

NOTE

On Robbins Boundary Conditions, Elliptic Equations, and FFT Methods

1. INTRODUCTION

Recently, Pickering and Harley [1, 2] described a method for solving problems in a rectangular region, using cartesian coordinates (x, y) with Robbins boundary conditions, where the partial differential equation is one of the general class of linear elliptic equations which is soluble directly by FFT methods. Such equations must have constant coefficients in the x -direction (assuming that the Fourier transform is applied in this direction) and other coefficients which are, at most, functions of y (see, for example, Le Bail [3], Pickering [4]). In [1, 2] and in the following discussion it is assumed that the Fourier transform is employed in the x -direction. In [1] Poisson's equation was considered, whereas in [2] an illustrative example from the more general class was discussed. In both [1, 2] FFT techniques were used, together with an iterative procedure for solving the associated system of linear equations, and problems were considered in which the Robbins conditions were applied on the boundaries $x = \text{const}$.

The purpose of this note is to extend the application of the methods discussed in [1, 2] to problems where the Robbins conditions are imposed on boundaries $y = \text{const}$. Thus we consider conditions of the form

$$\frac{\partial \phi}{\partial y} = c(x)\phi + d(x), \tag{1}$$

on $y = \text{const}$, where $c(x)$ and $d(x)$ are specified functions. The conditions on the boundaries $x = \text{const}$ may be any suitable conditions which allow the use of the Fourier transform in the x -direction (Dirichlet, Neumann, periodic; see, for example, [3, 4]). For suitable partial differential equations, problems with $c(x) = \text{const}$ lead to systems of linear equations which can be *directly* solved by discrete Fourier analysis, as considered, for example, by Hockney and Eastwood [5] and Beard and Hockney [6].

In [1] an example was considered with $c(x) = \text{const}$ but with a Robbins condition (with variable coefficients) on a boundary $x = \text{constant}$. Hitherto the more general case of condition (1) has not been considered in connection with equations which are soluble by FFT methods. Thus, for elliptic equations of the form indicated, any combination of Dirichlet, Neumann, and Robbins conditions on the boundary sides may be covered by

either direct or iterative application of FFT techniques. The case of (1) with a Robbins condition (with variable coefficients) on $x = \text{const}$ may be treated by a straightforward combination of the techniques described here and in [1]. Furthermore, Pickering and Harley [8, 9] have recently shown how linear elliptic equations with coefficients which depend on both x and y may be solved in rectangular domains with Dirichlet, Neumann, or periodic conditions using a combination of FFT methods and a form of D'Yakonov [7] iteration. It may be possible to extend these methods to the case of Robbins conditions using the techniques described here and in [1].

2. A MODEL PROBLEM

For simplicity, as an illustrative example we consider solving Poisson's equation

$$\frac{\partial^2 \phi}{\partial x^2} + \frac{\partial^2 \phi}{\partial y^2} = q(x, y), \tag{2}$$

in the unit square $0 \leq x, y \leq 1$ with Dirichlet boundary conditions $\phi = f(y)$, $\phi = F(y)$ on $x = 0, 1$, respectively, $\phi = g(x)$ on $y = 0$, and Robbins condition (1) on $y = 1$. It is perhaps worth noting here that the cyclic-reduction algorithms of Swarztrauber and Sweet [10] do not allow the use of a variable c in boundary condition (1), whereas some Chebyshev spectral methods can deal with this case [11, 12]. This latter point is discussed further in Section 4.

Using a square grid of size $h = 1/n$, Eq. (2) was approximated by the usual 5-point formula and the Robbins condition was approximated by the standard central-difference formula. Thus the full set of finite-difference equations for ϕ_{ij} , $i = 1, 2, \dots, n - 1$; $j = 1, 2, \dots, n$, may be written in the form

$$\mathbf{M}\phi_1 + \phi_2 = \mathbf{q}_1^*$$

$$\phi_{j-1} + \mathbf{M}\phi_j + \phi_{j+1} = \mathbf{q}_j^* \quad (j = 2, 3, \dots, n - 1) \tag{3}$$

$$2\phi_{n-1} + (\mathbf{M} + 2hC)\phi_n = \mathbf{q}_n^*,$$

where $\phi_j = (\phi_{1,j}, \phi_{2,j}, \dots, \phi_{n-1,j})^T$, $\mathbf{q}_1^* = \mathbf{q}_1 - g$, $\mathbf{q}_j^* = \mathbf{q}_j$ ($j = 2, 3, \dots, n - 1$), $\mathbf{q}_n^* = \mathbf{q}_n - 2hd$, $\mathbf{q}_j = (h^2q_{1,j} - f_j, h^2q_{2,j}, \dots, h^2q_{n-2,j}, h^2q_{n-1,j} - F_j)^T$, $\mathbf{g} = (g_1, g_2, \dots, g_{n-1})^T$, $\mathbf{d} = (d_1,$

$d_2, \dots, d_{n-1})^T$. The matrix \mathbf{M} is a tridiagonal matrix of order $n - 1$ with diagonal elements -4 and sub- and super-diagonal elements unity and $\mathbf{C} = \text{diag}(c_1, c_2, \dots, c_{n-1})$.

Following the usual procedure we expand both ϕ_j and \mathbf{q}_j^* in terms of the (orthogonal) eigenvectors \mathbf{x}_s of \mathbf{M} so that

$$\phi_j = \sum_{s=1}^{n-1} \alpha_{j,s} \mathbf{x}_s, \quad \mathbf{q}_j^* = \sum_{s=1}^{n-1} \beta_{j,s} \mathbf{x}_s, \quad (4)$$

and hence obtain the system of linear equations for the Fourier harmonic amplitudes $\alpha_{j,r}$ in the form

$$\begin{aligned} \lambda_r \alpha_{1,r} + \alpha_{2,r} &= \beta_{1,r} \\ \alpha_{j-1,r} + \lambda_r \alpha_{j,r} + \alpha_{j+1,r} &= \beta_{j,r} \quad (j = 2, 3, \dots, n-1) \\ 2\alpha_{n-1,r} + \lambda_r \alpha_{n,r} &= \beta_{n,r} - \frac{4}{n^2} \sum_{s=1}^{n-1} \chi_{r,s} \alpha_{n,s} \end{aligned} \quad (5)$$

for $r = 1, 2, \dots, n-1$, where λ_r denotes an eigenvalue of \mathbf{M} and

$$\chi_{r,s} = \sum_{\nu=1}^{n-1} c_\nu \sin \frac{\nu \pi s}{n} \sin \frac{\nu \pi r}{n}. \quad (6)$$

Terms arising from the Robbins boundary condition have been taken to the right-hand side in Eq. (5) because they will be treated as "correction terms" in the computational procedure. Equation (5) differ from the corresponding equations derived in [2] (their Eq. (18)) in that the correction terms in (5) occur in the last equation of each block, whereas in [2] every equation contains correction terms. Furthermore, for the current problem, the coefficient of $\alpha_{n,s}$ in (5) is $\chi_{r,s}$ which may be calculated via (6) in $O(n^2 \log_2 n)$ operations by computing the discrete sine-transform of $c_\nu \sin(\nu \pi r/n)$ for $r = 1, 2, \dots, n-1$. In [1, 2] the right-hand side of (6) was simply $(-1)^{r+s}$.

The block-Jacobi and block-Gauss-Seidel iterations derived from Eq. (5) may be written as

$$A_r \alpha_r^{(m+1)} = \beta_r - \frac{4}{n^2} \sum_{s=1}^{n-1} W_{r,s} \alpha_s^{(m)} \quad (7)$$

and

$$A_r \alpha_r^{(m+1)} = \beta_r - \frac{4}{n^2} \sum_{s=1}^{r-1} W_{r,s} \alpha_s^{(m+1)} - \frac{4}{n^2} \sum_{s=r}^{n-1} W_{r,s} \alpha_s^{(m)}, \quad (8)$$

respectively, for $m = 0, 1, 2, \dots$; $r = 1, 2, \dots, n-1$, where A_r denotes the coefficient matrix of the left-hand side of (5), $\alpha_r = (\alpha_{1,r}, \alpha_{2,r}, \dots, \alpha_{n,r})^T$, $\beta_r = (\beta_{1,r}, \beta_{2,r}, \dots, \beta_{n,r})^T$, and $W_{r,s}$ denotes an $n \times n$ matrix whose elements are all zero apart from the (n, n) element which is $\chi_{r,s}$. Both iterative procedures were used and a theoretical convergence analysis is presented in the following section.

Since A_r is tridiagonal one complete iteration requires $O(n^2)$ operations once the correction term on the right-hand side of (5) has been evaluated. For the Gauss-Seidel case this can be efficiently calculated by a procedure similar to that described in [1] (see their Eq. (23)) which also requires $O(n^2)$ operations per iteration. It thus follows that the overall operation count, $T(m, n)$, has the form

$$T(m, n) = A(m)n^2 + Bn^2 \log_2 n, \quad (9)$$

where A is proportional to the number of iterations and B is a constant. The second term in (9) arises partly from the necessary Fourier analysis and synthesis via (4) and partly from the evaluation of $\chi_{r,s}$. This latter contribution is about a quarter of the second term in (9) if the relevant symmetries are exploited in the evaluation of (6). The FFT procedures used were those of Cooley *et al.* [13] and Press *et al.* [14].

3. ANALYSIS OF CONVERGENCE

For the block-Jacobi method defined by Eq. (7) it is easy to show that the relevant iteration matrix is given by

$$-\frac{4}{n^2} \begin{pmatrix} A_1^{-1} W_{1,1} & A_1^{-1} W_{1,2} & \cdots & A_1^{-1} W_{1,n-1} \\ A_2^{-1} W_{2,1} & A_2^{-1} W_{2,2} & \cdots & A_2^{-1} W_{2,n-1} \\ \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots \\ A_{n-1}^{-1} W_{n-1,1} & A_{n-1}^{-1} W_{n-1,2} & \cdots & A_{n-1}^{-1} W_{n-1,n-1} \end{pmatrix}. \quad (10)$$

Each $n \times n$ block of (10) has zero elements everywhere, except for the last column which consists of the last column of A_r^{-1} multiplied by $\chi_{r,s}$. Furthermore, it can be shown that the blocks of (10) do not commute. Some considerable simplifications in the analysis can be made, however, by assuming that the function $c(x)$ is constant and the work presented here concentrates on that case. The Jacobi iteration matrix (10) is then given by the block-diagonal matrix

$$-\frac{2c}{n} \text{diag} (A_1^{-1} J, A_2^{-1} J, \dots, A_{n-1}^{-1} J), \quad (11)$$

where J is an $n \times n$ matrix whose only non-zero element is the (n, n) element which is unity. Thus, apart from the factor $-2c/n$, each block of (11) is an $n \times n$ matrix whose elements are zero except in the n th column which is equal to the n th column of A_r^{-1} , for $r = 1, 2, \dots, n-1$. Consequently (11) is an upper triangular matrix and, denoting the (n, n) element of A_r^{-1} by $\bar{a}_{n,n}^{(r)}$, it follows that the non-zero eigenvalues, μ_r , of (11) are given by

$$\mu_r = -\frac{2c}{n} \bar{a}_{n,n}^{(r)} \quad (r = 1, 2, \dots, n-1). \quad (12)$$

The elements $\bar{a}_{n,n}^{(r)}$ may be conveniently evaluated by making use of an explicit formula for the inverse of a tridiagonal matrix derived by Vatsya and Pritchard [15]. Using their Eq. (6), after some algebra, we find that

$$\mu_r = \frac{c \tanh[n \cosh^{-1}(1 + 2 \sin^2(\pi r/2n))]}{2n \sin(\pi r/2n)(1 + \sin^2(\pi r/2n))^{1/2}}, \quad (13)$$

where $\mu_r/c > 0$, for $r = 1, 2, \dots, n - 1$. It can be shown that $|\mu_r| > |\mu_{r+1}|$ for $r = 1, 2, \dots, n - 2$. Hence the spectral radius of (11) is $|\mu_1|$ and values of μ_1/c are shown (correct to four decimal places) in Table I for some values of n . As n increases it is clear that these values tend towards a limiting value of approximately 0.3171 (letting $n \rightarrow \infty$ in (13) gives $\mu_1 = \tanh \pi/\pi = 0.3171$ to four decimal places) and that for practical meshes, say $n \geq 32$, the spectral radius is essentially independent of n .

The above analysis indicates that, for constant c , convergence can be expected for $|c| < 3.15$ approximately and that the number of iterations required for convergence is essentially independent of grid size. For the more general case no such simple results have been obtained; the analysis for this case is made particularly complicated by the fact that for such problems the blocks of (10) do not commute. However, the general practical behaviour of the algorithm, described in the next section, exhibits very similar properties to those discussed and analysed above for the special case of constant c .

For the block-Gauss-Seidel iteration a similar analysis was attempted and, in the general case, similar difficulties were encountered as for the Jacobi problem. (For the case of constant c the Gauss-Seidel and Jacobi iteration matrices can be shown to be the same.) In practice, for variable c , the Gauss-Seidel iteration showed a small advantage over the Jacobi iteration and, for different grid sizes, the number of Gauss-Seidel iterations required for convergence was found to be effectively constant.

4. COMPUTATIONAL RESULTS AND DISCUSSION

Table II summarises the results obtained using the Gauss-Seidel iteration with

$$\begin{aligned} f(y) &= 1, & F(y) &= e^y, \\ g(x) &= 1, \\ c(x) &= 2x, & d(x) &= -xe^x, \end{aligned} \quad (14)$$

TABLE I
Values of μ_1/c Derived from Eq. (13)

n	4	8	16	32	64	128
μ_1/c	0.3035	0.3132	0.3161	0.3169	0.3171	0.3171

TABLE II
Numerical Results for the Test Example Using the Gauss-Seidel Iteration

n	RMS error	Error ratios	Iterations	Overall timing (CPU sec)	Timing ratios	% time iterations
4	5.88, -4		15	*		*
8	1.63, -4	2.77	16	*	*	*
16	3.81, -5	4.27	16	1.10, -1	4.00	*
32	8.95, -6	4.26	16	4.40, -1	4.13	88%
64	2.16, -6	4.15	16	1.81, 0	4.15	82%
128	5.37, -7	4.02	16	7.53, 0		81%

* Indicates a value too small to be reliably measured.

so that the analytic solution is

$$\phi = e^{xy}. \quad (15)$$

The computations were performed on a 486DX (33Mhz) PC using Fortran with approximately 12 decimal digit precision and convergence was assumed when successive estimates to all the Fourier harmonic amplitudes agreed to five significant figures.

The RMS error was calculated by comparison of the computed solution with (15) and, as can be seen from Table II, the ratio of RMS error values (for consecutive pairs of n values) tend towards 4, as is expected in view of the h^2 nature of our discretisation. The number of iterations is essentially independent of n and this was also found to be the case in test runs using the Jacobi method, as suggested by the analysis of the previous section. The spectral radii of the relevant iteration matrices were determined numerically for this test example and were found to vary very little with n . For $n \geq 32$ the Jacobi and Gauss-Seidel spectral radii were found to be, respectively, 0.353 and 0.330 (to three decimal places) and the ratio $\ln(0.330)/\ln(0.353) \approx 1.0647$ is in reasonable agreement with the ratio of the number of iterations required for convergence, namely $17/16 = 1.0625$.

Overall execution times are also given in Table II. These values do not include the preprocessing time required for the evaluation of (6) and an entry a in Table II indicates that an estimate time was too small to be reliably determined. For $n = 64$ and $n = 128$ the preprocessing times were found to be between 3% and 4% of the corresponding overall execution times for each of these cases and thus the preprocessing part of the algorithm contributes very little to the overall computational cost. The ratios of execution times for consecutive pairs of n values are approximately 4 which is as expected in view of the form of operations count (9) and the fact that the number of iterations required is effectively constant. Table II also shows

the percentage of the total execution time taken by the iterations. For $n \geq 64$ these figures indicate that the method takes approximately 4.2–4.3 times as long as direct FFT methods which do not require such iterations.

For this particular example Chebyshev spectral solvers based on matrix diagonalisation techniques provide an alternative approach. Haidvogel and Zang [11] considered the Poisson problem with homogeneous Dirichlet conditions and Haldenwang *et al.* [12] cover the case of general boundary conditions for two- and three-dimensional Helmholtz problems. Their results, together with the general theory of spectral methods, demonstrate that, for solutions which are infinitely differentiable, such techniques can achieve very accurate solutions of Poisson's equation considerably more efficiently than methods based on finite-differences (see, for example, [11, Table II]). However, the preprocessing cost associated with these methods is substantial (for example, more than two orders of magnitude greater than the time required for solving Poisson's equation [11, Table VI]). If very high accuracy is not required then the methods described in [1, 2] and the current paper provide an efficient approach to solving the relevant finite-difference problem, without incurring significant preprocessing costs.

Similar considerations are likely to apply to the solution of more general elliptic equations. Pseudospectral iterative methods can provide high accuracy solutions with relatively large preprocessing costs, whereas if lower accuracy can be tolerated the methods of [8, 9] may be used, together with the methods of [1, 2] and the current paper, if the boundary conditions are of Robbins form.

ACKNOWLEDGMENTS

The authors are grateful to a referee whose helpful comments were of assistance in improving an earlier draft of this note.

REFERENCES

1. W. M. Pickering and P. J. Harley, *IMA J. Numer. Anal.* **13**, 215 (1993).
2. W. M. Pickering and P. J. Harley, *Int. J. Comput. Math.* **45**, 243 (1992).
3. R. C. Le Bail, *J. Comput. Phys.* **9**, 440 (1972).
4. W. M. Pickering, *An Introduction to Fast Fourier Transform Methods for Partial Differential Equations, with Applications* (Research Studies Press, Wiley, New York, 1986).
5. R. W. Hockney and J. W. Eastwood, *Computer Simulation using Particles*, (Adam Hilger/IOP, Bristol/New York, 1988), p. 215.
6. S. J. Beard and R. W. Hockney, *Comput. Phys. Commun.* **36**, 25 (1985).
7. E. G. D'Yakonov, *Dokl. Akad. Nauk SSSR* **138**, 522 (1961).
8. W. M. Pickering and P. J. Harley, "Iterative Application of FFT Methods for the Solution of Linear Elliptic Partial Differential Equations," in *Advances in Numerical Methods and Applications, Proceedings of the 3rd International Conference, Sofia, Bulgaria, 1994*, (World Scientific, Singapore, 1994), p. 178.
9. W. M. Pickering and P. J. Harley, *Int. J. Comput. Math.* **55**, 211 (1995).
10. P. Swarztrauber and R. Sweet, *NCAR Technical Note/1A-109*, 1975, Boulder, CO, p. 139 (unpublished).
11. D. B. Haidvogel and T. A. Zang, *J. Comput. Phys.* **30**, 167 (1979).
12. P. Haldenwang, G. Labrosse, S. Abboudi, and M. Deville, *J. Comput. Phys.* **55**, 115 (1984).
13. J. W. Cooley, P. A. W. Lewis, and P. D. Welch, *J. Sound Vib.* **12**, 315 (1970).
14. W. H. Press, B. P. Flannery, S. A. Teukolsky, and W. T. Vetterling, *Numerical Recipes, the art of scientific computing* (Cambridge Univ. Press, Cambridge, 1986), Chap. 12.
15. S. R. Vatsya and H. O. Pritchard, *J. Comput. Math.* **14**, 295 (1983).

Received May 25, 1993; revised November 21, 1994

W. M. PICKERING
P. J. HARLEY

*School of Mathematics and Statistics
Applied Mathematics Section
University of Sheffield
P.O. Box 597, Sheffield S3 7RH
United Kingdom*