MATHEMATICS RESEARCH CENTER, UNITED STATES ARMY
/ THE UNIVERSITY OF WISCONSIN

Contract No.: DA-31-124-ARO-D-462

ACCELERATING THE CONVERGENCE

OF DISCRETIZATION ALGORITHMS

Victor Pereyra

MRC-Technical Summary Report #687
October 1966

Madison, Wisconsin

ABSTRACT

Acceleration techniques for discretization algorithms used in the approximate solution of nonlinear operator equations are considered.

Practical problems arising in the solution of large systems of nonlinear algebraic equations are discussed.

These techniques are applied to the approximate solution of mildly non-linear elliptic equations by finite differences, and several numerical examples are given.

# ACCELERATING THE CONVERGENCE OF DISCRETIZATION ALGORITHMS

Victor Pereyra

## Introduction

In a recent paper ( Pereyra [1966] ) we have developed the theory of the method of deferred corrections. Applying this procedure we were able to accelerate the convergence of discrete approximations to solutions of certain types of nonlinear operator equations in Banach spaces.

In Section I of this paper we present another method for accelerating convergence. This is the generalization of the well-known Richardson's "deferred approach to the limit" (Richardson [1910]) and we call it the method of successive extrapolations.

In Section II we state the principal results about the linear deferred correction method studied in our earlier paper. The analysis of both types of acceleration procedures is based on the general presentation given by Stetter [1965] for the discussion of the asymptotic behavior of the global discretization error in this kind of problem.

When applying these procedures to boundary value problems it is necessary to solve large systems of nonlinear algebraic equations. Recent papers by Bers [1953], Greenspan and Parter [1965], Ortega and Rockoff [1965], and Schechter [1962], have considered generalized Gauss-Seidel or relaxation methods for this. On the other hand, one can apply the standard Newton

method; this requires the solution of a large system of linear equations at each step, which may in turn require an iterative process. Thus there will be outer and inner iterations, and while both kinds of iteration have been studied extensively ( cf. Kantorovich and Akilov [1964], Varga [1962]) little is known about their combined behavior.

Bellman and Kalaba [1965] ( p. 118) and Ortega [1966] have pointed out that it is an open problem to decide how accurately to solve the linear equations at each Newton step in order to have an overall optimal method. In Section III we give a partial answer to this question. We prove there a theorem similar to Mysovskii's theorem ( cf. Kantorovich and Akilov [1964]) on Newton's method which gives a constructive (and quite simple) procedure for interrupting the inner iterations while still preserving the quadratic convergence of the outer iteration. In the Appendix we give some numerical results using this procedure.

While it is always interesting to have means for accelerating the convergence of an approximate method there are problems for which this is not only interesting but essential. This is the case with boundary value problems for partial differential equations. Many of the relevant points on the practical application of finite differences and acceleration techniques to linear partial differential equations have been made in an excellent paper by Fox [1950][*]. It is only in this and other papers by Fox and his collaborators in which it is possible to find any significant information at all about the numerical performance of acceleration procedures in multidimensional problems. The best

---

[*] We are grateful to Professor L. Fox for calling our attention to this reference.

reference to this work is Fox [1962]. Most of Fox's qualitative comments apply without essential changes to the nonlinear case we consider in section IV. There is, of course, the added complication of having to solve nonlinear difference equations and we provide also a quantitative (asymptotic) analysis.

In section IV we discuss in detail the application of the different acceleration techniques to the solution of mildly nonlinear elliptic equations. The numerical solution of this type of problem by means of finite differences has been studied by Bers [1953], and more recently by Parter [1965], and Greenspan and Parter [1965]. Other numerical results are reported in Greenspan [1964] and Greenspan [1965].

In section V we present some of the numerical results obtained in solving mildly nonlinear elliptic equations. These examples serve two purposes. On one hand they show the actual performance of the methods as applied to non-trivial problems. On the other hand they provide a way for comparing the two classes of techniques, and some comments on this are offered in section VI.

We are indebted to Professor Donald Greenspan and Professor Colin Cryer of the Mathematics Research Center for their continuous interest, support, and valuable suggestions.

# I. Acceleration techniques for nonlinear problems.

1. In this paper we will be dealing with the numerical solution of nonlinear operator equations by means of discretization algorithms. In particular we will be interested in the discussion of techniques for accelerating the convergence of such algorithms as the mesh size approaches zero.

2. A frequently used example of this family of problems is the two point boundary value problem. We treat this problem in detail in order to make the sense of the generalization to operator equations more comprehensible. Let

$$y'' - f(x, y) = 0 \qquad a \leq x \leq b$$

$$y(a) - \alpha = 0 \qquad\qquad\qquad\qquad (1.1)$$

$$y(b) - \beta = 0 \qquad ,$$

where $f_y > 0$.

We assume that $f$ is as smooth as is necessary to ensure the validity of all our expansions. In this case the continuous problem (1.1) has a unique solution $y(x)$ (cf. Henrici [1962]). In operator notation we can write (1.1) as:

$$F(y) = 0 \qquad\qquad\qquad (1.1')$$

where $F:D \rightarrow E$, and $D, E$ are contained in the Banach spaces $C^{(2)}[a,b]$ and $C[a,b] \times R^2$, respectively.

3. A very simple discrete version of (1.1), which allows us to obtain approximate values of $y(x)$ at points $x_i$, is given by:

$$h^{-2} \delta^2 Y_i - f(x_i, Y_i) = 0 \qquad 1 \leq i \leq n-1$$

$$Y_0 - \alpha = 0 ,$$

$$Y_n - \beta = 0 , \qquad\qquad\qquad (1.2)$$

where

$$x_i = a + ih , \quad (i = 0, 1, \ldots, n) ,$$

$$h = (b - a)/n ,$$

and where

$$\delta^2 Y_i = Y_{i+1} - 2Y_i + Y_{i-1} .$$

This is a system of $n - 1$ nonlinear equations which has to be solved for the $Y_i$. If a solution exists, we expect $Y_i - y(x_i)$ to be small in some sense. In operator notation we can write (1.2) as

$$\Phi_h(Y) = 0 \tag{1.2'}$$

where $\Phi_h : D_h \to E_h$, $D_h = E_h = R^{n+1}$.

4. For any sufficiently differentiable function $u(x)$ it is easy to show that the expansion

$$h^{-2} \delta^2 u(x_i) - f(x_i, u(x_i)) = u''(x_i) - f(x_i, u(x_i)) +$$

$$+ \sum_{j=1}^{N} \frac{2}{(2j+2)!} u^{(2j+2)}(x_i) h^{2j} + O(h^{2N+2}) \tag{1.3}$$

is valid. This asymptotic expansion in powers of $h$ shows the relationship between the continuous operator in (1.1) and the discrete operator in (1.2). The expression $h^{-2} \delta^2 u(x_i) - u''(x_i)$ is sometimes called the local discretization error or truncation error.

For the development of our theory, the existence of such a relationship is essential. In order to formalize this statement we have to introduce some new operators. These are the linear, bounded operators $\Delta_h$, $\Delta_h^0$ which relate the continous and discrete spaces:

$$
\begin{array}{ccc}
D & \xrightarrow{\ F\ } & E \\
\Delta_h \downarrow & \quad \Phi_h \quad & \downarrow \Delta_h^0 \\
D_h & \xrightarrow{\ \Phi_h\ } & E_h
\end{array}
.
$$

Then the generalization of (1.3) reads:

For each $y \in D$

$$\Phi_h(\Delta_h y) = \Delta_h^0 \{ F(y) + \sum_{j=1}^{N} h^{p_j} F_j(y) \} + O(h^{p_{N+1}}) , \qquad (1.3')$$

where the $F_j$ are given operators, and the $p_j$ are rational, positive numbers satisfying

$$0 < p_1 < p_2 < \ldots < p_{N+1} .$$

$F$ and $\Phi_h$ will be assumed to be at least twice Fréchet differentiable, and in general, the spaces $E, D, E_h, D_h$ can be arbitrary Banach spaces.

5. Let us assume for a moment that (1.1') and (1.2') have unique solutions $y, Y(h)$. In this case we can define

$$e(h) = Y(h) - \Delta_h y \qquad (1.4)$$

the global discretization error.

Under fairly general conditions Stetter [1965] has proved that if (1.3') holds (with $F_j$ independent of $h$) then also

$$e(h) = \Delta_h \sum_{j=1}^{N} e_j h^{p_j} + O(h^{p_{N+1}}) \qquad (1.5)$$

where the $e_j \in D$, and are independent of $h$. In particular this holds for problem (1.1) with the discretization (1.2).

One of the main assumptions in Stetter's theorem is that the operator $\Phi_h$ be stable. By this is meant that for any $e, V \in D_h$ there exists a constant $K$, which is independent of $e$ and $h$, such that

$$\| e \| \leq K \| \Phi_h'(V) e \| , \qquad (1.6)$$

where $\Phi_h'$ is the Fréchet derivative of $\Phi_h$.

6. The $e_j$ depend, in general, on the exact solution of (1.1'), which is of course unknown. Nevertheless, the knowledge that such an expansion exists is enough to allow us, in principle, to improve upon the basic approximate solution $Y(h)$.

We say that the method (1.2') is <u>convergent</u> for $h \to 0$ if

$$e(h) = o(1) \quad ,$$

and that it is <u>convergent of order $p$</u> if

$$e(h) = O(h^p) \quad . \tag{1.7}$$

It is well known that (1.2) is convergent of order 2. We will now seek methods for increasing the order of convergence, starting from the basic method (1.2') and assuming the properties (1.3'), (1.5) and (1.7), with $p$ equal to the $p_1$ of (1.3'). We will assume that $\Phi_h$ has the <u>mean value</u> property, i.e. that for each $V_1, V_2 \in D_h$, there exists a linear operator $M(V_1, V_2)$ such that

$$\Phi_h(V_1) - \Phi_h(V_2) = M(V_1, V_2)(V_1 - V_2)$$

and

$$M(V_1, V_2) - \Phi_h'(V) = o(1) \quad \text{for} \quad V_1, V_2 \to V \ .$$

Of course, in general, $M$ does not have to be unique. Finally we also assume that there is an $M$ of this kind which is stable.

7. A well known technique that goes back, at least, to a paper by Richardson [1910] is that of extrapolation to the limit.

The idea is that if we know the error of (1.2') to behave like (1.5) then, by computing several approximate solutions for different steps $h$, and combining them appropriately we can obtain a much more accurate solution.

For instance one could compute approximate solutions of (1.2) with $h = h_1, 2^{-1}h_1, \ldots, 2^{-N}h_1$. These particular grid sizes are such that the

grid points of the coarsest mesh ($h = h_1$) belong to all the other grids. It is

at these points that we will be able to obtain a more precise solution.

The name of Romberg is usually associated with the scheme used in

performing successive extrapolations. To be precise, in 1955 Romberg

presented a scheme for the numerical integration of continuous functions,

having as a basic method the trapezoidal rule and improving accuracy by

successive extrapolations. There is no conceptual change in the more

general situation occurring in the present context. Among the authors who

have recently contributed to extend the domain of applicability of this method

we can mention the names of Bauer, Rutishauser, Stiefel, Gragg, Bulirsch,

Stoer, and Laurent. Reference to their work can be found in Gragg [1965].

The only attempt to set up a general theory is due to Stetter, as we have

mentioned before. However, even in Stetter's paper it is left implicit how

the general successive extrapolation procedure is to be applied to functional

equations. Consequently, we feel justified in using some space in order

to state precisely <u>the method of successive extrapolations for functional</u>

<u>equations</u> (S. E. for short).

8.$^{(*)}$ To do so we assume that we are solving approximately equation

(1.1') by means of (1.2'), that (1.1') has a unique solution, and that (1.2')

has a unique solution that we can compute for any $h$ in discussion.

For any given $N$ and $h_1 > h_2 > \ldots > h_N > 0$, we assume that there exists

---

* A summary of the results of § 8 through § 11 and some of the numerical results
at the end of this paper have been presented to the SIAM National Meeting
of May 1966 in Iowa City.

-8-                                                                         #687

a family of operators $\psi_{h_j}$ $(j = 1, \ldots, N)$ $\psi_{h_j} : D_{h_j} \to D_{h_1}$, such that:

$$\text{for each } V \in D_{h_j}, \quad \psi_{h_j}(V) = \Delta_{h_1}(v) \tag{1.8}$$

where $v$ is any element of $\Delta_{h_j}^{-1}(V)$, and $\| \psi_{h_j} \| \leq 1$. These operators will be well defined if

for each $v_1, v_2 \in D$, and $j = 1, \ldots, N$,

$$\Delta_{h_j} v_1 = \Delta_{h_j} v_2 \text{ implies } \Delta_{h_1} v_1 = \Delta_{h_1} v_2 .$$

If the $D_{h_j}$ are finite dimensional subspaces of $D$ such that $\overline{\bigcup_j D_{h_j}} = D$, and the $\Delta_{h_j}$ are projections, then by taking $\psi_{h_j} = \Delta_{h_j}$ these properties are equivalent to saying that the $\psi_{h_j}$ form a Schauder basis for $D$.$^{(*)}$

If for problem (1.1), (1.2) we choose $\psi_{h_j}$ as the linear operator which maps $V \equiv (V_o, V_1, \ldots, V_{n_1 2^{j-1}})$ into the vector $W$ of $D_{h_1}$ defined by

$W \equiv (V_o, V_{2^{j-1}}, \ldots, V_{(n_1-1) 2^{j-1}}, V_{n_1 2^{j-1}})$, then we see that all the conditions are fulfilled.

9. Now we can easily prove the following

<u>Lemma 1.1</u> <u>If</u> $U(h_j)$ <u>is an approximate solution of order</u> $p_{\sigma+1}$ <u>and</u> $e(h)$ <u>has an expansion (1.5) up to the order</u> $p_\sigma$, <u>then</u>

$$\psi_{h_j} U(h_j) - \Delta_{h_1} u - \Delta_{h_1} \sum_{i=1}^{\sigma} h_j^{p_i} e_i = O(h_j^{p_{\sigma+1}}) . \tag{1.9}$$

* A sequence of finite dimensional subspaces $D_n$ of a Banach space $D$ and projections $P_n : D \to D_n$ is called a Schauder basis for $D$ if $D_n \subset D_{n+1}$, $\overline{\bigcup_n D_n} = D$, $\| P_n \| \leq M$, and $P_m P_n = P_m$ for $n \geq m$.

## Proof:

From (1.8) it follows that

$$\Delta_{h_1} u = \psi_{h_j}(\Delta_{h_j} u)$$

$$\Delta_{h_1} \sum_{i=1}^{\sigma} h_j^{p_i} e_i = \psi_{h_j}(\Delta_{h_j} \sum_{i=1}^{\sigma} h_j^{p_i} e_i) \ . \tag{1.10}$$

Thus, the left hand side of (1.9) can be written as:

$$\psi_{h_j}[U(h_j) - \Delta_{h_j} u - \Delta_{h_j} \sum_{i=1}^{\sigma} h_j^{p_i} e_i]$$

and because of our hypotheses this is in norm

$$\leq \| e(h_j) - \Delta_{h_j} \sum_{i=1}^{\sigma} h_j^{p_i} e_i \| \leq C_{\sigma+1} h_j^{p_{\sigma+1}}$$

as we wanted to prove.

10. Now we can state the S.E. method precisely. Let us first define

$$U_i^{(0)} = \psi_{h_i} U(h_i) \quad (i = 1, 2, \ldots, N)$$

$$T_i^{(0)} = \Delta_{h_1} u - U_i^{(0)}, \qquad r_1 = h_{i+1}/h_i, \tag{1.11}$$

$$g_{\nu,i}^{(0)} = 1, \qquad \rho_{\nu,i}^{(0)} = 1, \qquad \tau_{\nu,i}^{(0)} = 1 \ (\nu, i = 1, \ldots, N).$$

By Lemma 1.1 we have that

$$T_i^{(0)} = \Delta_{h_1} \sum_{\nu=1}^{N} e_\nu g_{\nu,i}^{(0)} h_i^{p_\nu} + O(h_i^{p_{N+1}}) \ . \tag{1.12}$$

With this we can recursively define

$$U_i^{(k)} = \frac{\rho_{k,i-1}^{(k-1)} \, r_{i-1}^{p_k} \, U_{i-1}^{(k-1)} - U_i^{(k-1)}}{\rho_{k,i-1}^{(k-1)} \, r_{i-1}^{p_k} - 1} \qquad (i = k+1, \ldots, N), \qquad (1.13)$$

$$g_{\nu,i}^{(k)} = \frac{\tau_{\nu,i-1}^{(k-1)} \, r_{i-1}^{p_k - p_\nu} - 1}{\rho_{k,i-1}^{(k-1)} \, r_{i-1}^{p_k} - 1} \, g_{\nu,i}^{(k-1)} \qquad (i, \nu = k+1, \ldots, N), \qquad (1.14)$$

$$\rho_{\nu,i}^{(k)} = g_{\nu,i+1}^{(k)} / g_{\nu,i}^{(k)}, \qquad \tau_{\nu,i}^{(k)} = \rho_{k+1,i}^{(k)} / \rho_{\nu,i}^{(k)} \qquad (i, \nu = k+1, \ldots, N). \qquad (1.15)$$

The next theorem will show that an expansion similar to (1.12) is valid for

$T_i^{(k)} = \Delta_{h_1} u - U_i^{(k)}$, and that this expansion starts with $\nu = k+1$.

<u>Theorem 1.2:</u> <u>With the definitions and hypotheses above, and if</u> $N$ <u>and the</u>

<u>rates</u> $r_i$ <u>are given, then for</u> $h_i \downarrow 0$ <u>the discretization errors</u> $T_i^{(k)}$ <u>satisfy</u>

$$T_i^{(k)} = \Delta_{h_1} \sum_{\nu=k+1}^{N} e_\nu g_{\nu,i}^{(k)} h_i^{p_\nu} + O(h_1^{p_{N+1}})$$

$$(k = 0, \ldots, N-1) \qquad (1.16)$$

$$(i = k+1, \ldots, N)$$

<u>where the</u> $g_{\nu,i}^{(k)}$ <u>do not depend on</u> $h_1$.

<u>Proof:</u> The proof is by induction on $k$. For $i = 2, 3, \ldots, N$, we have

$$T_i^{(1)} = \Delta_{h_1} u - U_i^{(1)} = \frac{1}{\rho_{1,i-1}^{(0)} \, r_{i-1}^{p_1} - 1} \left( \rho_{1,i-1}^{(0)} \, r_{i-1}^{p_1} T_{i-1}^{(0)} - T_i^{(0)} \right) =$$

$$= \Delta_{h_1} \sum_{\nu=1}^{N} \frac{e_\nu}{\rho_{1,i-1}^{(0)} \, r_{i-1}^{p_1} - 1} \left( \rho_{1,i-1}^{(0)} \, r_{i-1}^{p_1} g_{\nu,i-1}^{(0)} \, h_{i-1}^{p_\nu} - g_{\nu,i}^{(0)} \, h_i^{p_\nu} \right) + O(h_1^{p_{N+1}}).$$

$$(1.17)$$

For $\nu = 1$, the term in parentheses in (1.17) vanishes and for $\nu = 2, \ldots, N$

we have that every term becomes equal to

$$e_\nu \frac{\tau_{\nu, i-1}^{(0)} r_{i-1}^{p_1-p_\nu} - 1}{\rho_{1, i-1}^{(0)} r_{i-1}^{p_1} - 1} g_{\nu, i}^{(0)} h_i^{p_\nu} = e_\nu g_{\nu, i}^{(1)} h_i^{p_\nu} \quad ,$$

and (1.16) is proved for $k = 1$.

If we assume (1.16) to be true for $k < N-1$, a completely analogous

procedure permits us to pass to $k + 1$ and the induction argument is completed.

Corollary 1.3:   Under the hypotheses of Theorem 1.2 we have that after $k$

extrapolations the discretization error becomes

$$\Delta_{h_1} u - U_i^{(k)} = \Delta_{h_1} e_{k+1} g_{k+1, i}^{(k)} h_i^{p_{k+1}} + O(h_1^{p_{k+2}}) \qquad k = 0, \ldots, N-1$$
$$i = k, \ldots, N \qquad (1.18)$$

with

$$g_{k+1, i}^{(k)} = \prod_{j=0}^{k-1} \frac{\tau_{k+1, i-1}^{(j)} r_{i-1}^{p_{j+1}-p_{k+1}} - 1}{\rho_{j+1, i-1}^{(j)} r_{i-1}^{p_{j+1}} - 1} \quad . \qquad (1.19)$$

Proof.  Formula (1.18) is the same as (1.16).

From (1.14), for $\nu = k+1$ we obtain

$$g_{k+1, i}^{(k)} = \frac{\tau_{k+1, i-1}^{(k-1)} r_{i-1}^{p_k-p_{k+1}} - 1}{\rho_{k, i-1}^{(k-1)} r_{i-1}^{p_k} - 1} g_{k+1, i}^{(k-1)} = \ldots = \prod_{j=0}^{k-1} \frac{\tau_{k+1, i-1}^{(j)} r_{i-1}^{p_{j+1}-p_{k+1}} - 1}{\rho_{j+1, i-1}^{(j)} r_{i-1}^{p_{j+1}} - 1}$$

which proves (1.19).

11.  In certain very important special cases, (1.19) can be written more

explicitly.  For instance, if

$$r_i \equiv r \qquad (i = 1, \ldots, N)$$

$$(1.20)$$

$$p_\nu = \nu p_1 \qquad (\nu = 1, \ldots, N)$$

then the $g_{\nu, i}^{(k)}$ are independent of $i$ and consequently $\rho_{\nu, i}^{(k)} \equiv \tau_{\nu, i}^{(k)} \equiv 1$

for all values of the indices. Thus

$$g_{k+1, i}^{(k)} = \prod_{j=0}^{k-1} \frac{r^{p_{j+1} - p_{k+1}} - 1}{r^{p_{j+1}} - 1} = (-1)^k \, r^{\sum_{j=0}^{k-1} (j-k) p_1} =$$

$$= (-1)^k \, r^{-\frac{1}{2} k(k+1) p_1} = (-1)^k \, r^{-\frac{1}{2}(k+1) p_k} . \qquad (1.21)$$

The results of §10 and §11 are related to Theorem 1 in Bulirsch and Stoer [1964]; Theorem 1.2 contains one of the cases treated by Bulirsch and Stoer, while the results of §11 are contained in their other case. Of course, the discussion of §8 and §9, which allows the basic solutions to be in different linear spaces is new. The general algorithm of §10 is also new. W. Gragg has communicated to the author that he has a similar algorithm (unpublished).

II. <u>Deferred corrections</u>

1. A completely different technique for accelerating the convergence of discretization algorithms is based on an idea of Fox [1947], [1950], [1962]. The corresponding general procedure has been developed in Pereyra [1966][*]. Earlier contributors have been Volkov [1957], [1963], [1965], Bickley, Michaelson, and Osborne [1960], Henrici [1962], Lees [1966], and Pereyra [1965].

Let us discuss the basic idea when applied to the problem (1.1). Since

---

* From now on this paper will be referred to as P.

we know the expansion (1.3), it seems appropriate to use this information in order to improve our approximate solution. A way of doing this is to replace the higher derivatives of the solution by sufficiently accurate difference approximations. In this fashion we will obtain a more complicated basic method for which, in more general situations, we may not have sufficient theoretical results. Another possibility, and this is Fox's idea, is to first solve (1.2) and then to feed back this approximate solution into an appropriate difference approximation to the right hand side of (1.3). Hopefully, solving (1.2) with this new right hand side will yield an improved solution.

Passing now to our general problem (1.2') with the expansion (1.3'), we will discuss a method of deferred corrections which will allow us to obtain from an $h^{p_1}$-solution an $h^{\bar{p}}$ solution with $\bar{p} = p_1 + \min(p_1, p_2 - p_1)$.

2. The procedure we are going to present will be called the <u>one-step linear deferred correction</u> for operator equations (L.D.C.).

We assume that there exists an operator $S$ such that

$$\Delta_h^0 F_1 u - S(U) = O(h^{p*}) \qquad (2.1)$$

where $U$ is a $p_1$-approximate solution of (1.1') and $p* = \min(p_1, p_2 - p_1)$. Then

$$U_1 = U - h^{p_1} e* \qquad (2.2)$$

is an approximate solution of (1.1') of order $\bar{p} = \min(2p_1, p_2)$. Here $e*$ is the solution of

$$\Phi_h'(U) e* = -S(U) . \qquad (2.3)$$

For the proof of this statement see P. , Theorem 3.1.

In problem ( 1.1) we can take

$$[S(Y)]_i = -\frac{h^{-2}}{12} \delta^2 f(x_i, Y_i) \tag{2.4}$$

which has all the desired properties.

Observe that in order to obtain this improved solution we had only to solve the linear problem ( 2.3), which is of the same "size" ( same h) as the original one.

It is also clear that if Newton's method was used for the solution of (1.2') then ( 2.3) has the same structure as a Newton iteration. All these remarks shows that it is rather economical to gain $p*$ extra orders of accuracy by means of this method, whenever the computation of $S(U)$ is not very cumbersome.

An iterative deferred correction (I.D.C.) procedure has been also described in P but since its application presents several new problems we prefer to discuss it in a separate publication.

III. Incomplete nested iterations.

1. In solving problem ( 1.1') by means of ( 1.2') we typically find several nested sequences, which are generally generated by means of iteration procedures. It has been observed$^{(*)}$ that it is not always necessary to carry out the inner iterations to completion.

In this section we give a sufficient criterion for stopping the inner iterations without perturbing excessively the final results.

2. We now describe the different iterations.

---

*Cf. Douglas [1961], Henrici [1962], Pereyra [1965], [1966], and Ortega and Rockoff [1965].

(a) The first level iteration, $I_1$, corresponds to solving

$$\Phi_h(U) = 0 \qquad\qquad (3.1)$$

for a sequence of decreasing $h$. We assume that (3.1) has a unique solution $U(h)$ and that

$$U(h) - \Delta_h u = O(h^p), \qquad p > 0 .$$

(b) The second level iteration, $I_2$, corresponds to the solution of the non-linear problem (3.1) for each $h$. We concentrate on linearization methods and more specifically on Newton's method, i.e.: for a given initial value $V_o$, $I_2$ is the iteration

$$\Phi_h'(V_i)(V_{i+1} - V_i) = -\Phi_h(V_i) . \qquad\qquad (3.2)$$

We assume that this iteration is defined and that the $V_i$ converge to $U(h)$ as $i \to \infty$.

(c) In many problems the linear equation (3.2) will also have to be solved by an iterative procedure, $I_3$. For instance, if (3.2) is a large system of linear equations some of the standard iteration techniques may be used, generating a sequence $W_{i+1}^j$ which we assume, for $j \to \infty$, convergent to $W_{i+1} = V_{i+1} - V_i$, the exact solution of (3.2).

3. For an equation $G(u) = 0$ we define its <u>residual</u> for $v \in D$, as

$$r = \|G(v)\| . \qquad\qquad (3.3)$$

If $u$ is a solution of $G(u) = 0$ then $r = 0$. For a given $h$ we would like to control the iterations $I_2$ and $I_3$ by inspecting only the residuals of equations (1.1') and (3.2) for the successive iterates.

In $P$ we have already studied the effects of prematurely stopping the

iteration $I_2$. There we showed that, if $\Phi_h$ is stable, in order to obtain a $\widetilde{U}(h)$ such that $\widetilde{U}(h) - \Delta_h u = O(h^p)$ it is sufficient that $\Phi_h(\widetilde{U}(h)) = O(h^p)$, regardless of how $\widetilde{U}(h)$ has been obtained.

In anticipation of future needs we will ask that the residual of $\Phi_h(\widetilde{U}(h))$ be of the order $h^q$, with $q \geq p$. This will ensure that the error $U(h) - \widetilde{U}(h)$ does not interfere with the asymptotic expansion of $e(h)$. In fact, if we know that the exact solution $U(h)$ of (3.1) satisfies

$$e(h) = U(h) - \Delta_h u = \sum_{\nu=1}^{N} e_\nu h^{p_\nu} + O(h^{p_{N+1}})$$

and that

$$\Phi_h(\widetilde{U}(h)) - \Phi_h(U(h)) = O(h^q) \tag{3.4}$$

then the mean value property and (3.4) imply that

$$M(\widetilde{U}(h), U(h))(\widetilde{U}(h) - U(h)) = O(h^q)$$

and the stability of $M$ proves that

$$\widetilde{U}(h) - U(h) = O(h^q) \quad . \tag{3.5}$$

From the asymptotic expansion for $e(h)$ and (3.5) we obtain

$$\widetilde{U}(h) - \Delta_h u = \sum_{\nu=1}^{N} e_\nu h^{p_\nu} + O(h^q) + O(h^{p_{N+1}}) \quad . \tag{3.6}$$

4. Recalling that all these results are independent of the way in which $\widetilde{U}(h)$ has been obtained we will now seek a stopping procedure for $I_3$. To do this we will study the behavior of the sequence $\{V_i\}$ generated by

$$\Phi_h'(V_i)(V_{i+1} - V_i) = -\Phi_h(V_i) + \mathcal{E}_i \tag{3.7}$$

where the residual $E_i = \| \mathcal{E}_i \|$ satisfies certain conditions. The precise results are contained in the next theorem which is a slight variation of Mysovskikh's theorem ( cf. Kantorovich and Akilov [1964], Chap. XVIII, p. 717).

<u>Theorem 3.1</u>   <u>Let</u> f <u>be a twice continuously differentiable operator mapping the Banach space</u> D <u>into the B-space</u> E. <u>Let us assume that for given</u> $x_0 \epsilon$ D, $\mathcal{E}_0 \epsilon$ E, <u>and a sphere</u> $\Omega_0 = \{x: \| x - x_0 \| \leqq r \}$ <u>there exist constants</u> B, $\eta_0$, <u>and</u> $\alpha < \frac{1}{2}$ <u>such that the following conditions are satisfied:</u> $[f'(x)]^{-1}$ <u>exists on</u> $\Omega_0$ <u>and</u>

$$\| [f'(x)]^{-1} \| \leqq B \qquad , \quad x \epsilon \Omega_0, \; 0 < \frac{\alpha}{\sqrt{2(1-\alpha)} - 1} \leqq B \; , \qquad (3.8)$$

$$\| f(x_0) \| \leq \eta_0 \qquad , \qquad (3.9)$$

$$\| f''(x) \| \leqq K \qquad , \quad x \epsilon \Omega_0 \; , \qquad (3.10)$$

$$H_0 = B^2 \eta_0 K < 1 \qquad , \qquad (3.11)$$

$$\| \mathcal{E}_0 \| \leqq E_0 = \alpha \eta_0 H_0 ' \qquad , \qquad (3.12)$$

and

$$r' = \sqrt{2(1-\alpha)} \; B \eta_0 \sum_{k=0}^{\infty} H_0^{2^k - 1} < r \; . \qquad (3.13)$$

<u>Then the iteration</u>

$$f'(x_i)(x_{i+1} - x_i) = -f(x_i) + \mathcal{E}_i \qquad (3.14)$$

<u>is defined for all</u> $i \geqq 0$. <u>Here</u> $\mathcal{E}_i$ <u>is such that</u>

$$\| \mathcal{E}_i \| \leqq E_i = \alpha H_i \eta_i \qquad (3.15)$$

<u>and</u> $\eta_i$, $H_i$ <u>are defined recursively by</u>

$$\eta_i = H_{i-1} \eta_{i-1} \qquad , \qquad (3.16)$$

$$H_i = H_{i-1}^2 \qquad .$$

<u>For any choice of</u> $\mathcal{E}_i$ <u>satisfying</u> ( 3.15) <u>the sequence</u> $\{x_i\}$ <u>converges to a</u>

<u>solution</u> $x^* \epsilon \Omega_o$ <u>of</u> $f(x) = 0$. <u>The speed of convergence is given by</u>

$$\|x^* - x_i\| \leq B\eta_o \frac{H_o^{2^i - 1}}{1 - H_o^{2^i}} \sqrt{2(1-\alpha)} \quad . \tag{3.17}$$

<u>Proof.</u>  First of all we have that

$$\|x_1 - x_o\| \leq \|[f'(x_o)]^{-1}\| \; \|f(x_o)\| + E_o \leq B\eta_o (1 + \alpha H_o / B) \leq B\eta_o \sqrt{2(1-\alpha)}$$

and thus $x_1 \epsilon \Omega_o$ .

Since

$$f(x_1) = f(x_1) - f(x_o) - f'(x_o)(x_1 - x_o) + \mathcal{E}_o$$

we obtain

$$\|f(x_1)\| \leq \tfrac{1}{2} K \|x_1 - x_o\|^2 + E_o \quad ,$$

and thus

$$\|f(x_1)\| \leq KB^2 \eta_o^2 (1-\alpha) + \alpha \eta_o H_o = H_o \eta_o = \eta_1 \quad .$$

Now if we define

$$H_1 = B^2 \eta_1 K = B^2 \eta_o KH_o = H_o^2 < 1$$

and

$$E_1 = \alpha H_1 \eta_1$$

then we have reproduced the conditions of the theorem for the index $i = 1$.

Suppose now that ( 3.16) and the estimates for $\|x_i - x_{i-1}\|$ and $\|f(x_i)\|$

are valid for $i = 1, \ldots, k$.

We have as before

$$\|x_{k+1} - x_k\| \leq \eta_k B + E_k = B\sqrt{2(1-\alpha)} \; \eta_k .$$

But

$$\eta_k = H_{k-1}\,\eta_{k-1} = H_{k-1}H_{k-2}\,\eta_{k-2} = \cdots = \prod_{j=0}^{k-1} H_j\,\eta_o \ ,$$

and

$$H_j = H_{j-1}^2 = H_{j-2}^{2^2} = \cdots = H_o^{2^j} \ .$$

Therefore

$$\eta_k = \eta_o H_o^{2^k - 1} \ ,$$

and

$$\|x_{k+1} - x_k\| \leq \sqrt{2(1-\alpha)}\ B\eta_o H_o^{2^k - 1} \ .$$

Writing this inequality for all the indeces between 0 and $k$, and adding them up we get:

$$\|x_{k+1} - x_o\| \leq \sqrt{2(1-\alpha)}\ B\eta_o \sum_{j=0}^{k} H_o^{2^j - 1} < r \ ,$$

and hence $x_{k+1} \epsilon\ \Omega_o$.

Also, for any $p \geq 0$

$$\|x_{k+p} - x_k\| \leq \sqrt{2(1-\alpha)}\ B\eta_o \sum_{j=k}^{k+p} H_o^{2^j - 1} \ ,$$

which tends to zero for $k \to \infty$ since the series is convergent by hypothesis. Thus, $\{x_k\}$ is a Cauchy sequence and there exists $x^* \epsilon\ \Omega_o$ such that $x_k \to x^*$. It is clear that $x^*$ is a solution of $f(x) = 0$.

If we let $p \to \infty$ in the above inequality we obtain

$$\|x^* - x_k\| \leq \sqrt{2(1-\alpha)}\ B\eta_o H_o^{2^k - 1} \sum_{j=0}^{\infty} (H_o^{2^k})^j = \sqrt{2(1-\alpha)}\ B\eta_o \frac{H_o^{2^k - 1}}{1 - H_o^{2^k}}$$

which is ( 3.17) .

In order to apply this result to problem ( 3.1) we first observe that condition ( 3.8) is just the uniform stability of $\Phi_h$ in $\Omega_o$. The practical

procedure consists in finding the values of the initial constants B, $\alpha$, and K, and a suitable $V_o$ in order that (3.11) and (3.13) be fulfilled. Once this has been done we can compute $E_o$ and then start the iteration $I_3$ in order to solve (3.2) (for i = 0). The iteration $I_3$ will be stopped when

$$\| \mathcal{E}_o \| = \| \Phi_h'(V_o) W_1^j + \Phi_h(V_o) \| \leq E_o \qquad (3.18)$$

and then we will take as the next iterate

$$V_1 = V_o + \widetilde{W}_1 \quad ,$$

$\widetilde{W}_1$ being the last iterate in (3.18). With this we can continue the recursion computing successively

$$\eta_i = \eta_{i-1} H_{i-1} \quad ,$$

$$H_i = H_{i-1}^2 \quad ,$$

$$E_i = \alpha \eta_i H_i \quad , \qquad (3.19)$$

and a $\widetilde{W}_{i+1}$ from $I_3$ satisfying (3.18) for the subindex i. We will stop when

$$\eta_i \leq Ch^q \qquad (3.20)$$

where C is a small constant.

If we are interested in solving a general system of nonlinear equations it is worthwhile to observe that from (3.17) we obtain the error bound

$$\| x^* - x_i \| \leq B \eta_i \frac{1}{1-H_i} \sqrt{2(1-\alpha)} \qquad (3.21)$$

with very little extra computation.

## IV. Mildly nonlinear elliptic boundary value problems.

1. Generalization of problem ( 1.1) to two dimensions gives

$$\Delta u(x, y) = u_{xx} + u_{yy} = f(x, y, u) \qquad (x, y) \in \mathcal{D}$$

$$u(x, y) = g(x, y) \qquad (x, y) \in \partial \mathcal{D}$$

$$f_u(x, y, u) \geq 0 \qquad (x, y) \in \mathcal{D}, \ |u| < \infty ,$$

$$(4.1)$$

where $\mathcal{D}$ is a bounded, open, plane region and its boundary consists of continuous closed curves. As stated, problem ( 4.1) has a unique solution $u(x, y)$ which is in $C(\bar{\mathcal{D}}) \cap C^{(2)}(\mathcal{D})$ ( cf. Parter [1965]) ( *) .

The structure of solutions of ( 4.1) without restrictions on $f_u$ has been recently studied by Greenspan and Parter [1965], and Parter [1965]. Under more stringent regularity conditions on $f, g$ and the boundary curve $\partial \mathcal{D}$ it is possible to ensure, a priori, that $u \in C^{(p)}(\bar{\mathcal{D}})$, $p > 2$. For some interesting regions with piecewise analytic boundaries having special kinds of corners it is possible to subtract off singularities appearing in the solution and its derivatives ( at least in the linear case), thus obtaining an associated problem with a regular solution ( cf. Volkov [1963]). In this section we will be mainly concerned with the application of the highly accurate methods described in I and II to problems of the form ( 4.1) with solutions of class $C^{2N}(\bar{\mathcal{D}})$, $N \geq 2$.
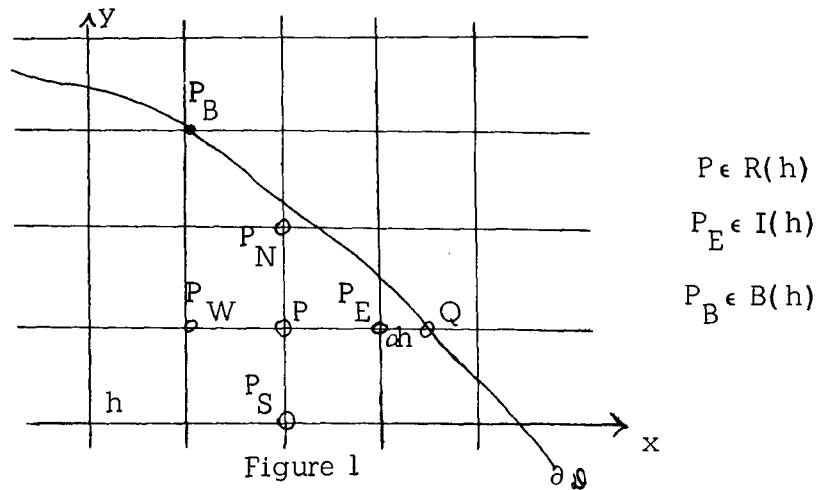
2. Let us first consider the case of a general, smooth boundary. As is customary, we take a square mesh of width $h$ covering $\bar{\mathcal{D}}$. Let

---

\* Since very little experimentation with these methods have been published, we emphasize that all the results of this section apply to the linear cases ( Laplace' s, Poisson' s equation, etc.).

$\mathcal{P} = ( P_1, \ldots, P_n)$ be the set of mesh points contained in $\bar{\mathcal{D}}$, ordered in some arbitrary but given manner. We subdivide $\mathcal{P}$ into three classes. The set of regular or interior points R( h) containing those mesh points whose four closest neighbors in the x and y directions belong to $\bar{\mathcal{D}}$; the set of irregular points I( h) containing those mesh points which are in $\mathcal{D}$ but not in R( h), and finally the set B( h) of boundary points, consisting of those mesh points belonging to $\partial\mathcal{D}$. Sometimes the neighbors of a point P in R( h) will be called $P_E$, $P_N$, $P_W$, and $P_S$ . See Figure 1 for an illustration.



P ∈ R( h)

$P_E$ ∈ I( h)

$P_B$ ∈ B( h)

Figure 1

Any function from $\mathcal{P}$ to R will be called a mesh function.

For any given mesh function $V( P) \equiv V_P$ we define the following finite difference operator

$$( \Phi_h ( V))_P = h^{-2} \{V( P_E) + V( P_N) + V( P_W) + V( P_S) - 4 V( P) \} - f( P, V( P))$$

$$P \in R( h) \quad , \qquad ( 4. 2)$$

$$( \Phi_h ( V))_P = g( P) \qquad P \in B( h) \quad .$$

For P ∈ I( h) we could choose either

$$( \Phi_h ( V))_P = g( Q) \qquad P \in I( h) \qquad ( 4. 2')$$

where $Q$ is the closest boundary point in the $x, y$ directions, or else the more accurate interpolation formula

$$( \Phi_h(V))_P = V(P) - \frac{\alpha(P)}{1+\alpha(P)} V(P_{NB}) - \frac{1}{1+\alpha(P)} g(Q),$$

$$P \in I(h) \qquad (4.2'')$$

where $Q$ is as before, $P_{NB}$ is the closest neighbor of $P$ on the same line as $Q$, and $\alpha(P) = \dfrac{\text{dist}(Q, P)}{h}$. For instance, if $P = P_E$ in Figure 1 then $P_{NB} = P$.

However, as was observed by Wasow [1955], these boundary approximations will not generally be good enough to provide an asymptotic expansion of the type (1.5), even if the exact solution is sufficiently differentiable. In Wasow's paper it was shown, on an unidimensional example, that, if the boundary was not treated carefully $e(h)$ was not differentiable with respect to $h$, and thus one could not expect an asymptotic expansion in powers of $h$. However, it was shown that if higher order interpolation were used at the boundary then one could recover the asymptotic expansion. From the point of view of Stetter's theorem, and always in the case of sufficiently smooth solutions, we see that the difficulty with the boundary stems from the values $\alpha(P)$, which depend on the mesh, appearing in the expansion (1.3') for the local discretization error. In fact, let $P_E \in I(h)$ and assume that (4.2'') is used with $P_W$ in place of $P_{NB}$, $Q \in \partial \Omega$ to the right of $P$, and $\alpha(P)$ as defined above (see Fig. 1; interpolation in other directions produces similar formulas). Then for any sufficiently differentiable function $v(x, y)$ we have

$$v(P) - \frac{\alpha(P)}{1+\alpha(P)} \ v(P_W) \ - \frac{1}{1+\alpha(P)} \ g(Q) = - \tfrac{1}{2}h^2 \alpha(P) v_{xx}(P) -$$

$$- \frac{h^3}{6} \alpha(P)(\alpha(P)-1) v_{xxx}(P) - \dots$$

and we see that for this part of the expansion (1.3') the operators $F_j(v)$ are <u>not independent of h</u> and Stetter's theorem does not apply.

3.  One way of avoiding this difficulty is to take higher order interpolation formulas, using only points belonging to $\bar{\aleph}$ . Ideally we would like to construct interpolation formulas using the closest neighbors of P on the plane mesh and perhaps some on the boundary $\partial \aleph$. However, except for the lowest orders, such as Mikeladzes' formula (cf. Panov[1963], p. 38), it is not simple to construct bivariate interpolation formulas for a general distribution of points and, when possible, it is difficult to assess the order of the discretization error. A feasible procedure is to use Newton's interpolation formula in one dimension, i.e., the generalization of (4.2") to higher order accuracy. Thus, to obtain an accuracy of order $h^{\bar{q}+1}$ we need to take P, Q and the $\bar{q}$ left, (say) neighbors of P, in order to form

$$(\Phi_{1h}(\Delta_h v))_P = v(P) - \prod_{k=1}^{\bar{q}} \frac{k}{\alpha+k} v(Q) - \sum_{k=1}^{\bar{q}} \binom{\bar{q}}{k} \frac{\alpha}{\alpha+k}(-1)^{k-1} v(P_{-k}) \quad (4.3)$$

which for $v \in C^{\bar{q}+1}(\bar{\aleph})$ has the residual

$$R_{\bar{q}}(\alpha) = \frac{\alpha h^{\bar{q}+1}}{\bar{q}+1} \ v_x^{(\bar{q}+1)}(\xi, y) \ , \quad\quad\quad (4.4)$$

where $v_x^{(\bar{q}+1)}$ indicates the $\bar{q}+1$ partial derivative of v in the direction of the interpolation.

Again we are faced with a difficulty. It is simple to see that for certain values of $\bar{q}$ and $\alpha$ we will have

$$\sum_{k=1}^{\bar{q}} \binom{\bar{q}}{k} \frac{\alpha}{\alpha+k} > 1 \quad , \tag{4.5}$$

and the matrix $\Phi_h'(V)$ will not be diagonally dominant. In Volkov [1957] it is proposed to take a multiple of the basic step $h$, $\bar{h}$, in order to make the sum in (4.5) less than 1, and it is shown that under certain geometrical conditions on the region $\mathcal{D}$ this can be done for all $h$ and $P \in I(h)$. The inconvenience of this procedure is that even for $\bar{q}$ small, say 6, and very good regions, like circles, it is necessary to take the step $h$ unreasonably small in order to have enough points to carry the process through.

In some of our numerical experiments we have used the interpolation formula (4.3) with the basic step $h$, disregarding (4.5). From the theoretical point of view we cannot, in these cases, apply the known results on convergence of $U(h)$ to $u$, etc. which depend on the diagonal dominance of $\Phi_h'$. On the other hand the numerical results seem to indicate that these properties still hold.

A modification of the difference correction, which resembles Fox's original method, gives quite good results even when the basic method uses the simple interpolation formula (4.3). This will be discussed in §6.

If the region $\mathcal{D}$ is rectangular, having grid lines for sides, then none of these difficulties appears since the set $I(h)$ is empty. Nevertheless, because the corners, in this case it is not easy to ensure a priori the sufficient differentiability of the solutions.

4.  Going back to (4.1) we see that the operator

$$F(u) \equiv \Delta u(x,y) - f(x,y,u) \qquad (x,y)\epsilon \mathcal{Q} \qquad\qquad (4.6)$$

maps $C(\bar{\mathcal{Q}}) \cap C^2(\mathcal{Q}) \cap \{u: u=g \text{ on } \partial\mathcal{Q}\}$ into $C(\mathcal{Q})$ . We take $D = C^{2N+3}(\bar{\mathcal{Q}})$ and $E = C^{2N+1}(\bar{\mathcal{Q}})$ . The discrete approximation has been already described in (4.2) and either (4.2'') or (4.3). Both $\Phi_h(V)$ (for (4.2'')) and $\Phi_{1h}(V)$ (for (4.3)) map the set $\mathcal{V}$ of mesh functions over $\mathcal{P}$ into itself. We take in $D_h = E_h = R^n$ the infinite norm

$$\|V(P)\| = \max_{j=1,\ldots,n} |V(P_j)| . \qquad\qquad (4.7)$$

Also $\Delta_h(u(x,y)) = \overset{o}{\Delta}_h(u(x,y)) = \{u(P)\}_{P\epsilon\mathcal{P}}$ . With these definitions we can write expansions of the form (1.3') for $\Phi_h$ and $\Phi_{1h}$ : for each $v\epsilon D$ ,

$$\Phi_h(\Delta_h v)(P) = \overset{o}{\Delta}_h\left\{ F(v)(P) + \sum_{j=1}^{N} h^{2j}\frac{(-2)}{(2j+2)!}(v_x^{(2j+2)}(P) + v_y^{(2j+2)}(P)) \right\}$$

$$+ O(h^{2N+1}) \qquad\qquad P\epsilon R(h) ; \qquad (4.8)$$

$$\Phi_h(\Delta_h v)(P) = -\frac{\alpha(P)}{1+\alpha(P)} \sum_{j=2}^{2N} h^j \frac{(\alpha(P)^{j-1}+(-1)^j)}{j!} v^{(j)}(P) + O(h^{2N+1})$$

$$P\epsilon I(h) \qquad\qquad (4.9)$$

where the derivatives are taken in the direction of the interpolation. Finally if we take $\bar{q} = 2N$ in (4.3) then

$$\Phi_{1h}(\Delta_h v)(P) = O(h^{2N+1}) , \qquad P\epsilon I(h) . \qquad (4.9')$$

By using (4.8), (4.9') and assuming that $\Phi_{1h}$ is stable and convergent of order 2 we can apply Stetter's theorem, obtaining an asymptotic expansion (1.5) for the discretization error. These assumptions are certainly satisfied

when $I(h)$ is empty, or when Volkov's procedure is used to obtain diagonal dominance in $\Phi_h'$.

In order to apply the procedure SE of Section I.10 we have to define the "projections" $\psi_{h_j}$ of I.8. Let us take a basic mesh size $h_1$ and define $h_j = 2^{-j}h_1$, i.e. $r = \frac{1}{2}$. With $h_j$ there is associated an $n_j$, the number of points in $\wp(h_j)$. Moreover, if $i > j$ then $\wp(h_i) \supset \wp(h_j)$ and in particular $\wp(h_1)$ is contained in all the refined meshes. Now we can define $\psi_{h_j}(V)$ for $V \epsilon D_{h_j}$, $V = (V_1, \ldots, V_{n_j})$ as

$$\psi_{h_j}(V) = W = (W_1, \ldots, W_{n_1}), \quad (j = 2, \ldots, N) \tag{4.10}$$

with $W_k = V(P_{s_k})$, $P_{s_k} \epsilon \wp(h_j)$ being the same point in the plane as $P_k \epsilon \wp(h_1)$. In other words $\psi_{h_j}(V)$ is the restriction of the mesh function $V$ to $D_{h_1}$. As an example consider the quarter circle in Fig. 2, $h_1 = \frac{1}{3}$, $h_2 = \frac{1}{6}$. There $n_1 = 11$, $n_2 = 35$. We denote by $P_i$ the points in $\wp(h_1)$ and $Q_j$ the ones in $\wp(h_2)$. Also $R(h_1) = \{P_6\}$, $I(h_1) = \{P_7, P_9, P_{10}\}$. If $U(h_i)$ satisfies

$$\Phi_{1h_i}(V) = O(h^{2N+1}) \tag{4.11}$$

(see III.3) then, as in I.10, we can define

$$U_i^{(0)} = \psi_{h_i} U(h_i) \quad (i = 1, 2, \ldots, N+1) \tag{4.12}$$

$$U_i^{(k)} = \frac{4^k U_i^{(k-1)} - U_{i-1}^{(k-1)}}{4^k - 1} \quad (k = 1, \ldots, N; \; i = k+1, \ldots, N+1) \tag{4.13}$$

and the result of I.11 applies. Recall that all $U_i^{(k)} \epsilon D_{h_1}$, i.e.: we obtain

-28-

more accurate approximations only at points belonging to the coarsest mesh.

The solutions of $\Phi_h(V) = O(h^q)$ will be obtained by a combination of

Newton's method and point successive overrelaxation, using the results

of III as the stopping procedure.

The Newton iteration is given in (3.2) and it will be stopped when

$$\| \Phi_h(V_i) \| \leq C h^q \quad .$$

On the other hand, if we consider the matrix $\Phi'_h(V_i) = D_i - E_i - F_i$,

where $D_i$ is diagonal, and $E_i$ and $F_i$ are respectively strictly lower and

upper triangular then, (3.2) can be approximately solved for $W_{i+1} = V_{i+1} - V_i$

by means of the iteration formula $(I_3)$

$$W_{i+1}^{(j+1)} = (I - \omega_i L_i)^{-1} \{ (1 - \omega_i) I + \omega_i U_i \} W_{i+1}^{(j)} +$$

$$+ \omega_i (I - \omega_i L_i)^{-1} D_i^{-1} \Phi_h(V_i), \qquad (4.14)$$

where $1 < \omega_i < 2$, and $L_i = D_i^{-1} E_i$, $U_i = D_i^{-1} F_i$.

As we said in III (see Theorem 3.1), $I_3$ will be interrupted when

$$E_j = \| \Phi'_h(V_i) W_{i+1}^{(j+1)} + \Phi_h(V_i) \| \leq \alpha H_i \eta_i \quad . \quad . \qquad (4.15)$$

It is well known that a bound $B$ for $\| [\Phi'_h(V_i)]^{-1} \|$ is

$$\| [\Phi'_h(V_i)]^{-1} \| \leq \frac{d^2}{16} = B , \qquad (4.16)$$

where $d$ is the diameter of the set $\mathcal{D}$. This bound $B$, together with $K$,

a uniform bound for $f_{uu}$, provides the necessary data to compute $H_o$
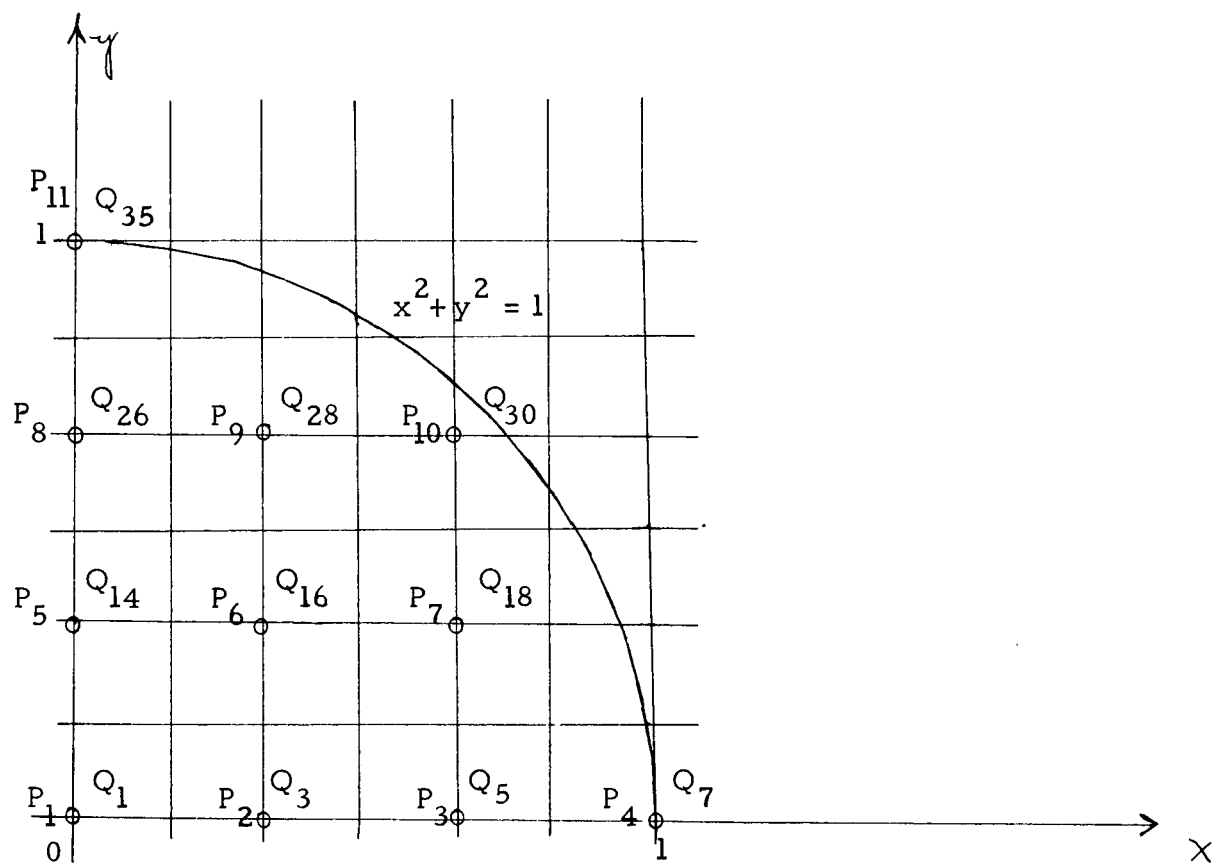
(see Theorem 3.1).

Figure 2

5. For the linear deferred correction L.D.C. described in II. 2 we need to determine the operator $S(U)$ of ( 2.1). Here $p^*$ will be 2 and

$$F_1(u) = -\frac{1}{12}(u_x^{(4)} + u_y^{(4)}) \quad (\text{assuming that we have used } \Phi_{1h} \text{ with } \bar{q} = 5).$$

We take $S(U)$ as the appropriate combination of symmetric fourth order differences of $U$ at all points at which this is possible, and unsymmetric formulas at the remaining points. For instance, in Figure 2, and for $h = \frac{1}{16}$ , we have that at $Q_{16}$ we can use symmetric fourth differences in every coordinate direction while at $Q_{30}$ we cannot. However, six point unsymmetric formulas, using the neighbors to the left and below can be applied. For these and other necessary difference approximations we refer to Ballester and Pereyra [1966].

6. As we mentioned before it is also possible to use the boundary interpolation ( 4.2'') if the L.D.C. is reiterated. In this case, the part of $F_1(v)$ corresponding to points $P \in I(h)$ becomes ( cf. ( 4.9)) :

$$F_1( v) ( P) = -\tfrac{1}{2} \alpha( P) ( v''( P) - \frac{h}{3} ( 1-\alpha( P)) v'''( P)) \qquad ( 4.17)$$

where the differentiation is in the direction of the interpolation. Observe that since $F_1(v)$ is not independent of $h$ the theory of II. 4 no longer applies. We can, of course, give an $O(h^2)$ discrete approximation to $F_1( v)$ for any function $v \in C^4$, namely,

$$S(\Delta_h v)_0 = -\frac{\alpha}{12} h^{-2} [ ( 7+5\alpha) v_0 - ( 12 +18\alpha) v_{-1} + 24\alpha v_{-2} + 2( 4-7\alpha) v_{-3}$$

$$+ 3( \alpha-1) v_{-4}] \; . \qquad ( 4.18)$$

Using ( 4.18) as $S( U)$, taking $U^{(0)}$ as the solution of $\Phi_h( V) = O(h^4)$ and

iterating according to the formulas

$$\Phi_h'(U^{(i)}) e^{(i)} = -S(U^{(i)}) \quad (i = 0, 1, \ldots) \qquad (4.19)$$

$$U^{(i+1)} = U^{(i)} - h^2 e^{(i)}$$

we have obtained, after a few steps, results equivalent to the ones obtained by means of the algorithm of § 5.

## V. Numerical examples.

1. In what follows we describe some of the numerical experiments we have carried out on the CDC 3600 Computer at the University of Wisconsin Computing Center. They refer to the problem discussed in Section IV and the methods of Sections I and II. Numerical results for two-point boundary value problems solved by means of L D.C. and LD.C. can be found in Pereyra [1965] and [1966].

The basic program solves a large system of nonlinear algebraic equations by Newton's method. At every Newton step the resulting system of linear equations is solved by point overrelaxation. No great effort has been made in order to find an optimal value for the overrelaxation parameter $\omega$. The comparison of the various methods is sensible only when the same $h$ and $\omega$ are involved (cf. Ortega and Rockoff [1965]).

For SE we decide a priori how many extrapolations we want to make and from the given basic step $h_1$ we deduce the smallest $h$ needed. We solve first for that $h$ using as an initial approximation a linear interpolation of the boundary values. For the larger $h$'s we use as starting values the solutions just obtained, considered at the relevant grid points. It may be

of interest to proceed in the reverse direction, beginning with the coarsest mesh, and using that solution and interpolating to fill the gaps, in order to start the iteration for the next mesh. Another possible modification is to use a different type of refinement, i.e. some sequence decreasing slower than $\left(\frac{1}{2}\right)^i$ .

2. <u>Problem 1</u>　Let $\bar{\Omega}$ be the quarter of the unit circle $x^2 + y^2 \leq 1$; $x, y \geq 0$ . Consider the problem

$$\Delta u = u^2 \quad \text{on } \Omega$$

$$u = 30/(x+2y+1)^2 \quad \text{on } \partial\Omega . \qquad (5.1)$$

The exact solution is $u(x, y) = 30/(x+2y+1)^2$; also $K = 2$ .

<u>Problem 2</u>　Same equation as in Problem 1 . $\bar{\Omega}$ the square $0 \leq x, y \leq 1$. This boundary value problem has been solved numerically in Greenspan [1964] with an $O(h^2)$ method.

<u>Problem 3</u>　$\bar{\Omega}$ as in Problem 1. Equation

$$\Delta u = \frac{\pi^2}{2} (x^2 + y^2) e^u \qquad \text{on } \Omega$$

$$u = -2 \log (\sin (\frac{\pi}{2} xy + \frac{\pi}{4})) \qquad \text{on } \partial\Omega ;$$

$K = \pi^2/2$ . The exact solution is the boundary function extended to $\bar{\Omega}$ .

<u>Problem 4</u>　$\bar{\Omega}$ as in Problem 2. Equation as in Problem 3.

3. We describe now the numerical results. In examples 1 through 4 the $\varepsilon_j^{(i)}$ are defined by

$$\varepsilon_j^{(i)} = \| u_j^{(i)} - \Delta_h u^* \| / \| \Delta_h u^* \| ,$$

where $u^*$ is the exact solution.

<u>Example 1</u>     Solution of Problem 1 by means of the S.E. method of Section I with $k = 3$ and using 6-th order interpolation at the boundary ( see ( 4.3)). The numerical results for the two basic mesh sizes $h_1 = \frac{1}{8}, \frac{1}{10}$ are given in Tables 1a and 1b . The lack of diagonal dominance at points in I( h) did not create any noticeable difficulty.

<u>Example 2</u>     Solution of Problem 3 by the same method used in Example 1. Numerical results are given in Tables 2a and 2b. Again no difficulty appeared because of the lack of diagonal dominance. In this example is seen even more clearly than in Example 1 the $h^2$, $h^4$, $h^6$ improvement of the successive extrapolates, despite the small difference between the two basic meshes. Let us remark that when performed without the 6-th order boundary interpolation no such an improvement was obtained. We also give in this case the corresponding computation times.

<u>Example 3</u>     A 4-step extrapolation with $h_1 = \frac{1}{8}$ was used on Problem 2. The numerical results are given in Table 3a. Also a 3-step extrapolation with $h_1 = \frac{1}{10}$ was performed. The corresponding numerical results are given in Table 3b.

<u>Example 4</u>     Problem 4 was solved:

a)  Using a 4-step extrapolation with $h_1 = \frac{1}{8}$ .

$U_4^{(4)}$ did not come out with the expected precision. The only reason we can suggest to explain this behavior is that the accuracy required in the inner iteration was beyond the machine's capability ( i.e. word length in simple precision) . In such a case the inner iterations will be "too incomplete" to furnish an accurate final result.

b)   Using a 5-step extrapolation with $h_1 = \frac{1}{4}$ .

Here the phenomena observed in (a) is more accentuated. $U_4^{(4)}$ is obtained with the expected accuracy (and it is better than $U_4^{(4)}$ in (a)), while $U_5^{(5)}$ is less accurate than $U_4^{(4)}$ and just slightly better than $U_3^{(3)}$. We attribute this to the same causes as in (a). It is necessary to observe that for $h_5 = \frac{1}{64}$ the resulting system has 3969 quite complicated nonlinear equations (see (4.2) and (5.2)).

c)   Using a 5-step extrapolation with $h_1 = \frac{1}{2}$ .

This gave the solution at the center of the square with an absolute error less than $10^{-8}$ in about 40 seconds of CDC 3600 computing time. In this case there was no trouble with the highest order extrapolate. The numerical results are given in Tables 4a, b, and c respectively.

Example 5    Solution of Problem 1 by means of L. D. C. procedure of II. 2 with $h^6$ interpolation at the boundary. The iteration of II. 6 is not theoretically necessary in this case; however, when applied, the second iteration gives some improvement over the first, especially for the smaller h's. We attribute this to the elimination of round-off. Using the notation of IV. 6 we denote by $\varepsilon^{(i)}$ the relative error

$$\varepsilon^{(i)} = \frac{\| U^{(i)}(h) - \Delta_h u^* \|}{\| \Delta_h u^* \|}$$

We count each iteration (4.18) as a Newton iteration.

Example 6    Solution of Problem 1 by means of L. D. C. with only $h^2$ boundary interpolation and iteration II. 6. The notation is as in Example 5 and the numerical results are given in Table 6. For $h = \frac{1}{40}$, $U^{(5)}$ still

gave some improvement: $\mathcal{E}^{(5)} = 10^{-6}$ , with 150 extra SOR iterations.

Example 7    Solution of Problem 4 by L. D. C. and iteration II. 6.
Numerical results are given in Table 7.

VI.  Comparisons between the different methods

Two basic requirements in the solution of boundary value problems by finite differences are the desired accuracy $\mathcal{E}$, and the desired "minimum definition" of the numerical solution, i.e. the minimum number of points (coarsest mesh) at which the solution is needed for practical purposes.

These requirements, and the available computing machinery will generally be the main guidelines in choosing a method of solution or, ultimately, in deciding if the problem can be solved at all.   The two kinds of methods of high order accuracy we have described and used in the former Sections have particular characteristics that make them convenient in a wide range of applications.

While our numerical examples are significant by themselves we have in mind in this discussion possible applications to problems in higher dimensions: three dimensional elliptic boundary value problems, parabolic equations, etc. for which storage and computation time are even more critical factors (cf. Forsythe and Wasow [1960], pp. 11-14  for some interesting figures and predictions for the case of linear equations; seven years after these figures have been published three dimensional problems still cannot be solved either very accurately or in too much detail ).

If high accuracy is required ( say $\mathcal{E} = 10^{-7}$ ) but not necessarily high definition ( i.e. basic mesh $h_1$ not too small) then it is clear that the usual

criticism about the standard $O(h^2)$ method is very much valid. In fact we see from Table 4a that in order to reduce $\mathcal{E}_i^{(0)}$ to less than $10^{-7}$ we would need to take a mesh $h < \frac{1}{800}$ (64000 points) which would give an accurate but exaggeratedly detailed discrete solution. On the other hand, by using three extrapolations and a finest mesh of $\frac{1}{64}$ ($\sim 4000$ points) we have obtained the solution at 49 points with the required precision. If we need not only high precision but also good definition then with some extra work the L.D.C. procedure gives us the solution at 900 points with accuracy $4 \times 10^{-7}$ by using $h = \frac{1}{32}$ (see Table 7). In conclusion, S.E. is of simpler application than L.D.C. and in principle can give more accuracy for a given basic mesh. However, as we see from some of our computation times, the most significant computation in S.E is the solution of the basic problem for the finest mesh. Thus it is fair to compare three or four extrapolation steps to one application of L.D.C. for the finest mesh used in S.E. In this case, at least in our examples, the precision obtained is comparable while we obtain much more detail from L.D.C.

Another factor which has not been considered here but which can be a source of difficulties is the ill-conditioning of the system of equations for small mesh sizes. This has been pointed out in Fox [1950] and it is known to generate difficulties in the solution of differential equations of higher orders, such as the biharmonic equation. This is then still another reason for keeping the step size reasonably large. Moreover, since here the difficulty is intrinsic to the problem it is unlikely that it will be solved by the appearance on the market of faster and larger computing machines.

It is our hope that the iterated deferred correction procedure will give both definition and accuracy for any given reasonable specification, thus providing a more flexible tool for the accurate solution of multidimensional boundary value problems.

Table 1 a

$$\Delta u = u^2, \quad u = \frac{30}{(x+2y+1)^2} \quad \text{on}$$

Method: Successive extra-
polations; 6-th order
boundary interpolation.

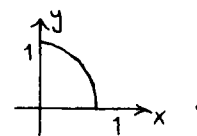| i | 1 | 2 | 3 |
|---|---|---|---|
| $h_i$ | $\frac{1}{10}$ | $\frac{1}{20}$ | $\frac{1}{40}$ |
| $n_i$ | 67 | 292 | 1214 |
| $\omega_i$ | 1.3 | 1.5 | 1.5 |
| $\varepsilon_i^{(0)}$ | $1.5 \times 10^{-3}$ | $4.3 \times 10^{-4}$ | $1.1 \times 10^{-4}$ |
| $\varepsilon_i^{(i)}$ | $1.5 \times 10^{-3}$ | $5.2 \times 10^{-5}$ | $1.7 \times 10^{-6}$ |
| Newton iter. (Total SOR iter.) | 4 ( 28) | 4 ( 171) | 6 ( 472) |

Table 1 b

| i | 1 | 2 | 3 |
|---|---|---|---|
| $h_i$ | $\frac{1}{8}$ | $\frac{1}{16}$ | $\frac{1}{32}$ |
| $n_i$ | 41 | 183 | 770 |
| $\omega_i$ | 1.3 | 1.4 | 1.5 |
| $\varepsilon_i^{(0)}$ | $2.3 \times 10^{-3}$ | $6.7 \times 10^{-4}$ | $1.7 \times 10^{-4}$ |
| $\varepsilon_i^{(i)}$ | $2.3 \times 10^{-3}$ | $1.7 \times 10^{-4}$ | $1.1 \times 10^{-5}$ |
| Newton iter. (Total SOR iter.) | 4 ( 43) | 3 ( 66) | 7 ( 306) |

## Table 2a

$$\Delta u = \frac{\pi^2}{2}(x^2+y^2)e^u, \quad u = -2\log\left(\sin\frac{\pi}{2}\left(xy+\tfrac{1}{2}\right)\right) \text{ on}$$

Method: Same as in Table 1.

| i | 1 | 2 | 3 |
|---|---|---|---|
| $h_i$ | $\frac{1}{10}$ | $\frac{1}{20}$ | $\frac{1}{40}$ |
| $\omega_i$ | 1.3 | 1.5 | 1.5 |
| $\varepsilon_i^{(0)}$ | $5.2\times10^{-4}$ | $1.4\times10^{-4}$ | $3.5\times10^{-5}$ |
| $\varepsilon_i^{(i)}$ | $5.2\times10^{-4}$ | $1.5\times10^{-5}$ | $2.2\times10^{-7}$ |
| Newton iter. (Total SOR iter.) | 2 (25) | 10 (102) | 6 (518) |

Time: 3' 25"

## Table 2b

| i | 1 | 2 | 3 |
|---|---|---|---|
| $h_i$ | $\frac{1}{8}$ | $\frac{1}{16}$ | $\frac{1}{32}$ |
| $\omega_i$ | 1.3 | 1.4 | 1.5 |
| $\varepsilon_i^{(0)}$ | $8.0\times10^{-4}$ | $2.2\times10^{-4}$ | $5.5\times10^{-5}$ |
| $\varepsilon_i^{(i)}$ | $8.0\times10^{-4}$ | $3.1\times10^{-5}$ | $8.1\times10^{-7}$ |
| Newton iter. (Total SOR iter) | 1 (28) | 2 (60) | 5 (333) |

Time: 1' 30"

#687

$$\Delta u = u^2, \quad u = \frac{30}{(x+2y+1)^2} \quad \text{on} \quad \quad \text{Method: Successive extrapolations.}$$

| $i$ | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| $h_i$ | $\frac{1}{8}$ | $\frac{1}{16}$ | $\frac{1}{32}$ | $\frac{1}{64}$ |
| $n_i$ | 49 | 225 | 961 | 3969 |
| $\omega_i$ | 1.7 | 1.7 | 1.8 | 1.8 |
| $\varepsilon_i^{(0)}$ | $2.3 \times 10^{-3}$ | $6.7 \times 10^{-4}$ | $1.7 \times 10^{-4}$ | $4.4 \times 10^{-5}$ |
| $\varepsilon_i^{(i)}$ | $2.3 \times 10^{-3}$ | $1.0 \times 10^{-4}$ | $1.9 \times 10^{-6}$ | $1.6 \times 10^{-8}$ |
| Newton iter. (Total SOR iter.) | 6 ( 48) | 3 ( 41) | 4 ( 66) | 5 ( 483) |

Time: 7 ' 21"

Table 3b

| $i$ | 1 | 2 | 3 |
|---|---|---|---|
| $h_i$ | $\frac{1}{10}$ | $\frac{1}{20}$ | $\frac{1}{40}$ |
| $n_i$ | 81 | 361 | 1521 |
| $\omega_i$ | 1.3 | 1.5 | 1.5 |
| $\varepsilon_i^{(0)}$ | $1.5 \times 10^{-3}$ | $4.3 \times 10^{-4}$ | $1.1 \times 10^{-4}$ |
| $\varepsilon_i^{(i)}$ | $1.5 \times 10^{-3}$ | $5.2 \times 10^{-5}$ | $8.0 \times 10^{-7}$ |
| Newton iter. (Total SOR iter.) | 4 ( 36) | 3 ( 97) | 10 ( 638) |

## Table 4a

$$\Delta u = \frac{\pi^2}{2}(x^2+y^2)e^u, \quad u = -2\log\left(\sin\frac{\pi}{2}(xy+\tfrac{1}{2})\right) \quad \text{on} \quad \text{Method: S. E.}$$

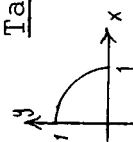| $i$ | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| $h_i$ | $\frac{1}{8}$ | $\frac{1}{16}$ | $\frac{1}{32}$ | $\frac{1}{64}$ |
| $\omega_i$ | 1.7 | 1.7 | 1.78 | 1.78 |
| $\varepsilon_i^{(0)}$ | $1.2\times10^{-3}$ | $3.1\times10^{-4}$ | $7.8\times10^{-5}$ | $2.0\times10^{-5}$ |
| $\varepsilon_i^{(i)}$ | $1.2\times10^{-3}$ | $3.3\times10^{-5}$ | $5.1\times10^{-7}$ | $1.5\times10^{-7}$ |
| Newton iter. (Total SOR iter.) | 4 (38) | 2 (42) | 4 (95) | 11 (664) |

Time: 10' 36"

## Table 4b

| $i$ | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| $h_i$ | $\frac{1}{4}$ | $\frac{1}{8}$ | $\frac{1}{16}$ | $\frac{1}{32}$ | $\frac{1}{64}$ |
| $\varepsilon_i^{(i)}$ | $3.8\times10^{-3}$ | $2.6\times10^{-4}$ | $6.5\times10^{-6}$ | $6.6\times10^{-8}$ | $1.4\times10^{-7}$ |

Time: 10' 44"

**Table 4c**

| i | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| $h_i$ | $\frac{1}{2}$ | $\frac{1}{4}$ | $\frac{1}{8}$ | $\frac{1}{16}$ | $\frac{1}{32}$ |
| $\mathcal{E}_i^{(i)}$ | $9.2\times10^{-3}$ | $2.0\times10^{-3}$ | $9.5\times10^{-5}$ | $1.7\times10^{-6}$ | $1.5\times10^{-8}$ |

Time: 40''

**Table 5**

$$\Delta u = u^2, \quad u = \frac{30}{(x+2y+1)^2} \quad \text{on}$$

Method: LDC + iter; 6-th order boundary interpolation.

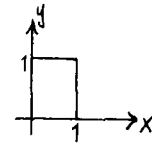| h | $\mathcal{E}^{(0)}$ | $\mathcal{E}^{(1)}$ | $\mathcal{E}^{(2)}$ | Newton iter. (Total SOR iter.) | 3 |
|---|---|---|---|---|---|
| $\frac{1}{10}$ | $1.5\times10^{-3}$ | $9.1\times10^{-5}$ | $8.0\times10^{-5}$ | 8 (153) | 1.4 |
| $\frac{1}{20}$ | $4.3\times10^{-4}$ | $1.25\times10^{-5}$ | $2.1\times10^{-6}$ | 13 (480) | 1.4 |
| $\frac{1}{40}$ | $1.1\times10^{-4}$ | $1.2\times10^{-6}$ | $1.8\times10^{-7}$ | 14 (888) | 1.5 |

Table 6

Same problem and method as in Table 5; $h^2$-interpolation at the boundary.

| $h$ | $\varepsilon^{(0)}$ | $\varepsilon^{(1)}$ | $\varepsilon^{(2)}$ | $\varepsilon^{(3)}$ | $\varepsilon^{(4)}$ | Newton iter. (Tot.SOR iter.) |
|---|---|---|---|---|---|---|
| $\frac{1}{10}$ | $2.1\times10^{-3}$ | $6.2\times10^{-4}$ | $4.3\times10^{-4}$ | $2.1\times10^{-4}$ | $2.5\times10^{-4}$ | 9(112) |
| $\frac{1}{20}$ | $5.4\times10^{-4}$ | $1.8\times10^{-4}$ | $6.7\times10^{-5}$ | $1.7\times10^{-5}$ | $1.7\times10^{-5}$ | 11 (576) |
| $\frac{1}{40}$ | $1.4\times10^{-4}$ | $4.7\times10^{-5}$ | $1.6\times10^{-5}$ | $5.2\times10^{-6}$ | $2.5\times10^{-6}$ | 9 (1011) |

Table 7

$$\Delta u = \frac{\pi^2}{2}(x^2+y^2)e^u, \quad u = -2\log\left(\sin\frac{\pi}{2}(xy+\tfrac{1}{2})\right) \text{ on }$$

Method: LDC iterated.

| $h$ | $\varepsilon^{(0)}$ | $\varepsilon^{(1)}$ | $\varepsilon^{(2)}$ | $\varepsilon^{(3)}$ | Newton iter. (Tot.SOR iter.) |
|---|---|---|---|---|---|
| $\frac{1}{2}$ | $5.4\times10^{-3}$ | $3.6\times10^{-3}$ | – | – | 3 (16) |
| $\frac{1}{4}$ | $2.6\times10^{-3}$ | $6.1\times10^{-4}$ | $5.1\times10^{-4}$ | – | 5 (50) |
| $\frac{1}{8}$ | $8.2\times10^{-4}$ | $7.5\times10^{-5}$ | $4.8\times10^{-5}$ | $4.5\times10^{-5}$ | 8 (132) |
| $\frac{1}{16}$ | $2.1\times10^{-4}$ | $7.4\times10^{-6}$ | $3.0\times10^{-6}$ | $2.8\times10^{-6}$ | 8 (170) |
| $\frac{1}{32}$ | $5.4\times10^{-5}$ | $5.8\times10^{-7}$ | $4.3\times10^{-7}$ | – | 6 (378) |

## Table 8

| | $\omega = 1$ | $\omega = 1.53$ | $\omega = 1.8$ |
|---|---|---|---|
| Method I | 98 ( 98) | 28 ( 28) | 71 ( 71) |
| Method II | 4 ( 204) | 4 ( 65) | 4 (148) |
| Method III | 15 (105) | 11 ( 30) | 14 ( 73) |
| Method IV | 25 ( 98) | 21 ( 30) | 33 (72) |

# APPENDIX

Bellman, Juncosa, and Kalaba [1961], Greenspan and Parter [1965], and Ortega and Rockoff [1965] have reported numerical experiments for the problem

$$( * ) \qquad \Delta u = e^{u} \, , \qquad u(x,y) = x + 2y \qquad \text{on the boundary of the}$$

unit square. They generate the nonlinear difference equations in the same way as described in Section IV. Greenspan and Parter gave results for two different schemes:

Method I: take one Gauss-Seidel sweep per Newton iteration;

Method II: solve the linear systems of Newton's method by SOR.

Ortega and Rockoff have investigated these for varying values of the overrelaxation parameter and the figures for Methods I and II in Table 8 are from Ortega [1966] who graciously made them available to us before their publication. We compare these two methods with our procedure (Theorem 3.1) for stopping the inner SOR iterations. Since the previous authors used a convergence criterion of the form

$$( ** ) \qquad \| V_{i+1} - V_i \| \leq 10^{-6}$$

(in the notation of Section III), we have to modify the convergence criterion given by Theorem 3.1.

Lemma  With the notation and hypotheses of Theorem 3.1, suppose that, given $\mathcal{E}$, we proceed as indicated by (3.19) until

$$\alpha H_{i-1} \eta_{i-1} = E_{i-1} < \alpha \, \mathcal{E}/4B \, . \qquad (1)$$

and then before computing $V_i$ we replace $E_{i-1}$ by $\widetilde{E}_{i-1} = \mathcal{E}/4B$. Taking also

$\widetilde{E}_i = \widetilde{E}_{i-1}$, then

$$\|\widetilde{V}_{i+1} - V_i\| \leqq \mathcal{E} \quad . \tag{2}$$

Proof: If $V_{i+1}$ were computed exactly $(E_i = 0)$ then we would have

$$\|V_{i+1} - V_i\| \leqq B\eta_i \leqq B((1-\alpha)H_{i-1}\eta_{i-1} + \widetilde{E}_{i-1}) \leqq \mathcal{E}/2 . \tag{3}$$

Let us call $\widetilde{V}_{i+1}$ to the solution obtained by allowing the residual to be in

norm at most as large as $\widetilde{E}_i$ . Then

$$\|\widetilde{V}_{i+1} - V_{i+1}\| \leqq B\|e_i\| \leqq B\widetilde{E}_i \leqq \mathcal{E}/4 . \tag{4}$$

Finally (2) follows from (3) and (4).

It is clear from this proof and Theorem 3.1 that $\|\Phi_h(V_i)\| \leqq \mathcal{E}$ is a

more convenient convergence criterion than (2), especially considering the

straightforward error estimation given by (3.21).

In Table 8 (p. 44) we give the results of Method I and II as applied to [(*)]

with a step size of $h = 0.1$. For the inner iterations of Method II ([**]) is

used as a convergence criterion.

For Methods III and IV we start with an empirical criterion for stopping

the inner iterations until $H_0 < 1$ when we can switch to the theoretical

criterion of Theorem 3.1. In Method III the starting procedure consists of

reducing the residual of the linear equations at the i-th Newton step below

$\widetilde{\mathcal{E}}_i = 4 \times (0.25)^i$, but allowing only a maximum of four SOR sweeps. In

Method IV we use Method I, which is known to be convergent, as a

starting procedure. The first figure in Table 8 indicates the number of Newton

iterations while the figure in parentheses indicates the total number of SOR

sweeps.

#687

# REFERENCES

1. Ballester, C. and V. Pereyra [1966] "On the construction of discrete approximations to linear differential expressions", MRC Technical Report #671, University of Wisconsin, Madison.

2. Bers, L. [1953] "On mildly nonlinear difference equations of elliptic type", Jour. Res. Nat. Bur. Stand. 51, 229-236.

3. Bickley, W.G., S. Michaelson and M.R. Osborne [1961] "On finite difference methods for the numerical solution of boundary value problems", Proc. Roy. Soc. London, A262, 219-236.

4. Bulirsch, R. and J. Stoer [1964] "Fehlerabschätzungen und Extrapolation mit rationalen Funktionen bei Verfahren von Richardson-Typus", Numer. Math. 6, 413-427.

5. Douglas, J. [1961] "Alternating direction iteration for mildly nonlinear elliptic difference equations", Numer. Math. 3, 92-98.

6. Forsythe, G. and W. R. Wasow [1960], Finite-Difference Methods for Partial Differential Equations, Wiley, New York.

7. Fox, L. [1947] "Some improvements in the use of relaxation methods for the solution of ordinary and partial differential equations", Proc. Roy. Soc. London, A190, 31-59.

8. _____, [1950] "The numerical solution of elliptic differential equations when the boundary conditions involve a derivative", Phil. Trans. Roy. Soc. London, A242, 345-378.

9. _____, [1962] (Editor) Numerical Solution of Ordinary and Partial Differential Equations, Pergamon Press, Oxford.

10. Gragg, W.B. [1965] "On extrapolation algorithms for ordinary initial value problems", J. SIAM Numer. Anal. Ser B., 2, #3, 384-403.

11. Greenspan, D. [1964] "On approximating extremals of functionals. Part I: The method and examples for boundary value problems", MRC Tech. Report # 466, University of Wisconsin, Madison.

12. _____, [1965] <u>Introductory Numerical Analysis of Elliptic Boundary Value Problems</u>, Harper & Row, New York.

13. _____, and S. Parter [1965] "Mildly nonlinear elliptic partial differential equations and their numerical solutions, II", Numer. Math., _7_, 129-146.

14. Henrici, P. [1962] <u>Discrete Variable Methods in Ordinary Differential Equations</u>, Wiley, New York.

15. Kantorovich, L.V. and G. P. Akilov [1964] <u>Functional Analysis in Normed Spaces</u>, McMillan, New York.

16. Ortega, J. [1966] "Relaxation methods for nonlinear equations", National SIAM Metting, Iowa City, Iowa, May 1966.

17. _____ and M. Rockoff [1965] "Nonlinear difference equations and Gauss-Seidel type iterative methods", Tech. Rep. 65-20- NsG-398, University of Maryland, College Park.

18. Parter, S. [1965] "Mildly nonlinear elliptic partial differential equations and their numerical solutions, I", Numer. Math., _7_, 113-128.

19. Pereyra, V. [1965] "The difference correction method for non-linear two-point boundary value problems", Tech. Report CS18, Stanford University.

20. _____, [1966] "On improving an approximate solution of a functional equation by deferred corrections", Numer. Math., _8_, 376-391.

21. Richardson, L. F. [1910] "The approximate arithmetical solution by finite differences of physical problems involving differential equations", Phil. Trans. Roy. Soc. London, _210_, 307-357.

22. Schechter, S. [1962] "Iteration methods for nonlinear problems", Trans. AMS, _104_, 179-189.

23. Stetter, H. [1965] "Asymptotic expansions for the error of discretization algorithms for nonlinear functional equations", Numer. Math., _7_, 18-31.

24. Varga, R. [1962] <u>Matrix Iterative Analysis</u>, Prentice Hall, New Jersey.

25. Volkov, E.A. [1957] "An analysis of an algorithm of heightened precision for the solution of Poisson's equation", Vych. Math., $\underline{1}$, 62-80 (Russian). AMS Translations, Series 2, $\underline{35}$, 117-136 (1964).

26. _____, [1963] "Methods of refinements using higher-order differences and $h^2$- extrapolation", Doklady, $\underline{150}$ (Russian), Soviet Math., $\underline{4}$, #3, 671-674 (1963).

27. _____, [1965] "Solution of Dirichlet problem by the method of refining with high-order differences", Doklady, $\underline{164}$, #3 (Russian), Soviet Math., $\underline{6}$, #5, 1234-1237 (1965).

28. Wasow, W. [1955] "Discrete approximations to elliptic differential equations", Jour. Appl. Math. Phys., $\underline{6}$, #2, 81-97.

29. Bellman, R.E. and R.E. Kalaba [1965], Quasilinearization and Nonlinear Boundary-Value Problems, American Elsevier Pub. Co., New York.

30. Lees, M. [1966] "Discrete methods for nonlinear two-point boundary value problems", in Numerical Solution of Partial Differential Equations (edited by J.H. Bramble, Academic Press), 59-72.

31. Panov, D.J. [1963] Formulas for the Numerical Solution of Partial Differential Equations by the Method of Differences. F. Ungar Publ. Co., New York.

32. Bellman, R., M. Juncosa, and R. Kalaba [1961], "Some numerical experiments using Newton's method for nonlinear parabolic and elliptic boundary-value problems", Comm. ACM $\underline{4}$, 187-191.