# Iterative procedures for solving finite-difference approximations to separable partial differential equations

*By* M. R. Osborne

We show that, for a class of separable partial differential equations of elliptic type in two independent variables, the optimum value of the over-relaxation parameter $\omega$ can be calculated when the eigenvalue of maximum modulus of a certain double eigenvalue problem is known. This generalizes the well-known case of Laplace's equation for a rectangle (for which explicit results are known), and provides a convenient method for the calculation of the optimum value of $\omega$. We extend some of our results to partial differential equations in three independent variables, and we also use our formalism to discuss the Peaceman–Rachford iteration.

## Introduction

In this paper we consider certain classes of iterative methods for solving finite-difference approximations to elliptic partial differential equations having the form

$$\frac{\partial^2 Q}{\partial x^2} + \frac{\partial^2 Q}{\partial y^2} + f(x)\frac{\partial Q}{\partial x} + g(y)\frac{\partial Q}{\partial y}$$
$$+ [p(x) + q(y)]Q = c(x, y) \quad (1)$$

with the boundary conditions given on a rectangle with sides parallel to the $x$ and $y$ axes, respectively. Our analysis will be worked for the case $Q$ prescribed on the boundary but applies (with obvious modifications) to the case $AQ + B\dfrac{\partial Q}{\partial n}$ prescribed, where $A$ and $B$ are constant on any one side. Typical equations to which our results are applicable are Laplace's equation for the rectangle and (in polar coordinates) the circular annulus, and Reynolds equation* for a rectangle.

By restricting our attention to rectangular regions and separable partial differential equations, we find that it is possible to represent the unknowns in our finite-difference equations as the components of a rectangular matrix (in contrast to the usual representation as the components of a vector), and that the finite-difference equations then take a very simple form [equation (5)]. This representation is given in section 1 where it is applied to analyze the iterative scheme of successive over-relaxation by points. Successive over-relaxation by lines is considered in section 2, and the Peaceman–Rachford iteration in section 3. Successive over-relaxation by points for three independent variables is considered in section 4. In section 5 we give a computational scheme for solving the eigenvalue problems derived in the preceding sections, and in section 6 we consider the relation between our method and that normally used to analyze the iterative procedures discussed here.

---

* This equation, which is important in the theory of film lubrication, has the form

$$\frac{\partial}{\partial x}\left[h^3(x)\frac{\partial P}{\partial x}\right] + h^3(x)\frac{\partial^2 P}{\partial y^2} = \frac{dh}{dx}.$$

## 1. Iteration by points

If we write $\phi_i$ for the vector whose components make up the $i$th line $(x = x_i, i = 1, 2, \ldots, n)$ of the solution, then the standard five-point difference approximation to (1) takes the form

$$h^{-2}\left(1 - \frac{hf_i}{2}\right)\phi_{i-1} + [G - (2h^{-2} - p_i)I]\phi_i$$
$$+ h^{-2}\left(1 + \frac{hf_i}{2}\right)\phi_{i+1} = C_i \quad (2)$$

where $G$ is an $(m \times m)$ tridiagonal matrix whose $j$th row is

$k^{-2}\left(1 - \dfrac{kg_j}{2}\right)$, $-2k^{-2} + q_j$, $k^{-2}\left(1 + \dfrac{kg_j}{2}\right)$, and where $h$ and $k$ are the mesh spacings in the $x$ and $y$ directions, respectively.

Writing $G - (2h^{-2} - p_i)I = L + X_i + U$ as the sum of a lower-triangular, a diagonal, and an upper-triangular matrix, we can represent the Young–Frankel over-relaxation scheme with iteration parameter $\omega$ as

$$h^{-2}\left(1 - \frac{hf_i}{2}\right)\phi_{i-1}^{(s+1)} + (L + \omega^{-1}X_i)\phi_i^{(s+1)}$$
$$+ [(1 - \omega^{-1})X_i + U]\phi_i^{(s)}$$
$$+ h^{-2}\left(1 + \frac{hf_i}{2}\right)\phi_{i+1}^{(s)} = C_i \quad (3)$$

where the upper suffix $s$ indicates the progress of the iteration.

To determine the rate of convergence of this iteration we require a knowledge of the eigenvalue of maximum modulus of the problem (as in Forsythe and Wasow, 1960, p. 214 and p. 247)

$$\lambda\left\{h^{-2}\left(1 - \frac{hf_i}{2}\right)\phi_{i-1} + (L + \omega^{-1}X_i)\phi_i\right\}$$
$$+ \{(1 - \omega^{-1})X_i + U\}\phi_i$$
$$+ h^{-2}\left(1 + \frac{hf_i}{2}\right)\phi_{i+1} = 0. \quad (4)$$

93

Before reducing (4) we note that the left-hand side of equation (2) can be described by an operator combining together consecutive $\phi_i$ plus an operation $G$ on the individual $\phi_i$. If we write $\Phi$ for the $m \times n$ matrix with columns $\phi_i$ then (2) becomes

$$G\Phi + \Phi F = C \qquad (5)$$

where $F$ is the $(n \times n)$ tridiagonal matrix whose $i$th row is

$$h^{-2}\left(1 + \frac{hf_{i-1}}{2}\right), \; -2h^{-2} + p_i, \; h^{-2}\left(1 - \frac{hf_{i+1}}{2}\right).$$

The representation of the difference equation in this form has been given before by Bickley and McNamee (1960).

If we set
$$\begin{aligned} F &= L_1 + D_1 + U_1, \\ G &= L_2 + D_2 + U_2, \end{aligned} \qquad (6)$$

then equation (4) becomes

$$\begin{aligned} &[\lambda L_2 + (\omega^{-1}\lambda + (1 - \omega^{-1}))D_2 + U_2]\Phi \\ &+ \Phi[L_1 + (\lambda\omega^{-1} + (1 - \omega^{-1}))D_1 + \lambda U_1] = 0. \end{aligned} \qquad (7)$$

Equation (7) can be further transformed by using the result (due to Friedman; see Forsythe and Wasow, 1960, p. 249) that there exist diagonal matrices $A$ and $B$ such that

$$\lambda L_2 + \zeta D_2 + U_2 = \sqrt{\lambda} A^{-1}\left\{L_2 + \frac{\zeta}{\sqrt{\lambda}}D_2 + U_2\right\}A,$$

and $L_1 + \zeta D_1 + \lambda U_1 = \sqrt{\lambda} B\left\{L_1 + \frac{\zeta}{\sqrt{\lambda}}D_1 + U_1\right\}B^{-1}.$
$$(8)$$

If we now set $\overline{\Phi} = A\Phi B$ we have reduced our problem to that of finding the eigenvalues of

$$\{L_2 + \sigma D_2 + U_2\}\overline{\Phi} + \overline{\Phi}\{L_1 + \sigma D_1 + U_1\} = 0 \qquad (9)$$

where
$$\sigma = \frac{\omega^{-1}\lambda + (1 - \omega^{-1})}{\sqrt{\lambda}}. \qquad (10)$$

The values of $\sigma$ which satisfy (9) can also be characterized as the solutions of the double eigenvalue problem

$$\begin{aligned} [L_1 + \sigma D_1 + U_1 + \gamma I]v_1 &= 0, \\ [L_2 + \sigma D_2 + U_2 - \gamma I]v_2 &= 0. \end{aligned} \qquad (11)$$

This follows from a result of Bickley and McNamee (1960, p. 118), who showed that the eigenvalue problem

$$R\Phi + \Phi S = \lambda\Phi \qquad (12)$$

where $R$ and $S$ can be diagonalized by similarity transformations, has its eigenvalues in the form $\lambda = \lambda_1 + \lambda_2$, where $\lambda_1$ is an eigenvalue of $R$, and $\lambda_2$ is an eigenvalue of $S$ (the eigenmatrix being the dyad formed from the corresponding eigenvectors of $R$ and $S$). In particular there is a non-trivial solution with $\lambda = 0$ only if $\lambda_1$ is an eigenvalue of $R$ and $-\lambda_1$ is an eigenvalue of $S$. Thus (11) just expresses the condition for (9) to have a non-trivial solution. The condition that the matrices

operating on $\overline{\Phi}$ in equation (9) can be diagonalized by similarity transformations is certainly satisfied in many cases of interest. A sufficient condition is that $h$ and $k$ can be chosen small enough to make $h|f_i| < 2$, $k|g_j| < 2$, $h^2|p_i| < 2$, and $k^2|q_j| < 2$ for all $i$, $j$, for then $D_1$ and $D_2$ are negative definite while the off-diagonal elements in $G$ and $F$ are one-signed.

In special cases the equations (11) can be further simplified. Assuming that the conditions in the previous paragraph are satisfied, we have for each of the matrix operators in (9)

$$U + L = DVMV^{-1} \qquad (13)$$

where $V$ is the matrix of (right) eigenvectors of $U + L$, and $M$ is the diagonal matrix of the eigenvalues $\mu$ taken with respect to the density matrix $D$. Therefore (in an obvious notation) (9) becomes

$$\begin{aligned} &D_2 V_2\{\sigma I + M_2\}V_2^{-1}\overline{\Phi} \\ &+ \overline{\Phi}D_1 V_1\{\sigma I + M_1\}V_1^{-1} = 0. \end{aligned} \qquad (14)$$

*Case* 1 Let $D_1$ be a constant multiple ($a_1$ say) of the unit matrix $I$. Then $D_1$ commutes with $V_1$ and equation (14) can be written (setting $\overline{\overline{\Phi}} = \overline{\Phi}V_1$)

$$\{L_2 + \sigma D_2 + U_2\}\overline{\overline{\Phi}} + a_1\overline{\overline{\Phi}}\{\sigma I + M_1\} = 0. \qquad (15)$$

From (15) we see that the $i$th column of $\overline{\overline{\Phi}}$ (say $\theta_i$) satisfies the equation

$$\begin{aligned} \{L_2 + a_1\mu_{1i}I + U_2 + \sigma(D_2 + a_1 I)\}\theta_i &= 0, \qquad (16) \\ i &= 1, 2, \ldots, n. \end{aligned}$$

Equation (16) is an eigenvalue problem which determines the possible values of $\sigma$ in this case.

*Case* 2 Let both $D_1$ and $D_2$ be constant multiples of the unit matrix ($D_2 = a_2 I$), and let $\overline{\overline{\Phi}} = V_2^{-1}\overline{\Phi}V_1$; then equation (14) can be written

$$a_2\{\sigma I + M_2\}\overline{\overline{\Phi}} + a_1\overline{\overline{\Phi}}\{\sigma I + M_1\} = 0. \qquad (17)$$

In this case $\theta_i$ satisfies

$$\{\sigma(a_1 + a_2)I + (a_2 M_2 + a_1\mu_{1i}I)\theta_i = 0, \qquad (18)$$

and as the matrices within the brackets are diagonal we have

$$\begin{aligned} -\sigma_{ij} &= \frac{a_2\mu_{2j} + a_1\mu_{1i}}{a_1 + a_2}, \\ i &= 1, 2, \ldots, n; j = 1, 2, \ldots, m. \qquad (19) \end{aligned}$$

A special case of equation (19) is well known. This is the case of Laplace's equation for a rectangle (Forsythe and Wasow, 1960, p. 255). Also in this reference (pp. 250–57) it is shown that the optimum value of $\omega$ depends only on the maximum value of $\sigma$. A method of calculating this maximum value of $\sigma$ from equation (11) is given in section 5. No discussion of equation (11) other than that necessary to characterize the maximum value of $\sigma$ and its attendant value of $\gamma$ is attempted.

94

## 2. Block relaxation by lines

The iterative scheme for successive over-relaxation by lines (Forsythe and Wasow, 1960, pp. 266–71) can be written

$$h^{-2}\Big(1 - \frac{hf_i}{2}\Big)\phi_{i-1}^{(s+1)} + \omega^{-1}[G - (2h^{-2} - p_i)I]\phi_i^{(s+1)}$$

$$+ (1 - \omega^{-1})[G - (2h^{-2} - p_i)I]\phi_i^{(s)}$$

$$+ h^{-2}\Big(1 + \frac{hf_i}{2}\Big)\phi_i^{(s)} = C_i, \quad (20)$$

and the eigenvalue problem determining the rate of convergence is

$$[\lambda\omega^{-1} + (1 - \omega^{-1})]G\Phi$$

$$+ \Phi\{\lambda U_1 + [\lambda\omega^{-1} + (1 - \omega^{-1})]D_1 + L_1\} = 0. \quad (21)$$

Applying a Friedman transformation to the post-multiplier in (21) gives us

$$\frac{\lambda\omega^{-1} + (1 - \omega^{-1})}{\sqrt{\lambda}}G\overline{\Phi} + \overline{\Phi}$$

$$\Big\{L_1 + \frac{\lambda\omega^{-1} + (1 - \omega^{-1})}{\sqrt{\lambda}}D_1 + U_1\Big\} = 0 \quad (22)$$

so that the eigenvalues are determined when the values of $\sigma$ for which the set of equations

$$\sigma G\overline{\Phi} + \overline{\Phi}\{L_1 + \sigma D_1 + U_1\} = 0 \quad (23)$$

have non-trivial solutions $\overline{\Phi}$ are known.

As before we have $\sigma$ determined by the double eigenvalue problem

$$\{\sigma G - \nu I\}v_2 = 0,$$

$$\{L_1 + \sigma D_1 + U_1 + \nu I\}v_1 = 0. \quad (24)$$

In the case $D_1 = a_1I$ we have, substituting $L_1 + U_1 = D_1 V_1 M_1 V_1^{-1}$ into the second equation (24),

$$V_1\{a_1 M_1 + (a_1\sigma + \nu)I\}V_1^{-1}v_1 = 0,$$

and this gives the eigenvalues $\nu$ in terms of $\sigma$ as

$$- \nu_i = a_1(\mu_{1i} + \sigma), \quad i = 1, 2, \ldots, n. \quad (25)$$

Substituting for $\nu$ from (25) into the first equation of (24) we have

$$\{\sigma(G + a_1I) + a_1\mu_{1i}I\}v_2 = 0. \quad (26)$$

Finally, if $D_2 = a_2I$ we have

$$G = L_2 + D_2 + U_2 = V_2\{a_2I + a_2M_2\}V_2^{-1},$$

and the eigenvalues $\sigma$ of equation (26) are given by

$$\sigma_{ij} = \frac{- a_1\mu_{1i}}{a_1 + a_2 + a_2\mu_{2j}}. \quad (27)$$

Again (27) is well known for the case of Laplace's equation for a rectangle (Forsythe and Wasow, 1960, p. 270).

## 3. The Peaceman–Rachford iteration

Our methods can also be applied to analyze the implicit alternating direction methods for solving (5), and here we consider the Peaceman–Rachford method (Forsythe and Wasow, 1960, pp. 272–82). The iteration scheme for the error matrix is

$$\Phi^{(s-\frac{1}{2})} = \Phi^{(s-1)} - \alpha_s(G\Phi^{(s-1)} + \Phi^{(s-\frac{1}{2})}F),$$

$$\Phi^{(s)} = \Phi^{(s-\frac{1}{2})} - \alpha_s(G\Phi^{(s)} + \Phi^{(s-\frac{1}{2})}F). \quad (28)$$

Eliminating $\Phi^{(s-\frac{1}{2})}$ gives

$$\Phi^{(s)} = \{I + \alpha_s G\}^{-1}\{I - \alpha_s G\}\Phi^{(s-1)}\{I + \alpha_s F\}^{-1}\{I - \alpha_s F\},$$

$$= \prod_{r=1}^{s}\{I + \alpha_r G\}^{-1}\{I - \alpha_r G\}\Phi^{(0)}\prod_{r=1}^{s}\{I + \alpha_r F\}^{-1}$$

$$\{I - \alpha_r F\}. \quad (29)$$

Introducing the spectral decompositions of these right and left multiplying matrices we may write (29) as

$$\Phi^{(s)} = V_2 D_2 V_2^{-1}\Phi^{(0)}V_1 D_1 V_1^{-1} \quad (30)$$

where $D_1$ and $D_2$ are the diagonal matrices of the eigenvalues. Setting $\overline{\Phi} = V_2^{-1}\Phi V_1$ we have

$$\overline{\Phi}^{(s)} = D_2\overline{\Phi}^{(0)}D_1, \quad (31)$$

and separating out the $i$th column $\phi_i^{(s)}$ we see that it satisfies

$$\phi_i^{(s)} = d_{1i}D_2\phi_i^{(0)} \quad (32)$$

where $d_{1i}$ is the $i$th diagonal element of $D_1$. Thus the rate at which $\overline{\Phi}^{(s)}$ tends to zero is determined by the maximum value of $|d_{2j}d_{1i}|$. If we call the eigenvalues of $F\lambda_1, \ldots, \lambda_n$, and those of $G\mu_1, \ldots, \mu_m$

then

$$d_{1i}d_{2j} = \prod_{r=1}^{s}\frac{1 - \alpha_r\lambda_i}{1 + \alpha_r\lambda_i}\frac{1 - \alpha_r\mu_j}{1 + \alpha_r\mu_j}. \quad (33)$$

Equation (33) is well known for Laplace's equation in a rectangle. An alternative derivation of (33) is quoted in Martin and Tee (1961).

## 4. Separable equations in three dimensions

Here we indicate the way in which our results extend to problems with more than two independent variables. We consider the partial differential equation

$$\nabla^2\phi + f(x)\frac{\partial\phi}{\partial x} + g(y)\frac{\partial\phi}{\partial y} + e(z)\frac{\partial\phi}{\partial z}$$

$$+ (p(x) + q(y) + r(z))\phi = c(x, y, z) \quad (34)$$

with the boundary conditions $\phi$ specified on the surface of a rectangular box with faces parallel to the coordinate planes. The standard seven-point finite-difference approximation to this equation can be written [using equation (5)]

$$\alpha^{-2}\Big(1 - \frac{\alpha e_k}{2}\Big)\Phi_{k-1} - (2\alpha^{-2} - r_k)\Phi_k$$

$$+ \alpha^{-2}\Big(1 + \frac{\alpha e_k}{2}\Big)\Phi_{k+1} + G\Phi_k + \Phi_k F = C_k \quad (35)$$

where $\Phi_k$ is the matrix of solution points $\phi_{ijk}$ for fixed $k$ (constant $z$), $k = 1, 2, \ldots, t$, and $\alpha$ is the spacing of the mesh in the $z$ direction.

The formula for the iterative solution of (35) by successive over-relaxation by points, assuming each plane $\Phi_k$ is relaxed before proceeding to $\Phi_{k+1}$, is

$$\alpha^{-2}\left(1 - \frac{\alpha e_k}{2}\right)\Phi_{k-1}^{(s)} - \omega^{-1}(2\alpha^{-2} - r_k)\Phi_k^{(s)}$$
$$- (1 - \omega^{-1})(2\alpha^{-2} - r_k)\Phi_k^{(s-1)}$$
$$+ \alpha^{-2}\left(1 + \frac{\alpha e_k}{2}\right)\Phi_{(k+1)}^{(s-1)} + [L_2 + \omega^{-1}D_2]\Phi_k^{(s)}$$
$$+ \Phi_k^{(s)}[U_1 + \omega^{-1}D_1] + [(1 - \omega^{-1})D_2 + U_2]\Phi_k^{(s-1)}$$
$$+ \Phi_k^{(s-1)}[L_1 + (1 - \omega^{-1})D_1] = C_k. \tag{36}$$

The rate of convergence of the iteration (36) depends on the eigenvalue of maximum modulus of the problem

$$\lambda\alpha^{-2}\left(1 - \frac{\alpha e_k}{2}\right)\Phi_{k-1} - (\lambda\omega^{-1} + (1 - \omega^{-1}))$$
$$(2\alpha^{-2} - r_k)\Phi_k + \alpha^{-2}\left(1 + \frac{\alpha e_k}{2}\right)\Phi_{k+1}$$
$$+ [\lambda L_2 + (\omega^{-1}\lambda + (1 - \omega^{-1}))D_2 + U_2]\Phi_k$$
$$+ \Phi_k[L_1 + (\lambda\omega^{-1} + (1 - \omega^{-1}))D_1 + \lambda U_1] = 0. \tag{37}$$

If we now apply the transformations (8) to equation (37) we have, writing $\overline{\Phi}_k = A\Phi_k B$,

$$\lambda\alpha^{-2}\left(1 - \frac{\alpha e_k}{2}\right)\overline{\Phi}_{k-1} - \sqrt{\lambda}\sigma(2\alpha^{-2} - r_k)\overline{\Phi}_k$$
$$+ \alpha^{-2}\left(1 + \frac{\alpha e_k}{2}\right)\overline{\Phi}_{k+1} + \sqrt{\lambda}\{[L_2 + \sigma D_2 + U_2]\overline{\Phi}_k$$
$$+ \overline{\Phi}_k[L_1 + \sigma D_1 + U_1]\} = 0. \tag{38}$$

Now let the normalized, right and left, bi-orthogonal eigenvectors, and the corresponding eigenvalues of $L_1 + \sigma D_1 + U_1$ be $v_{1q}$, $v_{1q}^*$, and $\mu_{1q}$ respectively, $q = 1, 2, \ldots, n$. Also let the normalized, right and left, bi-orthogonal eigenvectors and the corresponding eigenvalues of $L_2 + \sigma D_2 + U_2$ be $v_{2p}$, $v_{2p}^*$, and $\mu_{2p}$, respectively; $p = 1, 2, \ldots, m$. The conditions stated following equation (12) are sufficient to guarantee the existence of these, and also to guarantee that we can represent $\Phi_k$, $k = 1, \ldots, t$ in the form

$$\Phi_k = \sum_{p, q} A_{kpq} v_{2P} v_{1q}^{*T}. \tag{39}$$

Using the bi-orthogonality of the eigenvectors we have

$$A_{kpq} = v_{2p}^{*T}\Phi_k v_{1q}. \tag{40}$$

If we now substitute (39) into (38) we obtain, after premultiplying by $v_{2p}^{*T}$, postmultiplying by $v_{1q}$, and using the bi-orthogonality of the eigenvectors

$$\lambda\alpha^{-2}\left(1 - \frac{\alpha e_k}{2}\right)A_{(k-1)pq} + \sqrt{\lambda}(-\sigma(2\alpha^{-2} - r_k)$$
$$+ \mu_{2p} + \mu_{1q})A_{kpq} + \alpha^{-2}\left(1 + \frac{\alpha e_k}{2}\right)A_{(k+1)pq} = 0 \tag{41}$$

If we denote by $A$ the vector with components $A_{kpq}$, $k = 1, 2, \ldots, t$, and if we define $H$ as the tridiagonal matrix whose $k$th row is

$$\alpha^{-2}\left(1 - \frac{\alpha e_k}{2}\right), \ -(2\alpha^{-2} - r_k), \ \alpha^{-2}\left(1 + \frac{\alpha e_k}{2}\right),$$

and set $H = L_3 + D_3 + U_3$, then equation (41) is clearly the $k$th row of the matrix equation

$$\{\lambda L_3 + \sqrt{\lambda}(\sigma D_3 + (\mu_{2p} + \mu_{1q})I) + U_3\}A = 0. \tag{42}$$

Applying a Friedman transformation to equation (42) gives us

$$\{L_3 + \sigma D_3 + U_3 + (\mu_{2p} + \mu_{1q})I\}A = 0. \tag{43}$$

From equation (43) we see that $\sigma$ is determined by the solutions of a triple eigenvalue problem which is the natural extension of equation (11). We may write this as

$$\{L_1 + \sigma D_1 + U_1 - \lambda I\}V_1 = 0,$$
$$\{L_2 + \sigma D_2 + U_2 - \mu I\}V_2 = 0,$$
and $\{L_3 + \sigma D_3 + U_3 - \nu I\}V_3 = 0,$
where $\lambda + \mu + \nu = 0. \tag{44}$

The reduction of (44) in the case where some or all of $D_1$, $D_2$, $D_3$ are scalar multiples of the corresponding unit matrices follows as before.

## 5. Calculation of the maximum value of σ

In this section we suggest a method for solving the double eigenvalue problem given by equation (11). Our procedure is a straightforward application of the initial value techniques often used in solving finite-difference approximations to the eigenvalue problems of ordinary differential equations (see, for example, Fox, 1960). The method can be applied with obvious modifications to equations (24) and (44).

We write equation (11) in the form

$$[L_1 + \sigma D_1 + U_1 - \mu_1 I]v_1 = 0,$$
$$[L_2 + \sigma D_2 + U_2 - \mu_2 I]v_2 = 0,$$
$$\mu_1 + \mu_2 = 0. \tag{45}$$

Taking trial values of $\sigma$, $\mu_1$, and $\mu_2$ we fix the scale of $v_1$ and $v_2$ by setting $(v_1)_1 = (v_2)_1 = 1$, and we use each matrix equation as a three-term recurrence for the components of the corresponding $v$. By this means we solve the equations

$$[L_1 + \sigma D_1 + U_1 - \mu_1 I]v_1 = \beta_1 e_n,$$
and $[L_2 + \sigma D_2 + U_2 - \mu_2 I]v_2 = \beta_2 e_m \tag{46}$

where $e_n$ is the $n$-dimensional unit vector with one in the $n$th place, and $e_m$ is defined similarly. The required values of $\sigma$ and $\mu$ can now be characterized as solutions of the system

$$\beta_1(\sigma, \mu_1) = 0,$$
$$\beta_2(\sigma, \mu_2) = 0,$$
$$\mu_1 + \mu_2 = 0. \tag{47}$$

96

Corrections to our trial values of $\sigma$ and $\mu$ can now be calculated by applying Newton's method to (47). This gives

$$\frac{\partial \beta_1}{\partial \sigma}\Delta\sigma + \frac{\partial \beta_1}{\partial \mu_1}\Delta\mu_1 + \beta_1(\sigma, \mu_1) = 0,$$

$$\frac{\partial \beta_2}{\partial \sigma}\Delta\sigma + \frac{\partial \beta_2}{\partial \mu_2}\Delta\mu_2 + \beta_2(\sigma, \mu_2) = 0,$$

$$\mu_1 + \mu_2 + \Delta\mu_1 + \Delta\mu_2 = 0. \tag{48}$$

The partial derivatives of $\beta_1$ can be computed from the the final rows of variational equations

$$[L_1 + \sigma D_1 + U_1 - \mu_1 I]\frac{\partial v_1}{\partial \sigma} = \frac{\partial \beta_1}{\partial \sigma}e_n - D_1 v_1,$$

$$[L_1 + \sigma D_1 + U_1 - \mu_1 I]\frac{\partial v_1}{\partial \mu_1} = \frac{\partial \beta_1}{\partial \mu_1}e_n + v_1, \tag{49}$$

with the initial conditions $\left(\dfrac{\partial v_1}{\partial \sigma}\right)_1 = \left(\dfrac{\partial v_1}{\partial \mu_1}\right)_1 = 0$, and

there is a similar set of equations for the partial derivatives of $\beta_2$. We mention that by making a different choice of the scale in the calculation of $v_1$ and $v_2$ it is possible to by-pass the solution of the variational equations in calculating the partial derivatives of $\beta_1$ and $\beta_2$. However, this may lead to a computation which is numerically unstable (Osborne, 1962).

We now discuss further some properties of equation (11). We assume that the conditions stated following equation (12) are satisfied and, in addition, that the matrices $F$ and $G$ are negative definite (the most usual case). If we set $\sigma = 1$ in the eigenvalue problems

$$[L_1 + \sigma D_1 + U_1 - \mu_1 I]v_1 = 0,$$
$$\text{and} \qquad [L_2 + \sigma D_2 + U_2 - \mu_2 I]v_2 = 0 \tag{50}$$

it follows from the above assumptions that their maximum eigenvalues $\mu_1$ and $\mu_2$ are negative. We will show that as $\sigma$ is decreased these eigenvalues increase. It then follows that there is a value of $\sigma$ for which the maximum eigenvalue of one of the matrices is sufficiently positive to cancel the maximum eigenvalue of the other. This is the value of $\sigma$ required in the calculation of the optimum value of $\omega$.

The proof is straightforward, and we consider here $\mu_1$ the maximum eigenvalue of the first problem. As $\sigma$ is decreased we follow the curve $\beta_1(\sigma, \mu_1) = 0$. Using equation (48) we see that our result holds if $\partial\beta_1/\partial\sigma$ and $\partial\beta_1/\partial\mu_1$ have the same sign on this curve. To show this we use that the matrix on the left-hand side of equation (49) is singular on $\beta_1(\sigma, \mu_1) = 0$, so that the right-hand sides must be orthogonal to the eigenvector of the transposed matrix. To derive this vector we use that the tridiagonal matrix $L_1 + U_1$ has its off-diagonal elements onesigned, and so can be made symmetric by premultiplying by a diagonal matrix $X$ with positive ele-

ments. The desired eigenvector is $Xv_1$ giving the orthogonality conditions

$$v_1^T X\{\partial\beta_1/\partial\sigma\, e_n - D_1 v_1\} = 0,$$
$$\text{and} \qquad v_1^T X\{\partial\beta_1/\partial\mu_1 e_n + v_1\} = 0. \tag{51}$$

$$\text{Thus} \qquad \frac{\partial\beta_1}{\partial\sigma}\bigg/\frac{\partial\beta_1}{\partial\mu_1} = -\,v_1^T X D_1 v_1/v_1^T X v_1 \tag{52}$$

and the desired result follows from this as all the elements of $D_1$ are negative.

Another useful property of the solutions to equation (11) can be deduced by noting that the principal minors of the matrices occurring in equation (50) form Sturm sequences with respect to $\mu$ for fixed $\sigma$. It follows that the eigenvectors $v_1$ and $v_2$ associated with the maximum eigenvalues $\mu_1$ and $\mu_2$ have the property that all their components have the same sign. This suggests that a suitable choice of $\sigma$ and $\mu$ to start our iteration would be one which leads to positive vectors $v_1$ and $v_2$ (this property can also be used to verify that the solution produced by our iteration is the correct one). A choice compatible with the assumptions made about $F$ and $G$ in this section is $\sigma = 1$, $\mu_1 = \mu_2 = 0$.

To illustrate the calculation consider the case of Laplace's equation for which explicit results are available. Taking $h = k = 1$, $m = 19$, $n = 5$, $\sigma = 1$, and $\mu_1 = \mu_2 = 0$ we found the corrections to $\sigma$ and $\mu_1$ given in the Table below.

| $\Delta\sigma$ | $\Delta\mu_1$ |
|---|---|
| $-0\cdot0466\ 1654$ | $-0\cdot0781\ 9549$ |
| $-0\cdot0214\ 9499$ | $-0\cdot0354\ 4978$ |
| $-0\cdot0048\ 0259$ | $-0\cdot0076\ 7722$ |
| $-0\cdot0002\ 2850$ | $-0\cdot0003\ 3983$ |
| $-0\cdot0000\ 0052$ | $-0\cdot0000\ 0062$ |

This gives $\sigma = 0\cdot926\ 857$, $\mu_1 = -0\cdot121\ 663$ agreeing

to all figures with the values of $\frac{1}{2}\left(\cos\dfrac{\pi}{6} + \cos\dfrac{\pi}{20}\right)$, and

$\cos\dfrac{\pi}{6} - \cos\dfrac{\pi}{20}$ taken from Chambers' 6-figure tables.

## 6. Some connections with the general theory of over-relaxation

In the iterative procedures considered in this paper we may represent the unknowns (*a*) as a column vector, and (*b*) as a rectangular matrix. The representation (*b*) has an advantage for our purposes as separability of the partial differential equation is mirrored by independent operations on the rows and columns of the solution matrix, as in equations (5) and (29) for example. The representation (*a*) has the merit of much greater generality as it applies to non-rectangular regions.

In the general theory of successive over-relaxation (S.O.R.) a key result is that the components of the solution vector can be so arranged that the matrix of the set of linear equations is block-tridiagonal, with the additional property that the blocks on the leading diagonal are diagonal matrices. After this has been

H

done it is fairly easy to show that the study of S.O.R. can be reduced to a study of a simpler iteration called the *method of simultaneous displacements* (Forsythe and Wasow, 1960, pp. 247–50). The essential steps in this development are paralleled in our work in the passage from equation (7) to equations (9) and (10). Equation (9) is the eigenvalue problem for the method of simultaneous displacements translated into our formalism, and equation (10) which relates the eigenvalues of the S.O.R. iteration with the method of simultaneous displacements is the same in both treatments.

Certainly some at least of the results obtained in this paper can also be obtained by applying known theorems to the matrix of the eigenvalue problem for the method of simultaneous displacements. As example we give a derivation of equation (16), that for (24) being similar.

We take as our starting point equation (9)

$$[L_2 + \sigma D_2 + U_2]\Phi + \Phi[L_1 + \sigma D_1 + U_1] = 0.$$

We assume that the elements of the $i$th row of $L_1 + \sigma D_1 + U_1$ are $1_i^{(1)}$, $\sigma d_i^{(1)}$, and $u_i^{(1)}$, and that the $i$th column of $\Phi$ is the vector $Q_i$. With this notation (9) implies the equations

$$1_{i+1}^{(1)} Q_{i+1} + [L_2 + \sigma(D_2 + d_i^{(1)}I) + U_2]Q_i + u_{i-1}^{(1)} Q_{i-1} = 0,$$
$$i = 1, 2, \ldots, n. \quad (53)$$

Assuming that $D_2 = a_2 I$ we may write (53) as

$$\frac{1_{i+1}^{(1)}}{a_2 + d_i^{(1)}} Q_{i+1} + \left\{ \frac{L_2 + U_2}{a_2 + d_i^{(1)}} + \sigma I \right\} Q_i + \frac{u_{i-1}^{(1)}}{a_2 + d_i^{(1)}} Q_{i-1} = 0,$$
$$i = 1, 2, \ldots, n. \quad (54)$$

If we now define an *mn*-rowed vector by the equation

$$v = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ Q_{i-1} \\ Q_i \\ \cdot \\ \cdot \end{bmatrix} \quad (55)$$

then we see from equation (54) that $-\sigma$ is an eigenvalue, and $v$ an eigenvector of the block tridiagonal matrix

$$A = \begin{bmatrix} \dfrac{L_2 + U_2}{a_2 + d_1^{(1)}}, & \dfrac{1_2^{(1)}}{a_2 + d_1^{(1)}}I, & \cdots \\ \dfrac{u_1^{(1)}}{a_2 + d_2^{(1)}}I, & \dfrac{L_2 + U_2}{a_2 + d_2^{(1)}}, & \dfrac{1_3^{(1)}}{a_2 + d_2^{(1)}}I \\ & \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot \\ & \dfrac{u_{n-1}^{(1)}}{a_2 + d_n^{(1)}}I, & \dfrac{L_2 + U_2}{a_2 + d_n^{(1)}} \end{bmatrix}. \quad (56)$$

To calculate the eigenvalues of $A$ we apply a result due to Afriat (Afriat, 1954) which may be stated as follows. Let an $mn \times mn$ matrix $M$ have the property that it can be partitioned into blocks $M_{ij} = R_{ij}(B)$, $i, j = 1, 2, \ldots, n$, where $R_{ij}(B)$ is a rational function of the $m \times m$ matrix $B$, and let $\lambda$ be an eigenvalue of $B$, then the eigenvalues of the $n \times n$ matrix with elements $R_{ij}(\lambda)$ are eigenvalues of $M$.

As $D_2 = a_2 I$ is a multiple of the unit matrix, we have the eigenvalues of $L_2 + U_2$ given by equation (13) in the form $a_2\mu_{2j}$, $j = 1, 2, \ldots, m$. Applying Afriat's theorem to the matrix $A$ we have the $\sigma$ given by the eigenvalue problems

$$\left\{ \begin{bmatrix} \dfrac{a_2\mu_{2j}}{a_2 + d_1^{(1)}}, & \dfrac{1_2^{(1)}}{a_2 + d_1^{(1)}} & & \\ \dfrac{u_1^{(1)}}{a_2 + d_2^{(1)}}, & \dfrac{a_2\mu_{2j}}{a_2 + d_2^{(1)}}, & \dfrac{1_3^{(1)}}{a_2 + d_2^{(1)}} & \\ & \cdot \quad \cdot \quad \cdot \quad \cdot \quad \cdot & \\ & \dfrac{u_{n-1}^{(1)}}{a_2 + d_n^{(1)}}, & \dfrac{a_2\mu_{2j}}{a_2 + d_n^{(1)}} \end{bmatrix} + \sigma I \right\} t = 0. \quad (57)$$

Equation (57) is equivalent to the eigenvalue problem

$$[L_1 + \sigma(D_1 + a_2 I) + U_1 + \mu_{2j} a_2 I]t = 0$$

which is precisely the form taken by equation (16) in this case.

Strictly, as $A + \sigma I$ is not diagonally block-tridiagonal, the general theory of S.O.R. would associate a matrix other than $A$ with the eigenvalue problem for the method of simultaneous displacement, the actual form of this matrix depending on the ordering of the components of the solution vector. However, there will be a permutation matrix $P$ such that this matrix is equal to $PAP^T$ (Forsythe and Wasow, 1960, pp. 242–45), and this transformation does not change the eigenvalues of $A$.

### 7. In conclusion

In this paper we have considered only the five-point difference approximation to partial differential equations in two independent variables. For the more accurate nine-point formula to be applicable we must have $f(x)$ and $g(y)$ equal zero in equation (1). We find that the approach used in this paper works only for the case of line over-relaxation, and here we require the further condition that at least one of $p(x)$ and $q(y)$ must vanish identically.

We wish to acknowledge the assistance of Dr. D. W. Martin of the National Physical Laboratory who, besides providing the derivation of equation (16) given in section 6, has made many suggestions concerning presentation which have greatly improved the paper.

### References

AFRIAT, S. N. (1954). "Composite Matrices," *Quart. J. Math.*, Vol. 5, pp. 81–98.
BICKLEY, W. S., and McNAMEE, J. (1960). "Matrix and other Direct Methods for the Solution of Systems of Linear Difference Equations, *Phil. Trans.*, 1005, Vol. 252, pp. 69–131.

FORSYTHE, S. E., and WASOW, W. R. (1960). *Finite Difference Methods for Partial Differential Equations*, Wiley.

FOX, L. (1960). *Some Numerical Experiments with Eigenvalue Problems in Ordinary Differential Equations, Boundary Problems in Differential Equations* (R. E. Langer, editor), University of Wisconsin Press.

MARTIN, D. W., and TEE, G. J. (1961). "Iterative Methods for Linear Equations with Symmetric, Positive Definite Matrix," *The Computer Journal*, Vol. 4, pp. 242–54.

OSBORNE, M. R. (1962). "A Note on Finite-Difference Methods for solving the Eigenvalue Problems of Second-order Differential Equations," *Mathematics of Computation*, Vol. 16, pp. 338–46.

# The $LL^T$ and $QR$ methods for symmetric tridiagonal matrices

*By* James M. Ortega and Henry F. Kaiser

## Introduction

It is known (Rutishauser, 1958, 1959 and 1960) that for a positive definite real symmetric matrix $A_1$ the algorithm:

Decompose $A_i$ into $L_i L_i^T$; $\left.\begin{array}{c} \\ \\ \end{array}\right\}$ $i = 1, 2, \ldots$
Form $A_{i+1} = L_i^T L_i$

where the $L_i$ are lower triangular matrices and $L_i^T$ is the transpose of $L_i$, produces a sequence of matrices $A_i$ which are all similar to $A_1$ and which converge to a diagonal matrix $A$.

A similar algorithm due to Francis (1961):

Decompose $A_i$ into $Q_i R_i$; $\left.\begin{array}{c} \\ \\ \end{array}\right\}$ $i = 1, 2, \ldots$
Form $A_{i+1} = R_i Q_i$

where the $Q_i$ are orthogonal and the $R_i$ are upper triangular, also produces a sequence of matrices which tend to a diagonal matrix. This algorithm has advantage that the similarity transformations involved are orthogonal congruences.

Now when $A_1$ is tridiagonal, that is $a_{ij} = 0$ for $|i - j| > 1$, then the matrices $A_i$ produced by either algorithm are also tridiagonal [Rutishauser (1958) and Francis (1961)] and we might hope that these algorithms would be effective for the important problem of finding the eigenvalues of $A_1$. If we carry out the algorithms in a natural way, however, it would seem that each iteration would require $n$ square roots and consequently would be inefficient compared with the Sturm sequence process now in common use.

It is the purpose of this paper to show that both algorithms can be carried out using no square roots.* Limited experiments have shown that with accelerating techniques involving translations of the origin (see, e.g., Rutishauser, 1959 and 1960), these modified algorithms produce eigenvalues of high accuracy three to ten times as fast as the Sturm process.

## The modified $LL^T$ algorithm

Let $A_1$ be a real symmetric positive definite tridiagonal matrix with diagonal elements $a_1, \ldots, a_n$ and off-diagonal

---

* J. H. Wilkinson, of Teddington, England, has also noted that the $LL^T$ algorithm can be carried out without square roots for tridiagonal matrices [private communication].

---

elements $b_1, \ldots, b_{n-1}$, and let $A_2 A_3 \ldots$ be the matrices produced by the $LL^T$ algorithm. Our modified algorithm produces the diagonal elements and squares of the off-diagonal elements of the $A_i$, and clearly this is sufficient since the signs of the off-diagonal elements do not affect the eigenvalues.

Let $\bar{a}_1, \ldots, \bar{a}_n$ and $\bar{b}_1, \ldots, \bar{b}_{n-1}$ be the diagonal and off-diagonal elements of $A_2$; it will suffice to describe the transition from $a_1, \ldots, a_n$ and $b_1^2, \ldots, b_{n-1}^2$ to $\bar{a}_1, \ldots, \bar{a}_n$ and $\bar{b}_1^2, \ldots, \bar{b}_{n-1}^2$. Now $A_2 = L_1^T L_1$ where $L_1 L_1^T = A_1$. If we let $d_1, \ldots, d_n$ and $s_1, \ldots, s_{n-1}$ be the diagonal and sub-diagonal elements of $L_1$ and carry out the indicated matrix multiplications we then obtain:

$$
\left.
\begin{array}{l}
d_1^2 = a_1; \\
s_i^2 = b_i^2/d_i^2 \\
\bar{a}_i = d_i^2 + s_i^2 \\
d_{i+1}^2 = a_{i+1} - s_i^2 \\
\bar{b}_i^2 = d_{i+1}^2 s_i^2 \\
\bar{a}_n = d_n^2.
\end{array}
\right\} \quad i = 1, \ldots, n-1;
$$

To carry out the algorithm in practice only two temporary storage registers are required to retain the current $s_i^2$ and $d_i^2$; the old $a_i$ and $b_i^2$ are also replaced by $\bar{a}_i^2$ and $\bar{b}_i^2$ immediately, so that a total of only $2n + 1$ storage positions are needed. Each iteration requires no square roots, $n - 1$ divisions, $n - 1$ multiplications and $2(n - 1)$ additions.

Note that we have not included any acceleration procedure in the algorithm, and one must be added to make this a practical process.

## The modified $QR$ algorithm

Again let $a_1, \ldots, a_n$ and $b_1, \ldots, b_{n-1}$ be the diagonal and off-diagonal elements of $A_1$, and $\bar{a}_1 \ldots, \bar{a}_n$ and $\bar{b}_1, \ldots, \bar{b}_{n-1}$ the corresponding elements of $A_2$ where now

$$A_2 = R_1 Q_1$$
and
$$A_1 = Q_1 R_1$$

where $Q_1$ is orthogonal and $R_1$ upper triangular. As with the modified $LL^T$ algorithm it will suffice to describe the transition from $a_1, \ldots, a_n$ and $b_1^2, \ldots, b_{n-1}^2$ to $\bar{a}_1, \ldots, \bar{a}_n$ and $\bar{b}_1^2, \ldots, \bar{b}_{n-1}^2$.