itself). Instead, we have to conduct an analytical investigation to find a lower bound for \bar{S}_{LC} ; the examples can be found in the next section.

4. Examples. General

In the following sections the calculation of the local convexity level and construction of global criteria are illustrated by four examples of nonlinear regression model. In the first example the exact value \bar{S}_{LC} is calculated. In the second example the LCL is derived using some specific inequalities. In the third example we demonstrate how the concept of the LCL can be generalized to functions different from the SS, and a global criterion is formulated based on the general criterion for the global minimum formulated in the previous section. At last, the exponential models are considered: the formula for the LCL is provided and a global criterion is suggested. Almost all criteria use the repeated eigenvalue problem at the final stage, considered in the preceding section.

As the reader can realize, every type of nonlinear model requires some creativity in finding a lower bound for the Hessian of the SS and constructing a global criterion. In particular, at some point one has to apply certain inequalities which take into account specific structure of the considered nonlinear regression model.

5. Polylinear regression

This type of nonlinear regression is defined by the function

$$f(\bar{\cdot}; \frac{1}{2}) = X\bar{\cdot} + \frac{1}{2}u_i \frac{1}{2}Z\bar{\cdot}; \quad \bar{\cdot} \in 2R^m; \frac{1}{2} \in 2R^1 \quad (5.1)$$

where X and Z are fixed $n \in m$ design matrices, matrix Z has full rank, $m < n$; and u is a $n \in 1$ vector. This kind of model emerges in time series regression analysis where residuals follow the first-order autoregression process (e.g., Judge et al. 1982). More precisely, if $y_t = \beta_0 x_t + \epsilon_t$ is the original linear model and $\epsilon_t = \frac{1}{2}\epsilon_{t-1} + \eta_t$ where η_t is the uncorrelated error term with zero mean and constant variance, we can apply first-order differences with unknown coefficient $\frac{1}{2}$ to remove the existing correlation between ϵ_t and ϵ_{t-1} . This gives us a nonlinear regression $y_{t-1} - \frac{1}{2}y_{t-2} = \beta_0 x_{t-1} - \frac{1}{2}\beta_0 x_{t-2} + \eta_{t-1}$ which has a particular form (5.1).

For regression (5.1) we have $G = \partial f / \partial (\bar{\cdot}; \frac{1}{2}) = [X \quad \frac{1}{2}Z; u_i \quad Z^T]$; the $n \in (m + 1)$ matrix, and

$$\frac{1}{2}p^0 \partial^2 S = \partial^2 G(\bar{\cdot}) p = p^0 G^0(\bar{\cdot}) G(\bar{\cdot}) p_i \quad \times \quad (y_i - f_i(\bar{\cdot})) p^0 H_i p$$

where

$$H_i = \begin{pmatrix} 0 & z_i \\ z_i^0 & 0 \end{pmatrix}; \quad (5.2)$$

and z_i is the i th row of matrix Z ; $\bar{\cdot} = (\beta_0; \frac{1}{2})^0$. We notice that the denominator of (3.3) for this model does not depend on parameter, and consequently in order to find the minimum over $\bar{\cdot}$ it suffices to consider only the numerator. Let $p = (p_1^0; p_2^0)^0$ be fixed, then we have

$$\min_{\bar{\cdot}} p^0 G^0(\bar{\cdot}) G(\bar{\cdot}) p = \min_{\bar{\cdot}} \|G p\|^2 = \min_{\bar{\cdot}} \| (X - \frac{1}{2}Z) p_1 + (u_i - Z^T) p_2 \|^2$$

$$= \min_{\tilde{z}} \| (Xp_1 + p_2u) - Z(p_2 + \frac{1}{2}p_1) \|_2^2;$$

Since \tilde{z} and $\frac{1}{2}$ are unrestricted vector $\tilde{z} = \tilde{z}p_2 + p_1\frac{1}{2}$ covers the entire space R^m ; so that we may write

$$\min_{\tilde{z}} \| (Xp_1 + p_2u) - Z(p_2 + \frac{1}{2}p_1) \|_2^2 = \min_{\tilde{z}} \| Wp - Z\tilde{z} \|_2^2;$$

where $W = [X;u]$: However, for fixed Wp the minimum over \tilde{z} can be found via least squares:

$$\min_p p^0 G^{(0)} G^{(0)} p = p^0 W^0 (I - Z(Z^0 Z)^{-1} Z^0) W p;$$

Hence, finally

$$\bar{S}_{LC} = \min_p \frac{(p^0 W^0 (I - Z(Z^0 Z)^{-1} Z^0) W p)^2}{(p^0 H_1 p)^2}. \quad (5.3)$$

Finally, we notice that the minimum in (5.3) may be found by solving the repeated eigenvalue problem considered previously. Therefore, for the polylinear regression (5.1), the Hessian is positive definite if sum of squares less than \bar{S}_{LC} : Notice that $\bar{S}_{LC} > 0$ if and only if any linear combination of vector-columns of matrix $[X;u]$ cannot be expressed as a linear combination of vector-columns of matrix Z ; e.g., matrix $[X;u;Z]$ has full rank $2m + 1 + n$:

6. Power regression

We define power regression as a regression with function $f_i^{(p)} = (x_i^0)^p$; where $p \neq 0$ is given, and $x_1; \dots; x_n$ are linearly independent. $p = 1$ gives standard linear regression model, $p = -1$ corresponds to hyperbolic, or Michaelis-Menten model (Bates and Watts 1988): Power regressions with values $p = \pm 2$ are also popular (e.g., Ezekiel and Fox 1959). The parameter set is defined as

$$\begin{aligned} \mathcal{E} &= \{f^{(p)} \in R^m : x_1^0 \geq 0; \dots; x_n^0 \geq 0\} \quad \text{when } p > 0; \\ \mathcal{E} &= \{f^{(p)} \in R^m : x_1^0 > 0; \dots; x_n^0 > 0\} \quad \text{when } p < 0; \end{aligned}$$

As follows from its definition the parameter set is convex; also it will be assumed that it is not empty and has dimension m . For the sake of simplicity the following development assumes $y_i > 0; i = 1; \dots; n$; a similar analysis can be accomplished in case when some of y_i are negative.

For the power regression

$$\frac{1}{2p} \frac{\partial^2 S}{\partial \theta^2} = X^0 D^{(p)} X \quad (6.1)$$

where the $n \times m$ matrix X consists of vectors $f x_i g$; and $D^{(p)}$ is the $n \times n$ diagonal matrix with diagonal elements

$$D_{ii}^{(p)} = (2p - 1) d_i^{2(p-1)} - (p - 1) y_i d_i^{p-2} = p(y_i - e_i)^{2p-2} - (p - 1) e_i (y_i - e_i)^{2p-3} \quad (6.2)$$

where $d_i = \sum_{j=1}^p x_j$ and $e_i = y_i - d_i^p$ is the i th residual.

We aim to find local convexity level, i.e. level sets of SS where the Hessian, (6.1) is positive definite (diagonal elements of matrix D are positive). The analysis of positive definiteness of the Hessian depends on the value of p : We consider some special cases:

1. $1 \leq p \leq 1$: Since $(p-1)y_i \leq 0$ and $(2p-1)d_i^p \geq 0$ the SS is a strictly convex function on \mathbb{R} . Consequently, the SS is unimodal, and any local minimizer is the global one.

2. $p > 2$: The main idea is to find lower bound for (6.2) in terms of a linear function of the residual, e_i : For the first term of the right hand side of (6.2), applying the inequality $(1-x)^v \geq 1-vx$; $v \geq 1$; $x \leq 1$; we obtain

$$(y_i - e_i)^{2i \frac{2}{p}} = y_i^{2i \frac{2}{p}} \left(1 - \frac{e_i}{y_i}\right)^{2i \frac{2}{p}} \geq y_i^{2i \frac{2}{p}} \left(1 - \frac{2}{p} \frac{e_i}{y_i}\right) = y_i^{2i \frac{2}{p}} \left(1 - \frac{2(p-1)}{p} \frac{e_i}{y_i}\right).$$

For the second term of (6.2) we apply another elementary inequality:

$$1 - x(y_i - x)^v \geq \frac{1}{v+1} \left(1 - \frac{1}{v+1}\right)^v y_i^{v+1}; \quad 0 \leq x \leq y_i; y_i > 0; v > 0;$$

which implies

$$1 - (p-1)e_i(y_i - e_i)^{1i \frac{2}{p}} \geq \frac{p}{2} \left(1 - \frac{p}{2(p-1)}\right)^{1i \frac{2}{p}} y_i^{1i \frac{2}{p}}.$$

Combining the two inequalities we come to a lower bound for i th diagonal element of matrix D , and therefore

$$\frac{1}{2p} \frac{\partial^2 S}{\partial x_i^2} \geq T \sum_{j=1}^p y_i^{2i \frac{2}{p}} x_j x_i^0 + \frac{2}{p} \sum_{j=1}^p y_i^{1i \frac{2}{p}} e_j x_i x_i^0;$$

where

$$T = T(p) = \frac{1}{4} \frac{p-2}{p-1} \left(\frac{p-2}{p-1}\right)^{\frac{1}{p-2}}.$$

Hence, the LCL for the power regression in the case $p > 2$ is

$$S_{LC} = \frac{T^2 p^2}{4(p-1)^2} \cdot \min_{\mathbf{x}} \sum_{i=1}^n y_i^{2i \frac{2}{p}} x_i x_i^0 + \sum_{i=1}^n y_i^{1i \frac{2}{p}} x_i x_i^0; \quad i = 1, \dots, n.$$

Again, λ_{\min} is calculated by the repeated eigenvalue problem considered in the previous section.

3. $p < 0$: Let us introduce the following function of x :

$$L(x) = (y - x)^{1-2p} (1 - p y - (1 - 2p)x); \quad x < y$$

where $y > 0$ is fixed. It is possible to show that for this function

$$L(x) > L(0) + x \left(\frac{dL}{dx} \right)_{x=0} = \sum_{i=1}^n p y_i^{2i-2-p} \sum_{j=1}^n (1 - p) y_i^{1i-2-p} x_j \quad (6.3)$$

for all x . Hence, applying the inequality (6.3) to $\sum_{i=1}^n D_{ii}(\mathbb{R})$ we obtain

$$\sum_{i=1}^n \frac{1}{2p} \frac{\partial^2 S}{\partial x_i^2} \cdot \sum_{j=1}^n p y_i^{2i-2-p} x_j x_i^0 \sum_{j=1}^n (1 - p) y_i^{1i-2-p} e_j x_j x_i^0.$$

Therefore, the local convexity level is

$$S_{LC} = \frac{p^2}{9(1-p)^2} \cdot \min_{i=1, \dots, n} \sum_{j=1}^n y_i^{2i-2-p} x_j x_i^0 \sum_{j=1}^n y_i^{1i-2-p} x_j x_i^0; \quad (6.4)$$

Example. We consider an example from Bates and Watts (1988, p. 114) with Michaelis-Menten model (modified) $f(z; \mu) = \mu_1 z / (\mu_2 + z + \mu_3 z^2)$ which can be parameterized as power model with $p = 1$ as follows, $\mathbb{R}_1 = 1/\mu_1$; $\mathbb{R}_2 = \mu_2/\mu_1$; $\mathbb{R}_3 = \mu_3/\mu_1$ and $x_1 = 1$; $x_2 = 1/z$; $x_3 = z$: The LSE is $(1.13 \pm 10^{-5}; 2.08 \pm 10^{-3}; 1.81 \pm 10^{-7})$ with the minimum SS 1.38 ± 10^6 . The S_{LC} value for this problem, computed by formula (6.4), is 4.81 ± 10^6 :

7. Robust M-regression

It is well known that the LSE is not robust to outliers (Huber 1981). The method of robust M-regression is based on the minimization of sum of a function of residuals, different from quadratic,

$$W(\mathbb{R}) = \sum_{i=1}^n \frac{1}{2} (y_i - \mathbb{R}^0 x_i); \quad (7.1)$$

where $\frac{1}{2} = \frac{1}{2}(e)$ is a loss function with the properties: a) $\frac{1}{2}(0) = 0$; b) $\frac{1}{2}(e) = \frac{1}{2}(j e)$; c) $\frac{1}{2}(e) \geq 0$; $e \geq 0$: As was pointed out by Huber (1981), when $\frac{1}{2}$ is not convex, the minimization of $W(\mathbb{R})$ may cause problems because it may have several local minima: "Unless we are careful we may even get trapped in a local minimum of $W(\mathbb{R})$: The situation gets particularly acute in multiparameter regression". The reader is referred to Fig. 3 where a piece of surface (7.1) is displayed for a real life data (see Example below).

A number of authors addressed the question of optimal choice of the loss function and relevant comparisons have been made (e.g. Andrews et al. 1972). Particularly, it was suggested to choose $\frac{1}{2}$ close to quadratic function near zero with asymptotes at $j \rightarrow \pm \infty$: Probably, the simplest loss function of this kind is

$$\frac{1}{2}(e) = \frac{e^2}{c + e^2}; \quad (7.2)$$

where $c > 0$ is a given parameter which affects the sensitivity of the M-regression to outliers. Montgomery and Peck (1992), among others, consider other loss functions and discuss iterative algorithms for W minimization.

We start with finding regions where $W(\beta)$ is convex, then we construct sufficient global criteria. The Hessian matrix for function (7.1) is

$$\frac{\partial^2 W}{\partial \beta^2} = \sum_{i=1}^n \frac{1}{2} \frac{\partial^2}{\partial \beta^2} (y_i - \beta^T x_i)^2 \quad (7.3)$$

where in our case

$$\frac{1}{2} \frac{\partial^2}{\partial \beta^2} (e) = \frac{2c(c - 3e^2)}{(c + e^2)^3} \quad (7.4)$$

A simple global criterion for the M-regression can be constructed based on the GCCMI formulated in the previous section. Indeed, from (7.4) we notice that $\frac{\partial^2 W}{\partial \beta^2}$ is p.d. if $c - 3(y_i - \beta^T x_i)^2 > 0$ for all $i = 1, \dots, n$. Thus, on the convex set $L = \{\beta : |y_i - \beta^T x_i| < \sqrt{c/3}; i = 1, \dots, n\}$ the objective function W is convex. Now, to apply the GCCMI we need to find a lower bound for $W(\beta)$ on L . In fact, if $\beta \in L$ then for some j we have $|y_j - \beta^T x_j| < \sqrt{c/3}$ and consequently $W(\beta) \geq \frac{1}{2} (c - 3 \cdot \frac{c}{3}) = 1/4$. Therefore we infer that if β_* is a local minimizer with $W(\beta_*) < 1/4$ then β_* is the global minimizer of (7.1), i.e., the M-estimate.

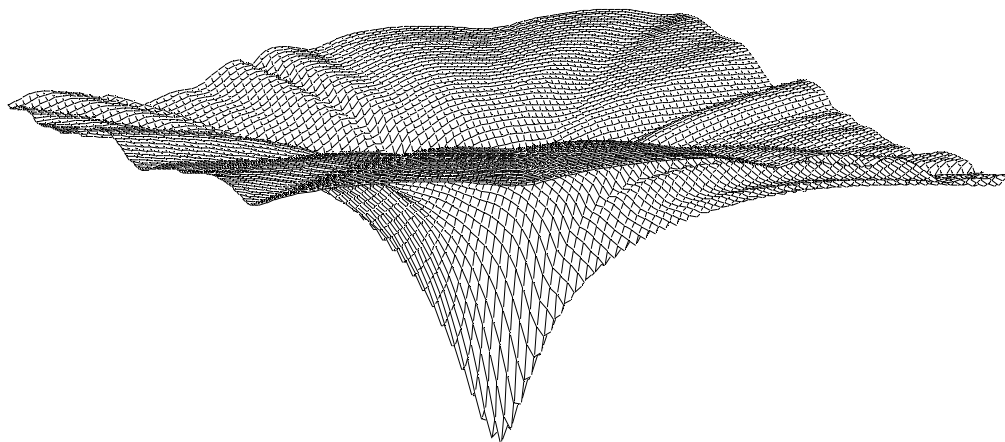


Fig. 3. A "mashroom" surface $W(\beta_1; \beta_2)$ for the M-regression in the neighborhood of a minimum (global?) for vending machine example.

Now we suggest a more powerful criterion. We start from seeking a lower linear bound for $\frac{1}{2}^{00}$. It is easy to show that for function (7.2)

$$\frac{1}{2}^{00}(0) = 2=c; \quad \frac{1}{2}^{00}(\overline{S^P_{c=3}}) = 0; \quad \frac{1}{2}^{00}(e) \neq 0 \text{ when } |e| \neq 1 :$$

The minimum of $\frac{1}{2}^{00}(e)$ is attained at $\overline{S^P_{c=3}}$ and equal $|e| \leq c$: We can find the minimal positive constant D such that $\frac{1}{2}^{00}(e) \geq \frac{1}{2}^{00}(0) - D |e|$ for all $e \in \mathbb{R}^1$: Geometrically, for positive e ; it corresponds to a line which goes through $(0; 2=c)$ and touches $\frac{1}{2}^{00}(e)$ at some point. Clearly, this point is determined by the equation $|\frac{1}{2}^{00}(e)| = (2=c - \frac{1}{2}^{00}(e)) = e$ which is equivalent to the cubic equation $u^3 + 4u^2 + 21u - 6 = 0$ where $u = e^2=c$: This equation has a unique real root u_* (approximately $u_* = 0.2708$): Hence $D = |\frac{1}{2}^{00}(\overline{P_{u_*c}})| = 3.492c^{1/3}$; and the lower bound for $\frac{1}{2}^{00}$ can be written as $\frac{1}{2}^{00}(e) \geq 2=c - 3.492c^{1/3} |e|$: Substituting it into (7.3) we come to a lower bound for the Hessian

$$\frac{\partial^2 W}{\partial \theta^2} \geq 3.492c^{1/3} - 0.5727 \overline{P_{c=3}}^T X X^T y - \theta^0 X_j y_j - \theta^0 X_j X_j X_j^T :$$

As follows from the previous section on the ellipsoid

$$E = \{x \in \mathbb{R}^m : (y_i - \theta^0 x_i)^2 \leq q\} \text{ where } q = (0.5727)^2 c^{1/3} \quad (7.5)$$

and $\theta^0_{\min} = \theta^0_{\min}(\overline{P_{x_i x_i^0 - j x_i x_i^0; i = 1; \dots; n}})$; the function $W(\theta)$ is convex.

Now, in order to apply GCGMI, we need to find a lower bound for W outside ellipsoid (7.5). For this purpose we make use of the following simple optimization problem

$$\overline{P_{(v_i - \theta^0 x_i)^2 \leq q}} \min_j v_j - \theta^0 x_j = \max(0; \zeta_j) \quad (7.6)$$

where $\zeta_j = \frac{q}{(q - S_{\min})x_j(X^0 X)^{-1}x_j} - j v_j - a_{OLS}^0 x_j$ and a_{OLS} is the OLS-estimate in regression $f v_j$ on $f x_j$; $q \geq S_{\min} = \min_j (v_j - a_{OLS}^0 x_j)^2$: Then if $\theta^0 \in E$; as follows from (7.6). letting $v_i = y_i$

$$W(\theta) = \sum_i \frac{e_i^2}{c + e_i^2} \geq \sum_i \frac{\max^2(0; \zeta_i)}{c + \max^2(0; \zeta_i)} = W_g \quad (7.7)$$

Based on this inequality the following sufficient global criterion can be formulated.

Criterion for the global minimum of $W(\theta)$. If for a local minimizer θ_* we have $W(\theta_*) < W_g$ then θ_* is the global minimizer, i.e., the true M-estimate.

Example. The above criterion is illustrated by a problem analyzing the vending machine service routes (Montgomery and Peck 1982, Example 9.4). The data consists of 25 observations where the dependent variable y is the delivery time (minutes), x_1 is the number of cases of product stocked and x_2 is the distance walked by the route driver (feet). To reduce

this problem to a two-parameter regression we eliminate the intercept term centering data: $y_i = y_i - \bar{y}$; $x_{i1} = x_{i1} - \bar{x}_1$; and $x_{i2} = x_{i2} - \bar{x}_2$. Then, the model is $y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i$. We apply the loss function (7.2) with $c = 400$. The objective function, $W(\beta_1; \beta_2)$ is quite "wavy" (see Fig. 2), and possibly has several local minima. The minimization of (7.1) by iterative weighted least squares (Huber 1981) leads to a local minimum $\beta_{\alpha 1} = 1.594$; $\beta_{\alpha 2} = 0.014$ with $W(\beta_{\alpha 1}; \beta_{\alpha 2}) = 0.544$. Is this the global minimum? To answer the question we apply the above criterion based on calculation (7.7). For our data $q = 795.2$ and $W_g = 1.13 > 0.544 = W(\beta_{\alpha 1}; \beta_{\alpha 2})$. Therefore, the answer is positive { the found local minimum is the global one.

8. Exponential regression models

This type of nonlinear regression is popular, for instance, in econometrics (Kobb-Douglas production function, Judge et al. 1982) and chemistry (monomolecular growth curves, Seber Wild 1989). The regression function of the exponential model is defined as $f_i(\beta) = \exp(\beta x_i)$ where $\beta \in \mathbb{R}^m$ and the system of vectors $f x_1; \dots; x_n g$ has rank m . The Hessian of the SS for this regression is given by

$$\frac{\partial^2 S}{\partial \beta^2} = 2 \sum_{i=1}^n (2e^{\beta x_i} - y_i) x_i x_i^T \quad (8.1)$$

We assume here $y_i > 0$ for all $i = 1; \dots; n$. As follows from (8.1), on the polyhedra

$$P = \{ \beta \in \mathbb{R}^m : \beta x_i > \ln y_i - \ln 2; \quad i = 1; \dots; n \} \quad (8.2)$$

the SS is convex. Now we apply GCGMI to construct a simple global criterion. For this we have to find a lower bound for the SS outside P . In fact, if $\beta \notin P$ then for some j we have $y_j - e^{\beta x_j} \geq 0.5 y_j > 0$ which implies

$$S(\beta) = \sum_{i=1}^n (y_i - e^{\beta x_i})^2 \geq (y_j - e^{\beta x_j})^2 \geq \frac{1}{4} y_j^2 \geq \frac{1}{4} \min_i y_i^2 g$$

Therefore, we can construct the following quick global criterion: if β_{α} is a local minimizer of the SS and $S(\beta_{\alpha}) < 0.25 \min y_i^2 g$; then β_{α} is the LSE.

Now we find a more precise LCL. Let us denote $e_i = y_i - e^{\beta x_i}$; then $2e^{\beta x_i} - y_i = y_i^2 - 3e_i y_i + 2e_i^2$; and

$$\frac{\partial^2 S}{\partial \beta^2} = \sum_{i=1}^n y_i^2 x_i x_i^T - \sum_{i=1}^n e_i y_i x_i x_i^T - 2 \sum_{i=1}^n e_i (y_i - e_i) x_i x_i^T \quad (8.3)$$

For the third term in (8.3) we apply an elementary inequality $e_i (y_i - e_i) \leq y_i^2/4$ which leads us to the following lower bound for the half Hessian

$$\frac{\partial^2 S}{\partial \beta^2} \geq \frac{1}{2} \sum_{i=1}^n y_i^2 x_i x_i^T - \sum_{i=1}^n e_i y_i x_i x_i^T$$

Therefore the LCL can be defined as

$$S_{LC} = \frac{1}{2} \min_i \sum_{j=1}^n y_i^2 x_i x_i^0 \quad i = 1; \dots; n : \quad (8.4)$$

Hence, for any parameter with the SS less than S_{LC} defined by (8.4) the Hessian (8.1) is positive definite.

Now we construct a global criterion. Since function e^v is convex $e^v \geq e^{v_0}(v - v_0) + e^{v_0}$ for any v and v_0 : Applying this inequality to (8.1) with $v = \beta^0 x_i$ and $v_0 = \ln y_i$ we obtain

$$\frac{\partial^2 S}{2 \partial \beta^2} \geq \sum_{i=1}^n y_i x_i x_i^0 - \sum_{i=1}^n (\ln y_i - \beta^0 x_i) y_i x_i x_i^0$$

Therefore on the ellipsoid

$$E = \{\beta \in \mathbb{R}^m : \sum_{i=1}^n (\ln y_i - \beta^0 x_i) < q\} \text{ where } q = \frac{1}{4} \min_i \left(\sum_{j=1}^n y_j x_j x_j^0 \right) \quad i = 1; \dots; n$$

the SS is convex. In order to apply GCGMI we need to find a lower bound for $S(\beta)$ when β lies outside E : For this purpose we again use the inequality (7.6), letting $v_i = \ln y_i$: Since $\sum_{j=1}^n \ln y_j - \beta^0 x_j \geq \sum_{j=1}^n \ln y_j - \beta^0 x_j$ for $\beta_j \geq 0$ implies $(y_j - e^{\beta^0 x_j})^2 \geq y_j^2 (1 - e^{-\beta_j})^2$ we obtain

$$S(\beta) \geq \sum_{\beta_j > 0} y_j^2 (1 - e^{-\beta_j})^2 = S_g; \quad \beta \notin E : \quad (8.5)$$

Hence, if β_* is a local minimizer in exponential model with the SS less than S_g then β_* is the global minimizer.

Example. We borrow data from Bates and Watts (1989, p. 281) on tetracycline concentration (1g/ml) versus time (hr). Thus, if y denotes tetracycline concentration and t time we suggest the following model

$$y_i = \exp(\beta_1 + \beta_2 t_i^2 + \beta_3 t_i) + \epsilon_i; \quad i = 1; \dots; n = 9:$$

The nonlinear estimation starting from the OLS-estimates in the linear regression $\ln y_i$ on $\beta_1 + \beta_2 t_i^2 + \beta_3 t_i$ comes to a local minimizer $\hat{\beta} = (1.33; 0.163; 1.56)$ with $SS_{\min} = 0.0142$. Is $\hat{\beta}$ the true LSE, i.e., did we find the global minimum of the SS? The χ^2 value is 0.355 and $S_g = 0.476$. Since $SS_{\min} < S_g$ we infer that $\hat{\beta}$ is the true least-squares estimate.

9. Conclusions

The existing algorithms for sum of squares minimization, such as Gauss-Newton or Marquardt, can find only a local minimum. Unlike linear regression the SS for nonlinear regression may have several local minima. Moreover, if a nonlinear model cannot be reparameterized to linear, the probability to have at least two local minima is positive. Therefore, almost

in all cases of nonlinear estimation there is a doubt that the model was estimated correctly, i.e., the found local minimizer is in fact the Least Squares Estimate. The current paper aims to provide some guidelines in constructing criteria which verify that the local minimum is the global one. Notice, these criteria are sufficient, not necessary. In other words if criteria fail it does not mean the SS has two local minima and the found minimum is not the global one. Our strategy is to express the properties of the SS in terms of descended levels:

1. Existence Level. This level was introduced in Demidenko (1981) and recently illustrated through different types of nonlinear models in Demidenko (1996). Hence, if \bar{S}_E is an EL and \mathbb{R}_0 is an initial parameter vector such that $S(\mathbb{R}_0) < \bar{S}_E$ then the LSE exists.
2. Local Convexity Level. This level is equal to the minimum of squared radius of the full curvature of the expectation surface. Hence, if S_{LC} is a level of LC and $S(\mathbb{R}) < S_{LC}$ then it is guaranteed that the Hessian of SS is positive definite at \mathbb{R} .
3. Global Minimum Level. If S_G is such a level and \mathbb{R}_* is a local minimizer such that $S(\mathbb{R}_*) < S_G$ then \mathbb{R}_* is the global minimizer, i.e., the LSE.

One need to understand that except rare cases we are not able to find the (exact) upper local convexity level because it would lead to another minimization problem, perhaps even more severe than the SS minimization itself. However, we suffice to find a lower bound. Much creativity should apply to find a satisfactory lower bound and examples considered in this paper hopefully illustrate this. However, much work should be done to develop these techniques for other, more complicated nonlinear regression models. Global criteria must be a part of software on nonlinear estimation.

ACKNOWLEDGMENT

This work was supported by National Cancer Institute grants CA52192 and CA61108.

10. Appendix. Proofs

The following auxiliary result will be used, Demidenko (p. 85, 1989).

Lemma 10.1. A multivariate regression is intrinsically nonlinear if and only if there exists \mathbb{R}_0 such that the system of $m(m+3)/2$ vectors from R^n

$$\left(\frac{\partial f(\mathbb{R}_0)}{\partial \mathbb{R}_j}; \frac{\partial^2 f(\mathbb{R}_0)}{\partial \mathbb{R}_j \partial \mathbb{R}_k}; j = 1, \dots, m; k = j \right) \quad (10.1)$$

has full rank.

Also, the following fundamental result of optimization theory will be used.

Lemma 10.2. Let $F(u)$ be a twice differentiable function of $u \in \mathbb{R}^m$ such that $\|u\| \leq 1$ implies $F(u) \leq 1$: If there exists u_0 such that $\nabla F(u_0) = 0$ and $\nabla^2 F(u_0) = \nabla^2 u^2$ is a negative definite matrix then $F(u)$ has at least two local minima on \mathbb{R}^m .

As pointed out by Mäkeläinen, Schmidt and Styan (1981) this result is a straightforward consequence of Morse theory, however they provide an elementary proof for a similar result. A more general results is proved in Demidenko (1989).

Proof of Theorem 2.1. The idea of the proof is to use the implicit function theorem by showing that the Jacobian matrix for an introduced mapping has full rank. It provides that, under conditions of the Theorem, observations y with multimodal SS fall an open set in \mathbb{R}^n : Clearly, that implies the probability to have at least two local minima of SS is positive, provided that the density of the error distribution is positive.

Now we realize this plan in some detail. Let us denote the $n \times m_\alpha$ matrix with the columns (10.1) as $M(\theta_0)$ where $m_\alpha = (m(m+3))/2$: As follows from Lemma 10.1 it has full rank. Without loss of generality one can assume that the upper $m_\alpha \times m_\alpha$ part of this matrix is nonsingular. Since the elements of matrix M are continuous functions of θ the upper matrix has full rank in some neighborhood $V = V(\theta_0)$: Further, let us part matrix $M(\theta)$ as follows

$$M(\theta) = M = \begin{pmatrix} G_1 & U_1 \\ G_2 & U_2 \end{pmatrix};$$

where G_1 is the $m_\alpha \times m$ matrix. Then $G = [G_1^0, G_2^0]^0$ is the matrix of first derivatives and $U = [U_1^0, U_2^0]^0$ is the matrix of second derivatives (argument θ is omitted). As was noticed above, $\det[G_1; U_1] \neq 0$: Further, let us consider the set of all symmetric positive definite $m \times m$ matrices M_+ : Clearly, any matrix D from M_+ is determined by $m(m+1)/2$ numbers. As follows from Lemma 10.2, the SS has at least two local minima if for some θ we have $\nabla S(\theta) = \nabla \theta = 0$ and $\nabla^2 S(\theta) = \nabla^2 \theta^2 = \sum_i D_i$. We aim to show that when θ takes values from V and D is p.d. the solutions of the two last equations, fall a certain n_j dimensional region in \mathbb{R}^n . Since the Hessian is a symmetric matrix we can apply the vech function (e.g., Magnus and Neudecker 1988) to $\nabla^2 S(\theta) = \nabla^2 \theta^2 = \sum_i D_i$ and therefore come to an equivalent vector equation

$$\sum_i \text{vech} \left(\frac{\nabla^2 f_i}{\nabla^2 \theta^2} \right) (y_i - f_i(\theta)) = \sum_i r_i + d$$

where $r = \text{vech}(G^0 G)$ and $d = \text{vech}(D)$: Thus, $\nabla S(\theta) = \nabla \theta = 0$ and $\nabla^2 S(\theta) = \nabla^2 \theta^2 = \sum_i D_i$ are equivalent to the following system of linear equations for y :

$$G_1^0 (y_1 - f_1) + G_2^0 (y_2 - f_2) = 0; \quad (10.2)$$

$$U_1 (y_1 - f_1) + U_2 (y_2 - f_2) = r + d; \quad \theta \in V; D \in M \quad (10.3)$$

where $y_1 = (y_1; \dots; y_{m_1})^0$; $y_2 = (y_{m_1+1}; \dots; y_n)^0$ and the similar partition for f : Now we construct the following mapping:

$$a : (\mathbb{R}^0; d^0; y_2^0) \mapsto (y_1^0; y_2^0); \quad \mathbb{R}^{2V}; D \in M_{+}; y_2 \in \mathbb{R}^{n_1 - m_1} \quad (10.4)$$

where y_1 is the solution to (10.2,10.3) given $\mathbb{R}; d; y_2$. This solution always exists because matrix $[G_1; U_1]$ has full rank. Our current aim is to find the Jacobian of y_1 as a function of $\mathbb{R}; d$ and show it is nonsingular. Differentiating (10.2) with respect to \mathbb{R} and d we obtain

$$G_1^0 \frac{\partial y_1}{\partial \mathbb{R}} + D = 0; \quad G_1^0 \frac{\partial y_1}{\partial d} = 0;$$

Differentiating (10.3) with respect to \mathbb{R} and d we obtain

$$U_1^0 \frac{\partial y_1}{\partial \mathbb{R}} + R = \frac{\partial r}{\partial \mathbb{R}}; \quad U_1^0 \frac{\partial y_1}{\partial d} = I; \quad \text{where } R = \frac{\partial U_1^0}{\partial \mathbb{R}} y_1 + \frac{\partial U_2(y_2; f_2)}{\partial \mathbb{R}}.$$

Combining these equations we obtain

$$\det \begin{pmatrix} \frac{\partial y_1}{\partial \mathbb{R}}; \frac{\partial y_1}{\partial d} \end{pmatrix} = \det ([G_1; U_1]^0)^{-1} \det \begin{pmatrix} D & 0 \\ \frac{\partial r}{\partial \mathbb{R}} & I \end{pmatrix} = \det ([G_1; U_1]^0)^{-1} \det(D) \neq 0;$$

which proves the Jacobian is nonsingular. Now it suffices to apply the implicit function theorem for mapping (10.4). It has a nonsingular Jacobian and, consequently, in some neighborhood has inverse. Therefore, the 'observations', $y = (y_1; y_2)$ fill an n_1 dimensional region and for each y from it the SS has at least two local minima. Since the density of the distribution of the error term is positive, the probability to fall in this region is positive as well.

References

- [1] Abramowitz, M., and Stegun, I.A. (1972). Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables. National Bureau of Standards. New York.
- [2] Amari, S.-I. (1985). Differential Geometrical Methods in Statistics. Lecture Notes in Statistics, 28, Springer, Berlin.
- [3] Andrews, D.F., Bickel, P.J., Hampel, F.R., Huber, P.J., Rogers, W.H., and Tukey, J.W. (1972). Robust Estimates of Location. Princeton University Press, Princeton. New Jersey.
- [4] Bates, D.M., and Watts, D.G. (1988). Nonlinear Regression Analysis and Its Applications. Wiley. New York.

- [5] Chong, E.K., and Zak, S.H. (1996). An Introduction to Optimization. Wiley. New York.
- [6] Demidenko, E. (1981). Linear and Nonlinear Regression (in Russian). Nauka. Moscow.
- [7] Demidenko, E. (1989). Optimization and Regression (in Russian). Nauka. Moscow.
- [8] Demidenko, E. (1996). On the existence of the least squares estimate in nonlinear growth curve models of exponential type. Communications in Statistics, Theory and Methods, 25, 159-182.
- [9] Ezekiel, M., and Fox, K.A. (1959). Methods of correlation and regression analysis, linear and curvilinear. Wiley. New York. 1959.
- [10] Floudas, C.A. and Pardalos, P.M. (1992). Recent Advances in Global Optimization. Princeton University Press. New Jersey.
- [11] Gallant, A.R. (1987). Nonlinear Statistical Methods. Wiley. New York.
- [12] Hansen, E. (1992). Global Optimization Using Interval Analysis. Dekker. New York.
- [13] Horst, R. and Pardalos, P.M. (1995). Handbook of Global Optimization. Kluwer. Dordrecht.
- [14] Horst, R. and Tuy, H. (1996). Global Optimization. Deterministic Approaches. Springer. New York.
- [15] Huber, P.J. (1981). Robust Statistics. Wiley. New York.
- [16] Judge, G.C., Hill, R.C., Griffiths, W.E., Lütkepohl, H., and Lee, T.-C. (1982). Introduction to the Theory and Practice of Econometrics. Wiley. New York.
- [17] Kearfort, R.B. (1996). Rigorous Global Search. Kluwer. Dordrecht.
- [18] Mäkeläinen, T., Schmidt, K., and Styan, G.P.H. (1981). On the existence and uniqueness of the maximum likelihood estimate of a vector-valued parameter in fixed-size samples. Annals of Statistics, 9, 758-767.
- [19] Mockus, J. (1989). Bayesian Approach to Global Optimization.. Theory and Application. Kluwer. Dordrecht.
- [20] Montgomery, D.C., and Peck, E.A. (1992). Introduction to Linear Regression Analysis. Wiley. New York.

- [21] Nakamura, T. (1984). Existence theorems of a maximum likelihood estimate from a generalized censored data sample. *Annals of the Institute of Statistical Mathematics*, 36, 375-393.
- [22] Ortega, J.M., and Rheinboldt, W.C. (1970). *Iterative Solution of Nonlinear Equations in Several Variables*. Academic Press. New York.
- [23] Pazman, A. (1984). Nonlinear least squares - uniqueness versus ambiguity. *Math. Operationsforschung Stat., ser. Statistics*, 15, 323-336.
- [24] Pazman, A. (1993). *Nonlinear Statistical Models*, Kluwer. Dordrecht.
- [25] Pinter, J.D. (1996). *Global Optimization in Action*. Kluwer. Dordrecht.
- [26] Rockafellar, R.T. (1970). *Convex Analysis*, Princeton University Press. New Jersey.
- [27] Seber, G.A.F., and Wild, C.J. (1989). *Nonlinear Regression*. Wiley. New York.
- [28] Webster, R. (1994). *Convexity*, Oxford University Press. Oxford.
- [29] Zhigljavsky, A.A. (1991). *Theory of Global Random Search*. Kluwer. Dordrecht.