

STICHTING
MATHEMATISCH CENTRUM
2e BOERHAAVESTRAAT 49
AMSTERDAM
REKENAFDELING

MR 104

A method for solving elliptic difference equations

by

P.J. van der Houwen



April 1969

Contents

1. Introduction	1
2. Mathematical preliminaries	3
3. A symmetric matrix representation for the model problem	7
4. Numerical results for the model problem	12
5. References	15

1. Introduction

In recent years a large number of methods have been developed to improve the rate of convergence of iterative processes for solving elliptic difference equations. In order to analyse these iterative methods one usually takes the discrete analogue of the Dirichlet problem for Poisson's equation on a square of side π , which serves as a "model problem" in elliptic difference equations. For the model problem one often may obtain an analytic expression for the rate of convergence so that one can learn about the properties of the iterative method. Until 1950, the rates of convergence of the known iterative procedures were of order h^2 , h being the mesh size of the grid used. Although Richardson had indicated as early as 1910 an iterative method with a potentially greater rate of convergence, he did not obtain a better method, because the parameters required for his method were not chosen optimally.

In 1950 Young, and independently Frankel, proposed a very powerful method with potentially a rate of convergence of order h . This method is known as the "method of successive overrelaxation of Young" (SOR method) or as the "extrapolated Liebmann method" as it was called by Frankel. In the paper of Frankel another method of the second order was described which also has a rate of convergence of order h . We shall call it "Frankel's method".

In 1953 Shortley applied Richardson's method with the optimal values of the required parameters. Asymptotically, the method has a rate of convergence of order h , however, for a large number of iterations the method turned out to be numerically unstable.

In 1955 Sheldon combined Richardson's method and the SOR method to obtain a process, the "method of symmetric successive overrelaxation", which experimentally proved to be of order \sqrt{h} for the model problem. We shall call the method the SSOR method.

A new approach in accelerating iterative methods was given by Peaceman and Rachford in 1955 and by Douglas and Rachford in 1956. Asymptotically, their theories result in a still greater rate of convergence, namely of order $1/\ln h^{-1}$. However, it was shown in 1959 by Birkhoff and Varga that

the theory only holds for the model problem, so that the actual value of the method is doubtful.

In 1958 the disadvantage of Richardson's method, namely its instability, was overcome by Stiefel, who introduced a second order version of the method which can be proved to be stable.

In the following years a large number of contributions to elliptic difference equations were made. The greater part of these are modifications or generalizations of the methods already mentioned.

In this paper a method is described based on the second order Richardson method and the SOR method. We give a detailed analysis of the method for the model problem showing that the rate of convergence is of order \sqrt{h} . On a computer we got rates of convergence which agree with the theoretically predicted rates of convergence.

The author acknowledges Mr. P. Beertema for writing the computer program by which the numerical results were obtained.

2. Mathematical preliminaries

In this section we review Richardson's method of second degree for solving matrix equations of the type

$$(2.1) \quad Lu = f,$$

where L is a symmetric matrix with positive eigenvalues, f is a known vector and u is the unknown vector. The method is defined by the recurrence relations

$$u_1 = u_0 - \omega_0(Lu_0 - f),$$

$$u_{k+1} = \alpha_k u_k + (1 - \alpha_k)u_{k-1} - \omega_k(Lu_k - f), \quad k = 1, 2, \dots,$$

$$(2.2) \quad \alpha_k = 2y_0 \frac{T_k(y_0)}{T_{k+1}(y_0)}, \quad k = 1, 2, \dots,$$

$$\omega_0 = \frac{2}{b+a}, \quad \omega_k = \frac{4}{b-a} \frac{T_k(y_0)}{T_{k+1}(y_0)}, \quad k = 1, 2, \dots,$$

where u_0 is an initial guess for the solution u , $y_0 = (b+a)/(b-a)$, $[a,b]$ is a positive interval containing all eigenvalues λ of L , and $T_k(y)$ is the Chebyshev polynomial of degree k . $T_k(y)$ satisfies the recurrence relation

$$(2.3) \quad T_0(y) = 1, \quad T_1(y) = y, \quad T_{k+1}(y) = 2yT_k(y) - T_{k-1}(y), \quad k = 1, 2, \dots$$

This iteration scheme was proposed by Stiefel [1958] and is called Richardson's method of second degree. A detailed discussion of Richardson's method and some accelerating procedures may be found in van der Houwen [1968], chapter IV. In the present investigation we are mainly interested in the average rate of convergence after K iterations of the iteration process for some operators L arising from the numerical solution of Poisson's equation. The rate of convergence, denoted by $R(K)$, is given by (cf. Forsythe and Wasow [1960], p. 231)

$$(2.4) \quad R(K) = \frac{1}{K} \ln T_K(y_0).$$

If $a \ll b$ it can be derived from the properties of $T_k(y)$ that

$$(2.4') \quad R(K) \sim 2\sqrt{\frac{a}{b}} - \frac{1}{K} \ln 2.$$

The quotient b/a is called the P-condition number of the matrix L and will be denoted by $P(L)$.

We now summarize the theory of the numerical solution of Poisson's equation needed in the following sections. A detailed treatment of this theory may be found in van der Houwen [1967].

Consider Poisson's equation

$$(2.5) \quad \Delta U + F = 0, \quad \Delta = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2}$$

on a square of side π with boundary conditions of the first kind. Let X_+ and Y_+ represent translations $+h$ along the x - and y -axis respectively. On a grid of square meshes of side h the following discrete analogue of (2.5) may be defined:

$$(2.6) \quad D(\gamma)u + f = 0,$$

where

$D(\gamma) = 1$ at the boundary points,

$D(\gamma) = L_1(X_+ + X_-)(Y_+ + Y_-) + L_2(X_+ + X_- + Y_+ + Y_-) + L_4$ at the internal points,

$$L_1 = \frac{2 - \gamma}{2h^2}, \quad L_2 = \frac{\gamma - 1}{h^2}, \quad L_4 = -\frac{2\gamma}{h^2}, \quad 1 \leq \gamma \leq 2,$$

and where u and f represent grid functions defined at the grid points. When $-f$ assumes the boundary values of U at the boundary points and f assumes the values of F at the internal grid points, then it can be proved that the solution of the boundary value problem (2.5) and the solution of the discrete problem (2.6) differ by a term $O(h^2)$ (cf. Forsythe and Wasow [1960], section 23). Problem (2.6) will be called the model problem.

The operator $D(\gamma)$ can be represented by a symmetric matrix operator. The eigenvalues of $D(\gamma)$ appear to be negative so that we define

$$(2.7) \quad L(\gamma) = -D(\gamma).$$

For $\gamma = 2$ this operator was considered by Frank [1960]. It is easily verified that

$$(2.8) \quad P(L(2)) \sim 4h^{-2},$$

so that

$$(2.9) \quad R(K) \sim h - \frac{\ln 2}{K} \sim h - \frac{0.693}{K}.$$

For $\gamma = 1$ it can be shown that (see van der Houwen [1967])

$$(2.10) \quad R(K) \sim \sqrt{2} h - \frac{0.693}{K}.$$

At this point we remark that the iteration scheme with $\gamma = 1$ may be interpreted as the iteration scheme with $\gamma = 2$ applied to a square which is rotated over 45° (see figure 2.1 and 2.2).

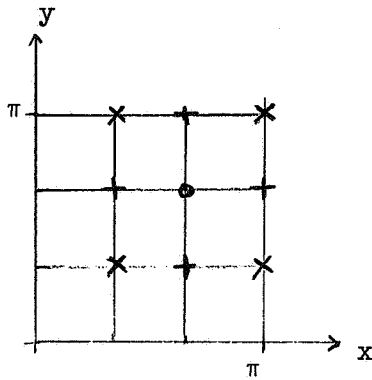


fig. 2.1

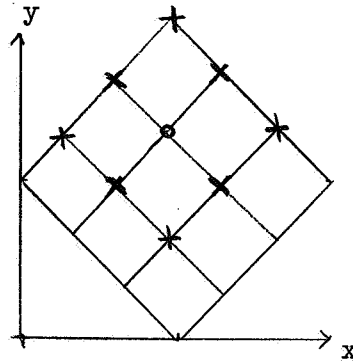


fig. 2.2

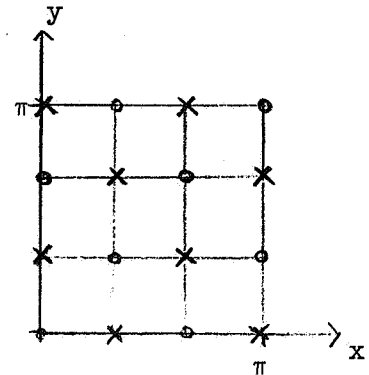


fig. 2.3

\times corresponds to $\gamma = 1$, $+$ corresponds to $\gamma = 2$

Further, we observe that using the operator $L(1)$ permits us to calculate u_k only in those net points (j,l) of figure 2.3 which are denoted by either a dot or a cross.

One point of departure in accelerating Richardson's method is to replace the matrix equation $L(\gamma)u = f$ by an equation

$$(2.11) \quad L'(\gamma)u = f'$$

which has the same solution u , but in which $L'(\gamma)$ has a lower P-condition number than $L(\gamma)$. Such an approach was given in van der Houwen [1967] for the model problem mentioned above. The operator $L'(\gamma)$ and the vector f' were defined by

$$(2.12) \quad L'(\gamma) = -(1 - D_1(\gamma))^{-1}D(\gamma), \quad f' = (1 - D_1(\gamma))^{-1}f,$$

where

$$(2.13) \quad \begin{aligned} D_1(\gamma) &= 0 \text{ at the boundary points,} \\ D_1(\gamma) &= p(L_1(X_+Y_- + X_-Y_-) + L_2(X_- + Y_-)) \text{ at the} \\ &\quad \text{internal points.} \end{aligned}$$

For $\gamma = 1$ and $\gamma = 2$ it is possible to choose the parameter p such that the eigenvalues of $L'(\gamma)$ are real. The operator $L'(2)$ then reduces to the operator occurring in Liebmann's method (compare Forsythe and Wasow [1960]). The P-condition number of this operator appears to be $\frac{1}{4}$ of the value of $P(L(2))$:

$$(2.14) \quad P(L'(2)) \sim h^{-2}.$$

Therefore, the asymptotic rate of convergence of Richardson's method with respect to $L'(2)$ is twice as large as the asymptotic rate of convergence of the scheme used by Frank. However, operators of type (2.12) are not symmetric and, in fact, the eigenvalues of $L'(1)$ and $L'(2)$ are very ill-conditioned, so that an arbitrary initial approximation u_0 will be a very poor approximation of the solution u . In Coolen and van der Houwen [1968] a method is given which eliminates the ill-conditioned eigenfunction components from the initial approximation. This preconditioning was found very successful, but the method explicitly uses the fact that the eigenvalues of the ill-conditioned components are known to be the greatest eigenvalues. In other cases, preconditioning may be less successful. Therefore, it is desirable to construct matrix representations with better conditioned eigenfunctions.

3. A symmetric matrix representation for the model problem

In this section a symmetric matrix representation of the boundary value problem is given for which the P-condition number is of order h^{-1} .

Let us define the operator

$$(3.1) \quad L'(\gamma) = -\frac{1}{2} [(1 - D_1(\gamma))^{-1} + (1 - D_2(\gamma))^{-1}] D,$$

where

$$(3.2) \quad \begin{aligned} D_2(\gamma) &= 0 \text{ at the boundary points,} \\ D_2(\gamma) &= p(L_1(X_- Y_+ + X_+ Y_-) + L_2(X_+ + Y_+)) \text{ at the} \\ &\quad \text{internal points.} \end{aligned}$$

This operator arises from averaging the operators which correspond to Gauss-Seidel's method starting from opposite corner points. We may expect that this averaging eliminates the ill-conditioning of the eigenfunctions.

Theorem 3.1.

Let $\gamma = 1$, then $L'(\gamma)$ is symmetric with eigenfunctions

$$(3.3) \quad e(n,m) = \sin njh \sin mlh, \quad n, m = 1, 2, \dots, \frac{\pi}{h} - 1$$

and eigenvalues

$$(3.4) \quad \lambda(n,m) = 4h^{-2} \frac{(2 - \Omega v \mu)(1 - v \mu)}{\Omega^2 v^2 - 4\Omega v \mu + 4}, \quad n, m = 1, 2, \dots, \frac{\pi}{h} - 1.$$

Here we have $v = \cos nh$, $\mu = \cos mh$ and $\Omega = 2ph^{-2}$.

Proof

It is easily verified that $D(1)$ has the eigenfunctions $e(n,m)$ defined above with eigenvalues

$$(3.5) \quad \delta(n,m) = -2h^{-2}(1 - \cos nh \cos mh) = -2h^{-2}(1 - v \mu).$$

From this it follows that

$$L'(1)e(n,m) = -\frac{1}{2} \delta(n,m)(v_1 + v_2),$$

where

$$v_1 = (1 - D_1)^{-1}e(n,m) = e(n,m) + D_1 v_1,$$

$$v_2 = (1 - D_2)^{-1}e(n,m) = e(n,m) + D_2 v_2.$$

By substituting (3.3) we find $L'(1)e(n,m) = \lambda(n,m)e(n,m)$ where $\lambda(n,m)$ is given by (3.4).

The eigenfunctions $e(n,m)$ are orthogonal and the eigenvalues $\lambda(n,m)$ are real. Therefore, $L'(1)$ is a symmetric matrix representation of the boundary value problem.

It may be remarked that for other values of γ it appeared not possible to derive simple expressions for $e(n,m)$ and $\lambda(n,m)$ as given above.

We desire to determine the relaxation factor Ω such that the P-condition number of $L'(1)$ is as small as possible. At first sight, one should choose Ω according to the theory of Young, that is one minimizes the condition number of the operator $L'(1)$ defined by (2.12). In Coolen and van der Houwen it was shown that the theory of Young yields

$$(3.6) \quad \Omega = 2\gamma^{-1}(1 - h\sqrt{2\gamma^{-1}}), \quad \gamma = 1, 2.$$

However, it will be shown here that for $\gamma = 1$ an exact analysis of (3.4) yields a slightly better value for Ω . Of course, when dealing with a more general problem than the model problem, one should take the value of Ω prescribed by the Young theory.

Let us suppose that Ω is close to but less than 2 and, temporally, let us assume that v and μ are continuous variables.

It is readily seen that the stationary points of $\lambda(n,m)$ satisfy the equations

$$\mu(2 + \Omega - 2\Omega v\mu)(4 - 4\Omega v\mu + \Omega^2 v^2) + 2\Omega(\Omega v - 2\mu)(1 - v\mu)(1 - \Omega v\mu) = 0,$$

$$v(2 + \Omega - 2\Omega v\mu)(4 - 4\Omega v\mu + \Omega^2 v^2) - 4\Omega v(1 - v\mu)(1 - \Omega v\mu) = 0.$$

For $v \neq 0$ the stationary points of $\lambda(n,m)$ are situated on the line

$$(3.7) \quad \mu = \frac{1}{4} \Omega v.$$

Therefore, the extrema of $\lambda(n,m)$ are reached at points of the curve (3.7), the μ -axis, or at points of the boundary of the (v,μ) -domain, i.e. the line segments $v = \pm \cos h$, $-\cos h \leq \mu \leq \cos h$ and $\mu = \pm \cos h$, $-\cos h \leq v \leq \cos h$ (see figure 3.1).

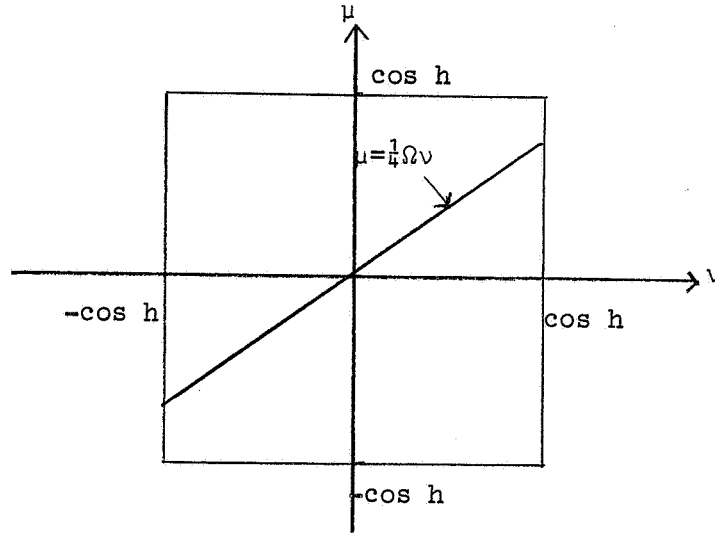


fig. 3.1

Along the curve $\mu = \frac{1}{4} \Omega v$ λ reaches a maximal value for $v = 0$ and along the μ -axis we have $\lambda = 2h^{-2}$. Therefore, we are only concerned with the values of λ at the boundary. Thus, returning to integer values for n and m , we have to consider the values:

$$(3.8) \quad \lambda\left(\frac{\pi}{2h}, 1\right), \lambda(n, 1) = \lambda\left(\frac{\pi}{h} - n, \frac{\pi}{h} - 1\right), \lambda(1, m) = \lambda\left(\frac{\pi}{h} - 1, \frac{\pi}{h} - m\right),$$

where $n, m = 1, 2, \dots, \frac{\pi}{h} - 1$.

In figure 3.2 the behaviour of the functions $\lambda(n, 1)$ and $\lambda(1, m)$ is illustrated for $\Omega \sim 2$.

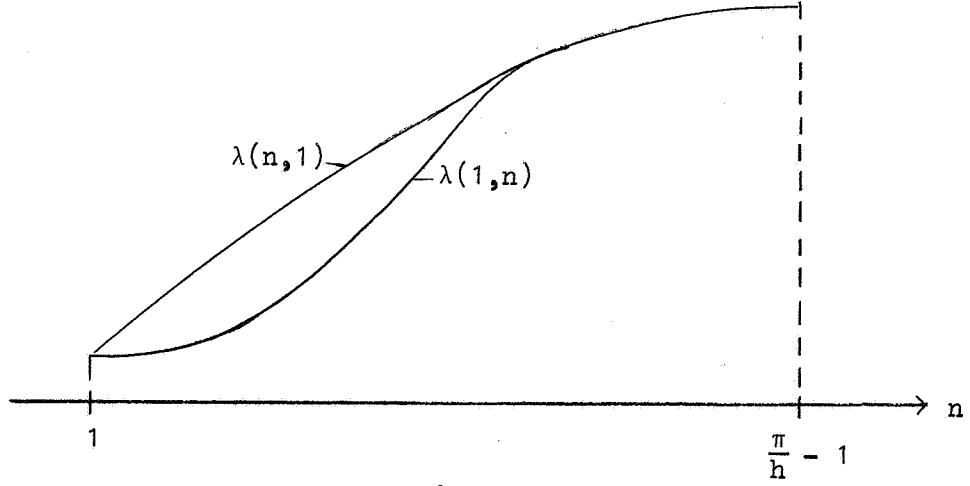


fig. 3.2

From this figure it may be concluded that the extrema of $\lambda(n, m)$ are bounded by the values of

$$\lambda\left(\frac{\pi}{2h}, 1\right) = 2h^{-2},$$

$$(3.9) \quad \lambda(1, 1) = 2h^{-2} \frac{2(1 - \cos^2 h)(2 - \Omega \cos^2 h)}{4 - 4\Omega \cos^2 h + \Omega^2 \cos^2 h} \sim 2h^{-2} \frac{2\eta(\epsilon + 2\eta)}{\epsilon^2 + 4\eta},$$

$$\lambda\left(1, \frac{\pi}{h} - 1\right) = 2h^{-2} \frac{2(1 + \cos^2 h)(2 + \Omega \cos^2 h)}{4 - 4\Omega \cos^2 h + \Omega^2 \cos^2 h} \sim 2h^{-2},$$

where $\eta = 1 - \cos^2 h \sim h^2$ and $\epsilon = 2 - \Omega$. The P-condition number is approximated by

$$(3.10) \quad P(L'(1)) \sim \frac{\epsilon^2 + 4\eta}{2\eta(\epsilon + 2\eta)}.$$

This expression is minimized by $\epsilon = 2\sqrt{\eta}$. Hence we find for Ω the approximate value (compare formula (3.6) for $\gamma = 1$)

$$(3.11) \quad \Omega \sim 2 - \epsilon = 2 - 2\sqrt{\eta} = 2 - 2 \sin h \sim 2 - 2h.$$

This value of Ω yields the condition number

$$(3.12) \quad P(L'(1)) \sim \frac{2}{\sin h + \sin^2 h} \sim 2h^{-1}.$$

A slightly larger condition number is found for the SOR value of Ω defined by (3.6), namely

$$(3.13) \quad P(L'(1)) \sim \frac{3}{2} \sqrt{2} h^{-1} \sim 2.1 h^{-1}.$$

The average rates of convergence of Richardson's method corresponding to (3.12) and (3.13) respectively, are

$$(3.14) \quad R(K) \sim \sqrt{2h} - \frac{\ln 2}{K} \sim 1.41 \sqrt{h} - \frac{0.693}{K},$$

$$(3.15) \quad R(K) \sim 2\sqrt{\frac{\sqrt{2}}{3} h} - \frac{\ln 2}{K} \sim 1.31 \sqrt{h} - \frac{0.693}{K}.$$

In order to compare the new method with other iterative processes, for instance Young's method, one must bear in mind that the method described above is twice as laborious per iteration. Thus, comparing (3.14) with the rate of convergence of Young's SOR method, i.e. $2\sqrt{2} h$ for large values of K (see Coolen and van der Houwen [1968]), we may conclude that the new method becomes faster if

$$\frac{1}{2} \sqrt{2h} > 2\sqrt{2} h,$$

or equivalently

$$(3.16) \quad h < \frac{1}{16} = .0625, \quad N > 50.$$

Thus, only for rather small values of h the new method will be advantageous.

4. Numerical results for the model problem

In this section the results are given of a number of experiments with the model problem on the EL X8 computer at the Mathematical Centre at Amsterdam.

In order to check the accuracy of the numerical solution we have chosen

$$(4.1) \quad F(x,y) = -2 \exp(x+y), \quad 0 < x < \pi, \quad 0 < y < \pi,$$

$$(4.2) \quad U(x,y) = \exp(x+y), \quad x = 0, \pi, \quad y = 0, \pi.$$

The inhomogeneous term (4.1) and the boundary function (4.2) give rise to the analytical solution

$$(4.3) \quad U(x,y) = \exp(x+y), \quad 0 \leq x \leq \pi, \quad 0 \leq y \leq \pi.$$

The iteration process was started by the following initial approximation:

$$(4.4) \quad u_0(j,1) = \left(\frac{1}{j} + \frac{1}{1} + \frac{1}{N-j} + \frac{1}{N-1} \right)^{-1} \left(\frac{u(j,0)}{1} + \frac{u(0,1)}{j} + \frac{u(N,1)}{N-j} + \frac{u(j,N)}{N-1} \right).$$

This net function equals at each point the average of its boundary neighbours (see figure 4.1).

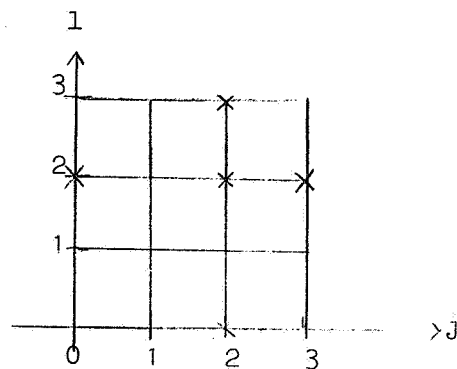


fig. 4.1

The method was applied for $N = 10$, $N = 20$ and $N = 40$. The following table contains the number of iterates K , the estimates a and b for the least and greatest eigenvalue of $L^v(1)$, the theoretical average rate of convergence $R(K)$ defined by formula (3.14), the experimental average rate of convergence $R^*(K)$ defined by

$$(4.5) \quad R^*(K) = -\frac{1}{K} \ln \left[\frac{\|L^v(1)u_K^*\|}{\|L^v(1)u_0\|} \right],$$

where u_K^* is the numerical net function and $\| \cdot \|$ denotes the maximum norm. Further we have listed the relative precision $P_r(K)$ after K iterates, defined by

$$(4.6) \quad P_r(K) = \left\| \frac{U - u_K^*}{U} \right\|$$

and the asymptotic relative precision $P_r(\infty)$.

TABLE 4.1

N	a	b	K	R(K)	R*(K)	P _r (K)	P _r (∞)
10	3.9780	23.5009	3	.6453	.8469	91 %	10.65%
			6	.7591	.8828	13 %	
			9	.7976	.8757	10.31%	
			12	.8169	.8645	10.63%	
			15	.8284	.8843	10.65%	
20	7.1301	87.4078	3	.3662	.5241	120 %	2.63%
			6	.4722	.6196	22 %	
			9	.5105	.6199	5.4 %	
			12	.5298	.5877	2.81%	
			18	.5491	.6011	2.63%	
40	13.7313	337.0450	5	.2740	.4173	108 %	.66%
			10	.3400	.4164	17 %	
			15	.3631	.4226	3.2 %	
			20	.3746	.4148	.78%	
			30	.3862	.4157	.65%	

In practice it is important to determine the value of K beforehand; one would desire that K is such that the errors due to discretization and iteration are comparable, i.e.

$$(4.7) \quad ||U - u|| \sim ||u - u_K||.$$

We can give a lower bound for K . According to Gerschgorin we have

$$(4.8) \quad ||U - u|| = ch^2,$$

where c is a constant. Further we have

$$(4.9) \quad ||u - u_K|| = 2\theta \exp(-2K\sqrt{\frac{a}{b}}) \cdot ||u - u_0||,$$

where $0 < \theta \leq 1$.

From (4.7), (4.8) and (4.9) it follows that

$$K \sim \sqrt{\frac{b}{a}} \left[\ln h^{-1} + \frac{1}{2} \ln\left(\frac{2\theta ||u - u_0||}{c}\right) \right].$$

Thus, using (3.12), we obtain

$$(4.10) \quad K \geq 1.41 \frac{\ln h^{-1}}{\sqrt{h}}.$$

For very small values of h ($\ln h^{-1} \gg \frac{1}{2} \ln(2\theta ||u - u_0||/c)$) the equality sign approximately holds. For instance, the case $N = 40$ requires at least 14 iterations.

From table 4.1 one may conclude that the numerical rates of convergence are in agreement with the theoretical rates of convergence predicted by the theory of the preceding sections.

5. References

- Coolen, T.M.T.,
Houwen, P.J. van der [1968]

Forsythe, G.E.,
Wasow, W.R. [1960]

Houwen, P.J. van der [1967a]

Houwen, P.J. van der [1968]

Stiefel, E.C. [1958]
- On the acceleration of Richardson's method IV. A non-symmetrical case. Report TW 109, Math. Centre, Amsterdam.
- Finite difference methods for partial differential equations. John Wiley & Sons, Inc., New York.
- On the acceleration of Richardson's method III. Applications. Report TW 108, Math. Centre, Amsterdam.
- Finite difference methods for solving partial differential equations. MC Tract 20, Math. Centre, Amsterdam.
- Kernel polynomials in linear algebra and their numerical applications. National Bureau of Standards Applied Mathematics Series 49, "Further Contributions to the Solution of Simultaneous Linear Equations and the Determination of Eigenvalues", U.S. Government Printing Office, Washington, D.C.