# Application of Stepwise Regression to Non-Linear Estimation

## R. I. Jennrich & P. F. Sampson

# Application of Stepwise Regression to Non-Linear Estimation*

## R. I. JENNRICH AND P. F. SAMPSON

### University of California, Los Angeles

This paper reviews three basic methods of non-linear least squares estimation and the best known modifications to the popular Gauss-Newton procedure. With regard to the Gauss-Newton procedure the paper identifies the problem of poor parameterization and gives simple examples to show how it may arise. A simple example is also given to show that one popular method of handling constrained boundaries can lead to erroneous results. A modification, based on stepwise linear regression techniques, to handle both of these problems is presented and discussed and an example is given to illustrate its effectiveness. The required fundamentals of stepwise linear regression are reviewed.

## 1. INTRODUCTION

Let $y_t$ denote a set of $n$ responses which have an assumed structure

$$(1) \qquad y_t = f_t(\theta) + e_t \qquad t = 1, \cdots, n.$$

The response function $f_t(\theta)$ is a known function of $t$ and a parameter vector $\theta$ which ranges over a parameter space $\Theta$ of the form

$$(2) \qquad \Theta = \{(\theta_1, \cdots, \theta_p): a_i \le \theta_i \le b_i ; i = 1, \cdots, p\}.$$

A parameter vector $\hat{\theta}$ is called a least squares estimate of some true and usually unknown parameter vector $\theta_0$ if $\hat{\theta}$ minimizes the residual sum of squares

$$(3) \qquad Q(\theta) = \sum_{t=1}^{n} (y_t - f_t(\theta))^2.$$

We will summarize briefly three methods of minimizing $Q(\theta)$. These are steepest descent, quadratic approximation, and the Gauss–Newton method.

The steepest descent method consists in moving the parameter vector $\theta$ a short distance in the direction opposite to that of the gradient vector $Q'(\theta) = ((\partial/\partial\theta_i)Q(\theta))$. The process is repeated until $\theta$ hopefully converges to a least squares estimate.

Quadratic approximation consists of approximating $Q(\theta)$ locally by a quadratic function and minimizing this function. The appropriate modification is made to the parameter vector and the whole process is repeated until it hopefully converges. Quadratic approximation is mathematically equivalent to using the Newton–Raphson method to solve the gradient equation $Q'(\theta) = 0$.

The Gauss–Newton method consists in linearly approximating the response function $f_i(\theta)$ viewed as a function of $\theta$ and solving the resulting linear least squares problem by the standard techniques of linear regression. To be more specific let $y$, $f(\theta)$, $e$ denote the vectors $(y_i)$, $(f_i(\theta))$, $(e_i)$ and let the vector $f_i'(\theta) = ((\partial/\partial\theta_i)f_i(\theta))$ denote the partial derivative of the vector valued function $f(\theta)$ with respect to the $i$th component of $\theta$. Equation (1) becomes $y = f(\theta) + e$. Replacing $f(\theta)$ in this equation by its local linear approximation

(4)          $$f(\theta + \Delta\theta) \approx f(\theta) + \Delta\theta_1 f_1'(\theta) + \cdots + \Delta\theta_p f_p'(\theta)$$

gives

(5)          $$y - f(\theta) = \Delta\theta_1 f_1'(\theta) + \cdots + \Delta\theta_p f_p'(\theta) + e$$

which may be viewed as a standard linear regression equation with the $f_i(\theta)$ playing the role of independent variables, the $\Delta\theta_i$ the role of regression coefficients and $y - f(\theta)$ the role of the dependent variable. Using this identification the $\Delta\theta_i$ are obtained by standard regression techniques. Replacing $\theta$ by $\theta + \Delta\theta$ the whole process is repeated until it hopefully converges.

Let $Q''(\theta)$ denote the matrix $((\partial/\partial\theta_i)\,(\partial/\partial\theta_j)Q(\theta))$. The modification $\Delta\theta$ of the vector $\theta$ prescribed by each of the three methods on a given step is for steepest descent:

(6)          $$\Delta\theta = -\alpha Q'(\theta), \qquad \alpha > 0$$

quadratic approximation:

(7)          $$\Delta\theta = -(Q''(\theta))^{-1}Q'(\theta)$$

Gauss–Newton:

(8)          $$\Delta\theta = -\tfrac{1}{2}(f'(\theta)^T f'(\theta))^{-1}Q'(\theta).$$

The step $\Delta\theta$ will in general be distinct in both length and direction for each of the three methods. The steepest descent equation prescribes only the direction and not the length of $\Delta\theta$. The latter is usually obtained by some form of a one dimensional search. Of the three the Gauss–Newton method seems to be by far the most popular probably because 1) it specifies the step size as well as its direction, 2) it does not require the computation of the $p(p + 1)/2$ second partial derivatives of each of the functions $f_i(\theta)$ which is required by the quadratic approximation method and 3) in the special case when the response function is linear in the parameters it converges in a single step.

## 2. Modifications of the Basic Gauss-Newton Procedure

The most commonly used, and perhaps most important, modification of the basic Gauss–Newton procedure is the use of partial steps $\pi\Delta\theta$ in place of the full step $\Delta\theta$ specified by (8). When the linear approximations on which the Gauss–Newton method is based fail the full step $\Delta\theta$ frequently results in an increase in the residual sum of squares $Q(\theta)$. It can be shown however that if $\theta$ is not already a minimizing value a sufficiently small step in the direction of $\Delta\theta$ will reduce the residual sum of squares and it is this observation which motivates the use of partial steps. The proportion $\pi$ of the step to be

used may be obtained by trying an arbitrary sequence of values say $1, \frac{1}{2}, \frac{1}{4}, \frac{1}{16}, \cdots$ until one which reduces the residual sum of squares is found or by some other one dimensional search procedure [4].

It can, and often does, happen that the vectors $\mathbf{f}'_1(\theta), \cdots, \mathbf{f}'_p(\theta)$ are, or are nearly, linearly dependent. This is equivalent to saying that the matrix $\mathbf{A}(\theta) = \mathbf{f}'^T(\theta)\mathbf{f}'(\theta)$ is, or is nearly, singular. Under these circumstances it may be difficult or impossible to compute the basic step $\Delta\theta$ accurately. One could describe the problem by saying the response function is poorly parameterized. One solution is to reparameterize ([1], p. 803) but this is often difficult and usually inconvenient. Even in the case of linear regression where it is comparatively simple to reparameterize most workers prefer to use the parameters which are natural to the problem at hand and, if necessary, employ a modified form of regression called "stepwise regression" to avoid problems of singularity. In this modification variables are selected one at a time in the order which "does the most good" in reducing the residual sum of squares. If singularity problems arise (i.e., the tolerance discussed below becomes too low) one or more of the variables may be omitted. In this paper we propose the local application of these same stepwise techniques to deal with the problem of poor parameterization which frequently arises in non-linear regression. It turns out that stepwise techniques also provide a convient and sensible way to handle boundary restrictions on the components of the parameter vector. This is done by simply not modifying a parameter which has a boundary value unless the resulting modification is toward the interior of the parameter region. This amounts to "solving the problem on the boundary" whenever it appears that the absolute best fit parameter vector lies outside the parameter region.

Another solution to the poor parameterization problem can be found in the modification of Marquardt [5] which consists in replacing the matrix $\mathbf{A}(\theta)$ in the Gauss–Newton method by $\mathbf{A}(\theta) + \lambda\mathbf{I}$ where $\mathbf{I}$ denotes an identity matrix. This method may be viewed as a compromise between the Gauss–Newton and steepest descent methods. When $\lambda = 0$ it gives a Gauss–Newton step. As $\lambda \rightarrow \infty$ the step shortens and its direction approches that of the steepest descent method. The point to note here is that even if the matrix $A(\theta)$ is singular, the matrix $\mathbf{A}(\theta) + \lambda\mathbf{I}$ is non-singular whenever $\lambda > 0$ and can be made arbitrarily well conditioned by choosing $\lambda$ sufficiently large. Values of $\lambda$ which are too large, however, lead to slow convergence. The usual practice is to adjust the value of $\lambda$ dynamically from one iteration to the next by a rule which is based on the past behavior of the algorithm.

## 3. Near-Singularity

There are many reasons why the sum of cross products matrix $\mathbf{A}(\theta)$ is frequently singular or nearly singular. In this section we will point out how this difficulty may arise when attempting to fit a sum of exponentials. Let

$$(9) \qquad f_t(\theta) = e^{\theta_1 t} + e^{\theta_2 t}.$$

In the notation of Section 1

$$(10) \qquad \mathbf{f}'_1(\theta) = (te^{\theta_1 t}) \quad \text{and} \quad \mathbf{f}'_2(\theta) = (te^{\theta_2 t}).$$

When $\theta_1 = \theta_2$ these two vectors are equal and thus linearly dependent. If the data $y_t$ follow a straight line or any concave-down path as $t$ ranges from 1 to $n$ it can be shown that the least squares values $\hat{\theta}_1$ and $\hat{\theta}_2$ must be equal and hence that $\mathbf{A}(\theta)$ is singular at a least squares solution and nearly singular in a neighborhood of it. Less precise but perhaps more general insight can be obtained by considering the slightly more general sum of exponentials

11)                    $$f_t(\theta) = \theta_1 e^{\theta_3 t} + \theta_2 e^{\theta_4 t}.$$

If the fitted response function $f_t(\hat{\theta})$ follows the path of a single exponential, say $\alpha e^{\beta t}$, it can be shown that $\mathbf{A}(\hat{\theta})$ is singular and if the fitted response function nearly follows such a path $A(\hat{\theta})$ is nearly singular. The same statement can be made when the fitted response follows or nearly follows any striaght line path.

A still less precise and more general statement is that singularity problems are likely to arise whenever there is a curve in the parameter space along which the response function $f_t(\theta)$ differs little from the least squares response $f_t(\hat{\theta})$; that is, whenever it is possible to vary $\theta$ substantially without substantially affecting the best fit. A sometimes useful, though obviously over-simplified, statement is that trouble arises when there are too many parameters.

## 4. STEPWISE REGRESSION

In this section we will review the fundamentals of stepwise regression. The method is based on the Gauss–Jordan pivot. Let $(a_{ij})$ be an $n$ by $n$ matrix whose $k$th diagonal element $a_{kk}$ is non-zero. The result of a Gauss–Jordan pivot on this element is a new matrix whose $ij$th element

(12)           $$\tilde{a}_{ij} \begin{cases} a_{ij} - a_{ik}a_{kj}/a_{kk} & i \neq k, \quad j \neq k \\ -a_{ik}/a_{kk} & i \neq k, \quad j = k \\ a_{kj}/a_{kk} & i = k, \quad j \neq k \\ 1/a_{kk} & i = k, \quad j = k. \end{cases}$$

Gauss–Jordan pivots commute. That is, a pivot on the $i$th diagonal element followed by a pivot on the $j$th diagonal element produces the same result as pivots on the $i$th and $j$th diagonal elements in reverse order. Moreover a Gauss–Jordan pivot is its own inverse. That is, two applications of the same pivot leaves the matrix $(a_{ij})$ unchanged. If $(a_{ij})$ is partitioned as the matrix on the left below

13)     $$\begin{bmatrix} \mathbf{A}_{11}, & \mathbf{A}_{12} \\ \mathbf{A}_{21}, & \mathbf{A}_{22} \end{bmatrix} \rightarrow \begin{bmatrix} \mathbf{A}_{11}^{-1}, & \mathbf{A}_{11}^{-1}\mathbf{A}_{12} \\ -\mathbf{A}_{21}\mathbf{A}_{11}^{-1}, & \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12} \end{bmatrix}$$

and if it is possible to pivot on all of the diagonal elements of $\mathbf{A}_{11}$ (i.e., no zeros are encountered) then $\mathbf{A}_{11}$ is non-singular and the result of pivoting on its diagonal elements is given on the right in (13). In particular if a matrix is pivoted on all of its diagonal elements the result is its inverse.

If a matrix is non-negative (i.e., positive semi-definite) then its diagonal

elements are non-negative and they remain so after any number of pivots. A non-negative matrix is singular if and only if for any order of pivoting at least one diagonal element becomes zero before all diagonal elements have been used as pivots. Let $\bar{a}_{ii}$ denote the value of the $i$th diagonal element of a non-negative matrix $(a_{ij})$ after pivoting on a number of other diagonal elements. The ratio $\bar{a}_{ii}/a_{ii}$ is called the tolerance of the $i$th diagonal element at this state. Since $\bar{a}_{ii}$ is obtained from $a_{ii}$ by a sequence of subtractions of non-negative quantities the tolerance of the $i$th diagonal element in a sense measures its precision. The use of an element with low tolerance as a pivot may well destroy the precision of the entire matrix. For this reason, in practice, elements with low tolerance (say $10^{-5}$ on an 8 decimal place computer) are not used as pivots.

Let $X = (x_{ij})$ denote an $n$ by $p$ matrix and let $Y = (y_{ij})$ denote an $n$ by $q$ matrix. The columns $x_i = (x_{ij})$ and $y_i = (y_{ij})$ of these matrices will be called variables. Assume the matrix $X^T X$ is non-singular and let the matrix on the right in (14) denote the result of pivoting the matrix on the left on the diagonal elements of $X^T X$

$$(14) \qquad \begin{bmatrix} X^T X, & X^T Y \\ Y^T X, & Y^T Y \end{bmatrix} \rightarrow \begin{bmatrix} (X^T X)^{-1}, & \beta \\ -\beta, & S \end{bmatrix}.$$

It follows from (13) and the standard formulas of regression theory that the $ij$th component $\beta_{ij}$ of the matrix $\beta$ is the regression coefficient for variable $x_i$ which results from regressing variable $y_j$ on the entire set of $x$-variables. That is

$$(15) \qquad \hat{y}_j = \beta_{1j} x_1 + \cdots + \beta_{pj} x_p$$

is the linear regression of $y_j$ on $x_1, \cdots, x_p$. Moreover the matrix $S = (s_{ij})$ is the sum of cross products matrix for the residuals $y_i - \hat{y}_i$. In particular its $i$th diagonal element is the residual sum of squares resulting from regressing $y_i$ on all of the $x$-variables. Thus all the essential building blocks of a standard regression analysis are readily available in the matrix on the right in expression (14).

Stepwise regression [3] is based on the observation that a variable can be moved from the set of $y$-variables to the set of $x$-variables or from the set of $x$-variables to the set of $y$-variables by simply pivoting on the corresponding diagonal element of the matrix on the right in expression (14). In standard applications one of the $y$-variables, say $y_d$, is the real dependent variable. The object at each step is to move that $y$-variable (other than $y_d$) into the set of $x$-variables (the set of variables regressed upon) which gives the greatest possible reduction in the residual sum of squares of the dependent variable $y_d$. (Equivalently this is the $y$-variable whose partial correlation with $y_d$, when partialed on the $x$-variables, is greatest; or the $y$-variable whose regression coefficient after entering the set of $x$-variables has the greatest $F$-value.) Since the reduction in the residual sum of squares resulting from entering the variable $y_i$ is $\beta_{id}^2/s_{ii}$ it is a simple matter to locate the required variable. The variable is entered providing the tolerance of the corresponding diagonal element of

the matrix on the right in expression (14) is above a specified threshold value. Thus the variable $y_i$ actually entered, if any, is the one among all $y_i \neq y_d$ , whose tolerance exceeds the prescribed threshold and whose $\beta_{id}^2/s_{ii}$ value is maximum. The process begins with all variables in the $y$-set. Entering one variable at a time into the $x$-set the process ends when no available variables satisfying the tolerance criterion remains or when the reduction in the residual sum of squares fails to exceed a specified threshold. (This latter threshold may alternatively be stated in terms of partial correlation or $F$-value thresholds).

Roughly speaking variables are entered in the order that (one step at a time at least) does the most good providing their entry does not do too much damage to the numerical presicion of the results. In ordinary applications to linear regression ([2], p. 233), regression coefficients and other intermediate results are produced at each step and variables whose removal does not greatly increase the residual sum of squares may be removed. In the applications below no such intermediate output is produced and variable removal is used only to prevent boundary violations.

## 5. DETAILS OF THE STEPWISE MODIFICATION

Augment the matrix $\mathbf{A}(\boldsymbol{\theta})$ by a $d = p + 1$st row and column so that its components become

$$a_{ij}(\boldsymbol{\theta}) = \sum_{t=1}^{n} \frac{\partial}{\partial \theta_i} f_t(\boldsymbol{\theta}) \frac{\partial}{\partial \theta_j} f_t(\boldsymbol{\theta}) \qquad i, j = 1, \cdots, p$$

(16)     $$a_{id}(\boldsymbol{\theta}) = a_{di}(\boldsymbol{\theta}) = \sum_{t=1}^{n} \frac{\partial}{\partial \theta_i} f_t(\boldsymbol{\theta})(y_t - f_t(\boldsymbol{\theta})) \qquad i = 1, \cdots, p$$

$$a_{dd}(\boldsymbol{\theta}) = \sum_{t=1}^{n} (y_t - f_t(\boldsymbol{\theta}))^2.$$

A given iteration step proceeds as follows. Identify the $d$th diagonal element of the augmented matrix $\mathbf{A}(\boldsymbol{\theta})$ with the dependent variable and perform Gauss–Jordan pivots in the order specified by stepwise regression. While the threshold tolerance should be a parameter of the computational procedure, a value of about $10^{-5}$ on an 8 decimal place computer will assure about 3 decimal places of precision in the pivots used and seems to be a satisfactory threshold. The components of $\Delta\boldsymbol{\theta}$ are given by

(17)     $$\Delta\theta_i = \begin{cases} \bar{a}_{id} & \text{if } \mathbf{A}(\boldsymbol{\theta}) \text{ has been pivoted on its } i\text{th diagonal element} \\ 0 & \text{otherwise.} \end{cases}$$

Here $\bar{a}_{id}$ denotes the value of the $id$th element of $\mathbf{A}(\boldsymbol{\theta})$ after pivoting. Next the maximum value $\alpha \leq 1$ for which $\boldsymbol{\theta} + \alpha \Delta\boldsymbol{\theta}$ satisfies the boundary conditions (2) is obtained. If $\alpha = 0$ then $\boldsymbol{\theta}$ lies on at least one coordinate boundary and the corresponding coordinate of $\Delta\boldsymbol{\theta}$ is positive if the boundary is an upper boundary and negative if it is a lower boundary. As soon as such a coordinate

of $\Delta\theta$ is encountered the corresponding diagonal element of $A(\theta)$ is unpivoted by means of an inverse Gauss–Jordan pivot. Using the resulting $\Delta\theta$ a new $\alpha$ is obtained as before and the process is continued until an $\alpha > 0$ is obtained or all diagonal elements are unpivoted. If an $\alpha > 0$ is obtained $\theta$ is replaced by $\theta + \alpha \Delta\theta$ completing the iteration step. If all diagonal elements are unpivoted (in the experience of the authors this has never happened) the entire iteration process is terminated. In effect what has been done is that the coordinates of $\theta$, which already have boundary values and whose modification would lead to boundary violations, are held fixed at their boundary values while a standard Gauss–Newton iteration step is performed on remaining coordinates which previously satisfied the tolerance criterion.

While this completes the description of the stepwise modification of a Gauss–Newton iteration step it is probably a good idea to also use some form of the partial step modification described in Section 2. Beginning with the step $\alpha \Delta\theta$ the value of $\alpha$ is further reduced, if necessary, until a reduction in the residual sum of squares is obtained.

The entire procedure then consists of selecting a starting value for $\theta$ and applying iteration steps until some scale free convergence criterion (such as a relative change in the residual sum of squares of less than $10^{-7}$) is satisfied.

## 6. Some Practical Considerations

In most implementations of the Gauss–Newton algorithm the user supplies one or more subroutines to evaluate the functions $f_i(\theta)$ (which almost always have the form $f_i(\theta) = f(x_i, \theta)$) and their partial derivatives. There is an advantage to doing this in a single subroutine since for many functions it is computationally only a little more expensive to evaluate the function and its partial derivatives than to evaluate the function alone. For example, if $f_i(\theta) = \theta_1 e^{\theta_2 t}$ then both of its partial derivatives $e^{\theta_2 t}$ and $t\theta_1 e^{\theta_2 t}$ can be obtained by one additional multiplication. This observation also suggests that extra evaluations of the residual sum of squares may be relatively expensive. In many applications they are nearly as expensive as an entire Gauss–Newton iteration step and for this reason the number of extra residual sum of squares evaluations in the partial step modification of Section 2 should be kept to a minimum.

A number of popular non-linear estimation routines handle boundary restrictions by computing an unconstrained step and projecting the results onto the appropriate boundary whenever a boundary violation would otherwise occur. When properly handled [6] this technique works for the steepest descent method. As can be seen from Figure 1 however it does not work for the Gauss–Newton method and most of its modified forms. The line in Figure 1 represents an upper boundary of a two-dimentional parameter region and the ellipse represents a constant value surface for $Q(\theta)$ corresponding to a response function which is linear in its parameters. Under these conditions the unconstrained minimum of $Q(\theta)$ occurs at the center $\theta_c$ of the ellipse which is outside of the parameter region. Moreover the constrained best fit occurs at the point of tangency $\theta$. Assume that an iteration step begins at any boundary point $\theta$. An unconstrained Gauss–Newton step will take it to $\theta_c$ and projection onto
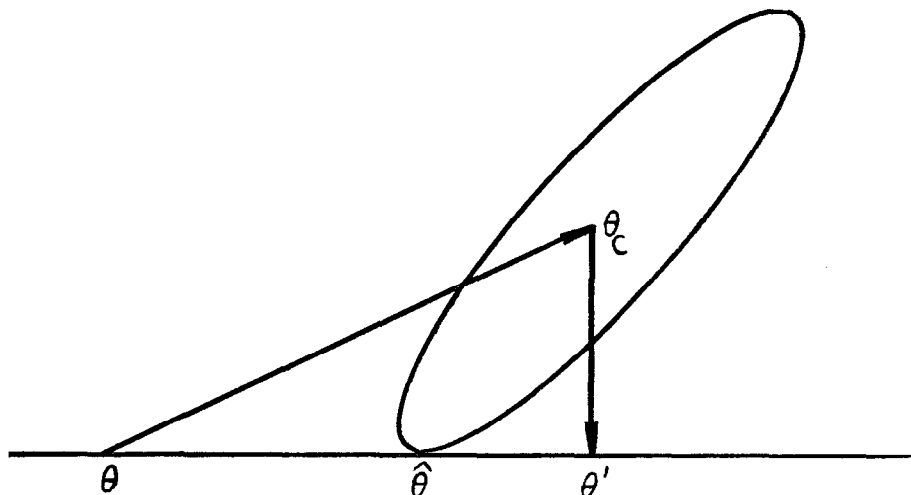
FIGURE 1—Figure showing the projection method of handling boundary restrictions is inappropriate when used with the Gauss-Newton method. The symbols are defined in Section 6.

the boundary to $\theta'$. Thus $\theta'$ is a limit point of such a procedure but, unfortunately, not a least squares solution.

Whenever possible it is a good idea to plot the original data and the fitted model or at least to list the residuals $y_t - f_t(\hat\theta)$. This gives some indication of the appropriateness of the assumed structure and may indicate a failure of the algorithm to converge to a minimizing solution (e.g., when it converges to a local, rather than a global, minimum). Statistical information, such as an estimate of the variance of the errors $e_t$ and estimates of the variances of the parameter estimates $\hat\theta_1, \cdots, \hat\theta_p$, is also helpful in interpreting and evaluating results (formulas for these and other statistical parameters can be found in the writeup of the computer program discussed in the next section).

## 7. EXAMPLE

To demonstrate that the stepwise modification discussed in the previous section works when the matrix $A(\theta)$ is singular or nearly singular we will use it to fit the response function (9) to the data

$$(18) \qquad\qquad y_t = 2 + 2t, \qquad t = 1, \cdots, 10.$$

For this problem, as pointed out in Section 3, the least squares values $\hat\theta_1$ and $\hat\theta_2$ are equal, $A(\theta)$ is singular at the least squares solution and nearly singular in a neighborhood of it.

Table 1 shows the convergence of the stepwise modification of the Gauss–Newton algorithm using a tolerance of $10^{-5}$. The accuracy of the final values $\hat\theta_1 = .25785$ and $\hat\theta_2 = .25780$ is demonstrated by the fact that they agree to four decimal places with the value $\hat\theta = .25783$ obtained by replacing $\theta_1$ and $\theta_2$ by a single parameter $\theta$ and solving the resulting well behaved one parameter least squares problem.

TABLE 1

*Convergence in the Example of Section 7 Using the Stepwise Modification*

| Iteration | Step Halving Freq. | $\theta_1$ | $\theta_2$ | Residual Mean Sq. |
|-----------|--------------------|--------|--------|-------------------|
| 0   | 0  | .30000 | .40000 | 521.41 |
| 1   | 2  | .35468 | .35515 | 429.84 |
| 2   | 0  | .35468 | .22546 | 88.154 |
| 3   | 4  | .34180 | .25697 | 83.743 |
| 10  | 0  | .22799 | .30565 | 21.330 |
| 30  | 0  | .25790 | .25774 | 15.545 |
| 60  | 10 | .25785 | .25780 | 15.545 |
| 100 | 10 | .25785 | .25780 | 15.545 |

The results of applying the Gauss–Newton algorithm without the stepwise modification (this was done using the stepwise alogrithm and setting the tolerance equal to zero) to the same problem are given in Table 2. The final residual mean square value 429.99 of Table 2 when compared with the corresponding value 15.545 of Table 1 indicates clearly that after 100 steps the iteration is far from a least squares solution. The iteration, unfortunately, did not obviously "blow-up" but instead looks like it has converged rapidly. After the second step, at the beginning of which the parameters have nearly equal values, only very small changes occur. While the step halving frequency is high (it had an imposed maximum of 11 corresponding to one 2048th of a full step) this is also a characteristic of convergence.

A computer program (7094 FORTRAN IV) which implements the stepwise method discussed in Section 5 has been written by one of the authors. This program was used to perform the calculations required in the previous example and has also been used quite successfully for the past two years on a large variety of problems arising in medical research. The program and its writeup may be obtained by writing to the Program Librarian at the Health Sciences Computing Facility, University of California, Los Angeles, California 90024.

TABLE 2

*Convergence in the Example of Section 7 Without the Stepwise Modification*

| Iteration | Step Halving Freq. | $\theta_1$ | $\theta_2$ | Residual Mean Sq. |
|-----------|--------------------|--------|--------|-------------------|
| 0   | 0  | .30000 | .40000 | 521.41 |
| 1   | 2  | .35468 | .35515 | 429.84 |
| 2   | 10 | .35503 | .35481 | 429.84 |
| 3   | 11 | .35468 | .35515 | 429.85 |
| 10  | 10 | .35503 | .35481 | 429.85 |
| 30  | 10 | .35503 | .35481 | 429.88 |
| 60  | 10 | .35504 | .35481 | 429.93 |
| 100 | 10 | .35505 | .35481 | 429.99 |

REFERENCES

1. Box, G. E. P., 1960. Fitting empirical data. *Annals of the New York Academy of Sciences, 86,* 792–816.
2. Dixon, W. J., 1965. Editor. *BMD Biomedical Computer Programs,* UCLA Student Store, 308 Westwood Blvd., Los Angeles, California 90024.
3. Efroymsen, M. A., 1960. Multiple regression analysis. *Mathematical Methods for Digital Computers.* Edited by A. Ralston and H. S. Wilf, John Wiley and Sons, New York, 191–203.
4. Hartley, H. O., 1961. The modified Gauss-Newton method for the fitting of non-linear regression functions by least squares. *Technometrics, 3,* 269–280.
5. Marquardt, D. W., 1963. An algorithm for the estimation of non-linear parameters. *Society for Industrial and Applied Mathematics Journal, 11,* 431–441.
6. Rosen, J. B., 1960. The gradient projection method for non-linear programming. Part I. Linear Constraints. *Society for Industrial and Applied Mathematics Journal, 8,* 181–217.