

MATHEMATICS RESEARCH CENTER, UNITED STATES ARMY  
THE UNIVERSITY OF WISCONSIN

Contract No. : DA-11-022-ORD-2059

ITERATIVE METHODS FOR SOLVING  
NONLINEAR LEAST SQUARES PROBLEMS

Victor Pereyra

MRC Technical Summary Report #622  
February 1966

Madison, Wisconsin

## ABSTRACT

In many technical applications it is desired to fit a nonlinear model to a set of observations. Several iterative techniques have been devised in order to determine a best set of parameters in the least squares sense.

In this paper we discuss conditions for convergence, and give error estimates for a class of methods, which includes as particular cases some well known techniques. It is shown that those methods can be considered as modified Newton's iterations for a suitable functional equation, and then a general theorem, first indicated by Bartle, is proved and applied to this particular case. The hypotheses are set in such a way that their checking by an automatic computer is made possible.

Some numerical examples are given. The main aim is to show that the automatic error estimation procedure works, rather than attempting to optimize the computational scheme.

# ITERATIVE METHODS FOR SOLVING NONLINEAR LEAST SQUARES PROBLEMS

Victor Pereyra

## 1. Introduction.

In many technical applications it is desired to fit a nonlinear model to a set of observations. When the best fit is sought in the least squares sense the problem can be stated as follows:

Given the nonlinear transformation  $F(\underline{x}) = \underline{y}$  between the finite dimensional Euclidean spaces  $E^n$  and  $E^m$  ( $n \leq m$ ), and the vector of "observations"  $\underline{b} \in E^m$ , find a vector  $\underline{x}^* \in E^n$  which minimizes the  $L_2$ -norm of  $F(\underline{x}) - \underline{b}$ .

We consider in this paper a general class of iterative methods for finding stationary points of  $\|F(\underline{x}) - \underline{b}\|_2^2$ . If we call  $\underline{f}(\underline{x}) = F(\underline{x}) - \underline{b}$ , these methods are of the form:

$$\underline{x}_{\nu+1} = \underline{x}_{\nu} - [\underline{T}_{\nu}]^{-1} \underline{f}'(\underline{x}_{\nu})^T \underline{f}(\underline{x}_{\nu}) \quad (\nu = 0, 1, \dots), \quad (1.1)$$

where the  $\underline{T}_{\nu}$  are linear, nonsingular transformations of  $E^n$  in itself, and  $[\underline{f}']^T$  is the transpose of the Jacobian matrix of the transformation  $\underline{f}$ . In components

$$[\underline{f}'(\underline{x})^T \underline{f}(\underline{x})]_j = \sum_{i=1}^m \frac{\partial f_i}{\partial x_j} f_i.$$

By choosing  $\underline{T}_{\nu}$  appropriately we can obtain some well known methods used in the solution of nonlinear least squares problems. For instance, if

---

Sponsored in part by the Mathematics Research Center, United States Army under Contract No. : DA-11-022-ORD-2059, and in part by NASA Research Grant NGR-50-022-028.

$\underline{T}_\nu = \underline{f}'(\underline{x}_\nu)^T \underline{f}'(\underline{x}_\nu)$  then the resulting iteration corresponds to the Gauss-Newton method ( see Moore and Zeigler [ 7 ] or Hartley [ 4 ] and references therein).

If  $\underline{T}_\nu = \underline{f}'(\underline{x}_0)^T \underline{f}'(\underline{x}_0)$  for all  $\nu$ , then we could call this the simplified Gauss-Newton method. In general, schemes like those of Powell [ 9 ], and Jakovlev [ 5 ], which replace  $\underline{f}'(\underline{x}_\nu)^T \underline{f}'(\underline{x}_\nu)$  by approximate expressions could be considered into the class of methods described by ( 1. 1 ).

Sufficient conditions for the convergence of the Gauss-Newton method were obtained in Zadunaisky and Pereyra [ 11 ] by applying a standard fixed point theorem.

Observing that to find the stationary points of  $\| \underline{f}(\underline{x}) \|^2_2$  is equivalent to finding the zeros of its gradient, the problem is reduced to the solution of the  $n \times n$  system of nonlinear equations

$$\varphi(\underline{x}) = \underline{f}'(\underline{x})^T \underline{f}(\underline{x}) = \underline{0} . \quad (1. 2)$$

If we now regard the iteration ( 1. 1 ) as a method for solving ( 1. 2 ) then a general theorem by Bartle [ 1 ] can be applied in order to obtain sufficient conditions for convergence, and error bounds for the successive approximations. We include a proof of this theorem since in Bartle's paper it is only vaguely indicated.

Considering the solution of overdetermined systems of nonlinear equations of the form  $\underline{f}(\underline{x}) = \underline{0}$ ,  $\underline{f} : E^n \rightarrow E^m$ , Ben-Israel [ 2 ], [ 3 ] has studied the convergence of generalized versions of Newton's method and the simplified ( or modified ) Newton method. The generalizations consist in replacing the inverses of the Jacobian matrix of  $\underline{f}(\underline{x})$ , which appear in the standard nonsingular  $n \times n$  case, by their pseudoinverses. This in particular permits considerations of the case in which  $\text{rank } \underline{f}'(\underline{x}) < n$ .

In this paper we will only be concerned with the case of full rank ( $=n$ ).

Since in that case

$$[\underline{f}'(\underline{x})^T \underline{f}'(\underline{x})]^{-1} \underline{f}'(\underline{x})^T \equiv \text{pseudoinverse of } \underline{f}'(\underline{x}) \equiv [\underline{f}'(\underline{x})]^\dagger, \quad (1.3)$$

we see that the Gauss-Newton method coincides with Ben-Israel's generalized Newton-Raphson's method.

The conditions obtained from Bartle's theorem for the general iteration (1.1) are specialized to the Gauss-Newton method in Theorem 4.1, and the representation (1.3) for the pseudoinverse of  $\underline{f}'(\underline{x})$  allows us to give an error estimation procedure which can be implemented on a digital computer. Ben-Israel's theorem and our result do not seem to be comparable. On one side we require conditions on the second derivatives of the function  $\underline{f}(\underline{x})$  while he does not. On the other hand we require only the boundedness of  $\|[\underline{f}'(\underline{x}_0)^T \underline{f}'(\underline{x}_0)]^{-1}\|$  while he uses the condition  $\|[\underline{f}'(\underline{x})]^\dagger - [\underline{f}'(\underline{y})]^\dagger\| \leq N\|\underline{x}-\underline{y}\|$  which is obviously more difficult to verify.

An implementation of the error estimation procedure is briefly explained in Section 5. Finally we present in Section 6 two numerical examples.

## 2. Approximate Newton type iteration.

As we said in the Introduction, the problem of finding the stationary values of  $\| \underline{F}(\underline{x}) - \underline{b} \|_2^2$  is equivalent to that of finding the solutions of the  $n \times n$  system of nonlinear equations:

$$\underline{\varphi}(\underline{x}) = \underline{f}'(\underline{x})^T \underline{f}(\underline{x}) = \underline{0} \quad (2.1)$$

where  $\underline{f}(\underline{x}) = \underline{F}(\underline{x}) - \underline{b}$ , and  $\underline{f}'(\underline{x})$  is its Jacobian matrix. If we put  $\underline{f}'(\underline{x})^T \underline{f}'(\underline{x}) \equiv \underline{N}(\underline{x})$  and assume that  $\underline{f}(\underline{x})$  is twice Fréchet differentiable in a certain region  $\Omega \subset E^n$ , then the Fréchet derivative (Jacobian) of  $\underline{\varphi}(\underline{x})$  can be written as

$$\underline{\varphi}'(\underline{x}) = \underline{N}(\underline{x}) + [\underline{f}'(\underline{x})^T]^T \underline{f}(\underline{x}) \quad (2.2)$$

This is easily obtained by applying the operational calculus with Fréchet derivatives (cf. Dieudonné [12]).

In the benefit of those readers not familiar with this calculus we will obtain this formula by operating on the components of the vector functions involved. We will use tensor notation with the summation convention. On the first place

$$\underline{\varphi}(\underline{x}) = \frac{\partial f_i}{\partial x_j} f_i,$$

and

$$\underline{\varphi}'(\underline{x}) = \frac{\partial}{\partial x_k} \left( \frac{\partial f_i}{\partial x_j} f_i \right) = \frac{\partial f_i}{\partial x_j} \frac{\partial f_i}{\partial x_k} + \frac{\partial^2 f_i}{\partial x_k \partial x_j} f_i.$$

Thus, going back to the matrix notation we obtain (2.2).

The standard Newton-Kantorovich method for solving (2.1) is:

$$\underline{x}_{\nu+1} = \underline{x}_{\nu} - [\underline{\varphi}'(\underline{x}_{\nu})]^{-1} \underline{\varphi}(\underline{x}_{\nu}), \quad (2.3)$$

and the approximate Newton iteration (see Bartle [1]) is obtained if  $\underline{\varphi}'(\underline{x}_\nu)$  is replaced by linear, nonsingular operators  $\underline{T}_\nu$  which are close to  $\underline{\varphi}'(\underline{x}_0)$  in some sense. But this is essentially what is stated in equation (1.1). Thus we can apply Bartle's theorem to that iteration obtaining sufficient conditions for its convergence. We can write those conditions in the following form:

Theorem 2.1. Let  $\underline{T}_\nu$  be a sequence of linear nonsingular operators from  $E^n$  into itself, such that, for  $\underline{x}_0 \in E^n$  and  $\rho > 0$  the sphere  $S(\underline{x}_0, \rho) \subset \Omega$ ,

$$(a) \quad \|\underline{T}_\nu^{-1}\| \leq \lambda,$$

$$(b) \quad \|\underline{T}_\nu - \underline{\varphi}'(\underline{x}_0)\| \leq \epsilon,$$

(c) if  $\underline{x}, \underline{y} \in S(\underline{x}_0, \rho)$  then

$$\|\underline{\varphi}(\underline{x}) - \underline{\varphi}(\underline{y}) - \underline{\varphi}'(\underline{x}_0)(\underline{x} - \underline{y})\| \leq \beta \|\underline{x} - \underline{y}\|,$$

$$(d) \quad \|\underline{\varphi}(\underline{x}_0)\| \leq \eta,$$

$$(e) \quad k = \lambda(\beta + \epsilon) < 1, \quad r = \frac{\lambda\eta}{1-k} \leq \rho.$$

Under these conditions the iteration (1.1) is well defined and converges to a solution  $\underline{x}^*$  of  $\underline{\varphi}(\underline{x}) = \underline{0}$ . Furthermore,  $\|\underline{x}^* - \underline{x}_0\| \leq r$  and is the only solution contained in this sphere. The rapidity of the convergence is given by  $\|\underline{x}^* - \underline{x}_\nu\| \leq k^\nu r$ .

Proof: First of all from (1.1), (a), and (d) it follows that

$$\|\underline{x}_1 - \underline{x}_0\| \leq \lambda \|\underline{\varphi}(\underline{x}_0)\| \leq \lambda\eta \leq \rho. \quad (2.4)$$

Furthermore, by (1.1), (b), and (c)

$$\begin{aligned}
\|\underline{\varphi}(\underline{x}_1)\| &= \|\underline{\varphi}(\underline{x}_1) - \underline{\varphi}(\underline{x}_0) - \underline{T}_0(\underline{x}_1 - \underline{x}_0)\| \leq \\
&\leq \|\underline{\varphi}(\underline{x}_1) - \underline{\varphi}(\underline{x}_0) - \underline{\varphi}'(\underline{x}_0)(\underline{x}_1 - \underline{x}_0)\| + \|(\underline{\varphi}'(\underline{x}_0) - \underline{T}_0)(\underline{x}_1 - \underline{x}_0)\| \leq \\
&\leq (\beta + \epsilon) \|\underline{x}_1 - \underline{x}_0\| . \tag{2.5}
\end{aligned}$$

With formulas (2.4) and (2.5) we have started an induction argument.

Assume that, for  $\nu = 1, \dots, n$

$$\begin{aligned}
(i_\nu) \quad &\|\underline{x}_\nu - \underline{x}_0\| \leq \rho , \\
(ii_\nu) \quad &\|\underline{x}_\nu - \underline{x}_{\nu-1}\| \leq \lambda \|\underline{\varphi}(\underline{x}_{\nu-1})\| , \\
(iii_\nu) \quad &\|\underline{\varphi}(\underline{x}_\nu)\| \leq (\beta + \epsilon) \|\underline{x}_\nu - \underline{x}_{\nu-1}\| .
\end{aligned}$$

From (1.1) and (a) we obtain

$$(ii_{n+1}) \quad \|\underline{x}_{n+1} - \underline{x}_n\| \leq \lambda \|\underline{\varphi}(\underline{x}_n)\| ,$$

and by  $(iii_\nu)$ ,  $(ii_\nu)$ ,  $(iii_{\nu-1})$ ,  $\dots$

$$\|\underline{x}_{\nu+1} - \underline{x}_\nu\| \leq \lambda(\beta + \epsilon) \|\underline{x}_\nu - \underline{x}_{\nu-1}\| \leq \dots \leq [\lambda(\beta + \epsilon)]^\nu \|\underline{x}_1 - \underline{x}_0\| < k^\nu \|\underline{x}_1 - \underline{x}_0\| .$$

Furthermore

$$\|\underline{x}_{n+1} - \underline{x}_0\| \leq \sum_{\nu=0}^n \|\underline{x}_{\nu+1} - \underline{x}_\nu\| \leq \sum_{\nu=0}^n k^\nu \|\underline{x}_1 - \underline{x}_0\| \leq \frac{\lambda \eta}{1-k} \leq \rho$$

which is  $(i_{n+1})$ .

Having proved this we can use (c) on

$$\begin{aligned}
\|\underline{\varphi}(\underline{x}_{n+1})\| &= \|\underline{\varphi}(\underline{x}_{n+1}) - \underline{\varphi}(\underline{x}_n) - \underline{T}_n(\underline{x}_{n+1} - \underline{x}_n)\| \leq \\
&\leq \|\underline{\varphi}(\underline{x}_{n+1}) - \underline{\varphi}(\underline{x}_n) - \underline{\varphi}'(\underline{x}_n)(\underline{x}_{n+1} - \underline{x}_n)\| + \\
&+ \|(\underline{\varphi}'(\underline{x}_n) - \underline{T}_n)(\underline{x}_{n+1} - \underline{x}_n)\| \leq (\beta + \epsilon) \|\underline{x}_{n+1} - \underline{x}_n\|
\end{aligned}$$



thus obtaining (iii<sub>n+1</sub>) which completes the induction argument.

With this we can establish that the sequence  $\{\underline{x}_\nu\}$  generated by (1.1) is a Cauchy sequence and simultaneously we can obtain the error estimation.

In fact, for any  $p > 0$

$$\|\underline{x}_{\nu+p} - \underline{x}_\nu\| \leq \sum_{i=1}^p \|\underline{x}_{\nu+i} - \underline{x}_{\nu+i-1}\| \leq k^\nu \sum_{i=0}^{p-1} k^i \|\underline{x}_1 - \underline{x}_0\| \leq \frac{\lambda \eta k^\nu}{1-k} = k^\nu r .$$

Consequently, there exists  $\underline{x}^* \in S(\underline{x}_0, r)$  such that  $\underline{x}_\nu \rightarrow \underline{x}^*$ , and furthermore

$$\|\underline{x}^* - \underline{x}_\nu\| \leq k^\nu r .$$

If  $\underline{x}^{**}$  is another solution of  $\underline{\varphi}(\underline{x}) = \underline{0}$  satisfying  $\|\underline{x}^{**} - \underline{x}_0\| \leq \rho$  we can write

$$\begin{aligned} \|\underline{x}^{**} - \underline{x}^*\| &= \|\underline{T}_0^{-1} \underline{T}_0(\underline{x}^{**} - \underline{x}^*)\| \leq \lambda \|\underline{\varphi}(\underline{x}^{**}) - \underline{\varphi}(\underline{x}^*) - \underline{T}_0(\underline{x}^{**} - \underline{x}^*)\| \leq \\ &\leq \lambda(\beta + \epsilon) \|\underline{x}^{**} - \underline{x}^*\| < \|\underline{x}^{**} - \underline{x}^*\| \end{aligned}$$

which is impossible unless  $\underline{x}^{**} = \underline{x}^*$ .

3. Another set of sufficient conditions.

Sometimes it is computationally advantageous to replace hypothesis (c) in Theorem 2.1 by other conditions which, though giving a slightly less general result, are more easy to handle.

Lemma 3.1. Suppose that  $\gamma, K_2 > 0$  exist such that

$$(c_1) \quad \|\underline{f}'(\underline{x})^T \underline{f}(\underline{x}) - \underline{f}'(\underline{x}_0)^T \underline{f}(\underline{x}_0)\| \leq \gamma, \quad \underline{x} \in S(\underline{x}_0, \rho),$$

$$(c_2) \quad \|\underline{N}'(\underline{x})\| \leq K_2, \quad \underline{x} \in S(\underline{x}_0, \rho),$$

are satisfied, then (c) follows with  $\beta = K_2 \rho + \gamma$ .

Proof: We will first calculate a bound for  $\|\underline{\varphi}'(\underline{x}) - \underline{\varphi}'(\underline{x}_0)\|$ ,  $\underline{x} \in S(\underline{x}_0, \rho)$ . By

(2.2) and the hypotheses

$$\begin{aligned} \|\underline{\varphi}'(\underline{x}) - \underline{\varphi}'(\underline{x}_0)\| &\leq \|\underline{N}(\underline{x}) - \underline{N}(\underline{x}_0)\| + \|\underline{f}'(\underline{x})^T \underline{f}(\underline{x}) - \underline{f}'(\underline{x}_0)^T \underline{f}(\underline{x}_0)\| \leq \\ &\leq K_2 \|\underline{x} - \underline{x}_0\| + \gamma \leq K_2 \rho + \gamma. \end{aligned}$$

If we call  $K_2 \rho + \gamma = \beta$  then by Bartle's Lemma 1:

$$\|\underline{\varphi}(\underline{x}) - \underline{\varphi}(\underline{y}) - \underline{\varphi}'(\underline{x}_0)(\underline{x} - \underline{y})\| \leq \beta \|\underline{x} - \underline{y}\|$$

which is (c).

4. The Gauss-Newton method.

By taking  $\underline{T}_\nu = \underline{N}(\underline{x}_\nu)$  in (1.1) we obtain the well known Gauss-Newton method for solving nonlinear least squares problems. We will express now the conditions (a) and (b) of Theorem 2.1 in terms of some of the quantities used in Section 3. If  $(c_1)$  and  $(c_2)$  hold then

$$\| [\underline{f}'(\underline{x}_0)^T]' \underline{f}(\underline{x}_0) \| \leq \gamma \quad (4.1)$$

implies condition (b) of Theorem 2.1 with  $\epsilon = \beta$ . This follows easily from

$$\begin{aligned} \| \underline{N}(\underline{x}_\nu) - \underline{\varphi}'(\underline{x}_0) \| &\leq \| \underline{N}(\underline{x}_\nu) - \underline{N}(\underline{x}_0) \| + \| [\underline{f}'(\underline{x}_0)^T]' \underline{f}(\underline{x}_0) \| \leq \\ &\leq K_2 \rho + \gamma = \beta, \quad (\| \underline{x}_\nu - \underline{x}_0 \| \leq \rho). \end{aligned}$$

We want to show now that the existence and uniform boundedness of  $\underline{N}(\underline{x})^{-1}$  ( $\underline{x} \in \underline{S}(\underline{x}_0, \rho)$ ) is a consequence of the existence and boundedness of  $\underline{N}(\underline{x}_0)^{-1}$ . To do so we will state without proof a well known result of the theory of matrices:

Lemma 4.1. Let  $\underline{B}$  and  $\underline{C}$  be  $n \times n$  matrices. Assume that

- (i)  $\underline{B}$  is nonsingular and  $\| \underline{B}^{-1} \| \leq \alpha$ ;
- (ii)  $\| \underline{C} - \underline{B} \| \leq \delta$ ;
- (iii)  $\alpha \delta < 1$ ,

then  $\underline{C}$  is nonsingular and

$$\| \underline{C}^{-1} \| \leq \frac{\alpha}{1 - \alpha \delta} . \quad (4.2)$$

From this we can prove the following:

Lemma 4.2. If

(a<sub>1</sub>)  $\underline{N}(\underline{x}_0)$  is nonsingular and  $\|\underline{N}(\underline{x}_0)^{-1}\| \leq \frac{1}{2}\lambda$  then  $\underline{N}(\underline{x})^{-1}$  ( $\underline{x} \in S(\underline{x}_0, \rho)$ ) exists and  $\|\underline{N}(\underline{x})^{-1}\| < \lambda$ .

Proof: Take in Lemma 4.1

$$\underline{B} = \underline{N}(\underline{x}_0), \quad \underline{C} = \underline{N}(\underline{x}), \quad \alpha = \frac{1}{2}\lambda \quad \text{and} \quad \delta = K_2 \rho .$$

Since  $\alpha\delta = \frac{1}{2}K_2\rho$  we have by (3.1) that  $\alpha\delta < \frac{1}{4}$  and thus  $\underline{N}(\underline{x})$  is nonsingular.

Furthermore

$$\|\underline{N}(\underline{x})^{-1}\| \leq \frac{\alpha}{1-\alpha\delta} \leq \frac{2}{3}\lambda < \lambda .$$

Collecting these results together we can state the following theorem.

Theorem 4.3. With the same notation as above and for  $\underline{x}_0 \in \Omega$ , let us assume that

$$\|[\underline{f}'(\underline{x})^T]^t \underline{f}(\underline{x}) - [\underline{f}'(\underline{x}_0)^T]^t \underline{f}(\underline{x}_0)\| \leq \gamma, \quad \underline{x} \in S(\underline{x}_0, \rho) \subset \Omega, \quad (4.3)$$

$$\|\underline{N}'(\underline{x})\| \leq K_2, \quad \underline{x} \in S(\underline{x}_0, \rho), \quad (4.4)$$

$$\|[\underline{f}'(\underline{x}_0)^T]^t \underline{f}(\underline{x}_0)\| \leq \gamma, \quad (4.5)$$

and that  $\underline{N}(\underline{x}_0)$  is nonsingular. Define

$$\lambda \equiv 2 \|\underline{N}(\underline{x}_0)^{-1}\|, \quad (4.6)$$

and assume that  $k = 2\lambda(K_2\rho + \gamma) < 1$ . Define

$$0 < r \equiv \lambda \|\underline{\varphi}(\underline{x}_0)\| / [1 - 2\lambda(K_2\rho + \gamma)]. \quad (4.7)$$

Assume further that  $r \leq \rho$ , then the sequence  $\{\underline{x}_\nu\}$  defined by

$$\underline{x}_{\nu+1} = \underline{x}_\nu - \underline{N}(\underline{x}_\nu)^{-1} \underline{\varphi}(\underline{x}_\nu) \quad (4.8)$$

converges to the unique solution  $\underline{x}^*$  of  $\underline{\varphi}(\underline{x}) = 0$  in the sphere  $S(\underline{x}_0, r)$ .

Moreover, the rate of convergence is estimated by

$$\|\underline{x}^* - \underline{x}_\nu\| \leq k^\nu r. \quad (4.9)$$

## 5. Automatic error estimation.

A Fortran 63 program has been written for the CDC 1604 computer at the University of Wisconsin Computing Center, which implements the Gauss-Newton method with the error estimation procedure given in Theorem 4.3.

The hypotheses are checked automatically as the iteration proceeds, and as soon as they are satisfied it can be ensured that the process converges and the bound (4.9) used to estimate the norm of the error.

As seen in step V of the procedure described below, a relaxation technique is used in order to prevent divergence in the earlier stages of the iteration (cf. Hartley [4]).

An interesting feature is the use of interval arithmetic (cf. Moore [6], Reiter [10]) to compute  $\gamma$  and  $K_2$  in Theorem 4.3. Given the region  $\Omega$  as an  $n$ -dimensional hypercube  $\Sigma(\tilde{\underline{x}}, \rho)$  with edge  $2\rho$  and center  $\tilde{\underline{x}}$ , we compute (4.3) and (4.4) in interval arithmetic with the argument  $\underline{x} = (\tilde{x}_i - \rho, \tilde{x}_i + \rho)$  and from there we can obtain  $\gamma$  and  $K_2$  immediately.

A further sophistication would be to use a program to generate the code for the partial derivatives which are needed in the discussion. We have not done this since our test cases were very simple, but such a program is available (Reiter [10]) and it has been successfully applied to the solution of some complicated systems of nonlinear equations by Newton's method.

We will briefly describe now the computational scheme. We use in this description informal Algol. For details on the Algorithmic Language Algol, cf. Naur [8].

The notation is the same as in Section 4. The norm used is the  $L_\infty$  norm

except for the residual  $\|f\|_2^2$  where, of course, we use the  $L_2$  norm.

Let  $\underline{x}_0$  and  $\rho$  be given.

begin procedure error;

I:  $\lambda_0 := 2 \| \underline{N}(\underline{x}_0)^{-1} \|$  ;

if  $\underline{N}(\underline{x}_0)$  is singular then go to error 1;

II:  $\Delta \underline{x}_0 := \underline{N}(\underline{x}_0)^{-1} \underline{f}'(\underline{x}_0)^T \underline{f}(\underline{x}_0)$ ;  $\bar{\rho} := \rho$  ;

III:  $\tilde{\gamma} := \max_{\underline{x} \in \Sigma(\underline{x}_0, \bar{\rho})} ( \| [\underline{f}'(\underline{x})]^T \underline{f}(\underline{x}) - [\underline{f}'(\underline{x}_0)]^T \underline{f}(\underline{x}_0) \| )$  ;

$\tilde{\gamma} := \| [\underline{f}'(\underline{x}_0)]^T \underline{f}(\underline{x}_0) \|$  ;

$\gamma_0 := \max(\tilde{\gamma}, \tilde{\gamma})$  ;

IV: if  $2\gamma_0\lambda_0 < 1$  then go to Rest of error procedure else

$\bar{\rho} := 0.1\bar{\rho}$  ; if  $\bar{\rho} > \lambda_0 \| \underline{f}(\underline{x}_0) \|$  then go to III;

V: comment this is the relaxation routine;

for  $i := 0$  step 1 until p do begin

comment p is a given integer;

$\underline{x}_1^i = \underline{x}_0 - 2^{-i} \Delta \underline{x}_0$  ;

if  $\| \underline{f}(\underline{x}_1^i) \|_2 < \| \underline{f}(\underline{x}_0) \|_2$  then begin

$\underline{x}_0 := \underline{x}_1^i$  ; go to I end end ;

Rest of error procedure:  $K_2 := \max_{\underline{x} \in \Sigma(\underline{x}_0, \bar{\rho})} ( \| \underline{N}'(\underline{x}) \| )$  ;

$k := 2 \lambda (K_2 \rho + \gamma_0)$  ;

VI: if  $k < 1$  then begin  $R_0 := \frac{\lambda \|\varphi(\underline{x}_0)\|}{1 - k}$  ;

if  $R_0 > \rho$  then go to V;

$r := R_0$ ;  $q := \lceil \frac{\ln(r) - \ln(\epsilon)}{\ln(2)} \rceil + 1$  ;

comment  $\epsilon$  is the desired accuracy;

for  $i := 0$  step 1 until  $q$  do

$\underline{x}_{i+1} := \underline{x}_i - \underline{N}(\underline{x}_i)^{-1} \underline{f}'(\underline{x}_i)^T \underline{f}(\underline{x}_i)$  end else

go to V;

error 1: end .

## 6. Test cases.

The following test cases have been run on the CDC 1604:

$$1) F_i(x_1, x_2, x_3) = g(t_i, \underline{x}) = x_2 e^{x_1 t_i} + x_3,$$

$10 \leq m \leq 20$ ,  $|t_i| \leq 10$ ,  $y_i$  taken from tables with different accuracies

(for given  $\underline{x}^*$ ).

Sample results are given in Table 1 for  $y_i = 0.1 e^{t_i} - 5$  truncated at the fifth significant figure;  $m = 10$ ,  $-2 \leq t_i \leq 2.5$ .

The conditions were usually fulfilled when the iterates were quite close to the exact solution.

$$2) g(t_i, x_1, x_2, x_3) = x_2 \sin(x_1 t_i) + x_3.$$

Taking  $y_i = \sin t_i$  to 3D for  $t_i = (0.105, 0.25, 0.4, 0.55, 0.7, 0.9, 1.1, 1.25, 1.35, 1.45, 1.55, 1.57, 1.6)$ ,  $m = 13$ . We have obtained the results shown in Table 2.

In the second example we see that before the 3rd iteration the condition  $k < 1$  is not fulfilled. Then  $2\lambda\gamma$  becomes less than one and simultaneously  $k < 1$  is also fulfilled. Thus  $r$  can be calculated and it comes to be  $< \rho$ . The error estimations are then read in the corresponding column. In this example we have, instead of using the second part of VI, chosen to proceed as if every iteration were the first, since in this way the error estimation is much better than the one obtained directly from the theorem. Thus, in the 4th iteration we have two error estimations:  $7.9 \times 10^{-5}$  obtained by considering the third iterate as  $\underline{x}_0$ , and  $3.5 \times 10^{-10}$  if we recompute everything anew. For the 5th iterate



we give in parentheses the error estimation obtained if the 4th iterate is taken as  $\underline{x}_0$ .

The disparity showed by these estimations stem from the fact that in this case the convergence is faster than linear since  $\|\underline{f}(\underline{x}^*)\|_2^2$  itself is quite small. This does not have to be the case in more real problems in which the model cannot be expected to reproduce the observations very accurately.

The conclusion we can draw from these and other experiments we have carried out is that the most difficult part in this error procedure is the choice of an appropriate  $\rho$ .

Also it is clear that in its present form the conditions are fulfilled only when we are quite close to the exact solution. On the other hand, we see from the definition of  $r$  that this quantity can be obtained with very little effort if  $k$  is ignored and, from a practical point of view, small values of  $r$  could be enough assurance of convergence, and they could be used as error estimators without any further check. In this case the only extra computation would be that of the step I of the error procedure.

Acknowledgement. The author is grateful to the referees for their most valuable comments and suggestions.

TABLE 1

iter.	$x_1$	$x_2$	$x_3$	$\ \underline{\varphi}(\underline{x})\ $	$\ \underline{f}(\underline{x})\ _2^2$	$\rho$	$2\lambda\gamma$	k	r
0	0.8	0.2	-4.5	17.75	2.08	1	$2.1 \cdot 10^6$	-	-
1	0.95520183514	0.08528910008	-4.9876152319	4.7	0.33	1	$8.1 \cdot 10^6$	-	-
2	1.01305373832	0.09946287544	-4.9995610656	0.6	0.039	1	$6.9 \cdot 10^6$	-	-
3	1.00025049364	0.09996678830	-4.9999737544	$6.6 \cdot 10^{-3}$	$4.6 \cdot 10^{-4}$	1	$6.7 \cdot 10^6$	-	-
4	0.99995249354	0.10000901153	-5.0000080515	$1.3 \cdot 10^{-6}$	$2.1 \cdot 10^{-5}$	$10^{-4}$	2.3	-	-
5	0.99995243331	0.10000903204	-5.0000080677	$3.7 \cdot 10^{-10}$	$2.1 \cdot 10^{-5}$	$10^{-6}$	0.025	0.075	$5.6 \cdot 10^{-9}$

TABLE 2

iter.	$x_1$	$x_2$	$x_3$	$\ \underline{\varphi}(\underline{x})\ $	$\ \underline{f}(\underline{x})\ _2^2$	$\rho$	$2\lambda\gamma$	k	r	error estimate
0	0.9	0.9	0.1	0.146	0.0339	1	2300	-	-	-
1	1.0261610437	0.98376291581	-0.0001976660	0.0695	0.00228	1	950	-	-	-
2	1.0003514003	0.99898261866	-0.0003389551	0.00629	0.00109	1	1029	-	-	-
3	0.99999658214	0.99979735620	-0.0003507988	$1.99 \cdot 10^{-6}$	0.00109	$10^{-3}$	0.066	0.9	$8.8 \cdot 10^{-5}$	$8.8 \cdot 10^{-5}$
4	0.99999671131	0.99979751510	-0.0003507976	$7.72 \cdot 10^{-11}$	0.00108	$8.8 \cdot 10^{-5}$	0.0058	0.079	$3.5 \cdot 10^{-10}$	$7.9 \cdot 10^{-5}$
5	0.99999671136	0.99979751510	-0.0003507976	$7.72 \cdot 10^{-11}$	0.00108		0.0058			$7.1 \cdot 10^{-5}$ ( $2.8 \cdot 10^{-11}$ )

## REFERENCES

- [ 1 ] Bartle R. G. , "Newton's method in Banach spaces", Proc. A. M. S. 6 (1955), pp. 827-831.
- [ 2 ] Ben-Israel A. , "A modified Newton-Raphson method for the solution of equations", Israel J. Math. 3 (1965), pp. 94-98.
- [ 3 ] \_\_\_\_\_ "On Newton's method for the solution of systems of equations", to appear in J. Math. Anal. Appl.
- [ 4 ] Hartley H. O. , "The modified Gauss-Newton method for the fitting of nonlinear regression functions by least squares", Technometrics 3 (1961), pp. 269-280.
- [ 5 ] Jakovlev M. N. , "On the solution of nonlinear equations by iterations", Dokl. Akad. Nauk SSSR 156 (1964), pp. 522-524.
- [ 6 ] Moore R. E. , "The automatic analysis and control of error in digital computation based on the use of interval numbers", in Error in Digital Computation, Vol. I (ed. by L. B. Rall) (1965), Wiley, New York.
- [ 7 ] Moore R. H. and Zeigler R. K. , "The solution of the general least squares problem with special reference to high speed computers", Los Alamos Sc. Lab. LL-2367 (1959), 38 pp.
- [ 8 ] Naur P. , "Revised report on the algorithmic language Algol 60", Comm. ACM 6 (1963), pp. 1-17.
- [ 9 ] Powell M. J. D. , " A method for minimizing a sum of squares of nonlinear functions without calculating derivatives", Comp. J. 7 (1965), pp. 303-307.
- [ 10 ] Reiter A. , "Compiler of differentiable expressions (CODEX)", MRC Computer Program #1; "Interval arithmetic package (INTERVAL)", MRC Computer Program #2. University of Wisconsin, Madison (1965).
- [ 11 ] Zaduanisky P. and Pereyra V. , "On the convergence and precision of a process of successive differential corrections", to appear in Proc. IFIPS 65.
- [ 12 ] Dieudonné J. , "Foundations of Modern Analysis", Academic Press, New York (1960).