

NAG Toolbox

nag_pde_2d_ellip_fd (d03eb)

1 Purpose

`nag_pde_2d_ellip_fd` (d03eb) uses the Strongly Implicit Procedure to calculate the solution to a system of simultaneous algebraic equations of five-point molecule form on a two-dimensional topologically-rectangular mesh. ('Topological' means that a polar grid, for example (r, θ) , can be used, being equivalent to a rectangular box.)

2 Syntax

```
[t, itcoun, itused, resids, chngs, ifail] = nag_pde_2d_ellip_fd(n1, a, b, c, d, e, q, t, aparam, itmax, itcoun, ndir, ixn, iyn, conres, conchn, 'n2', n2)
```

```
[t, itcoun, itused, resids, chngs, ifail] = d03eb(n1, a, b, c, d, e, q, t, aparam, itmax, itcoun, ndir, ixn, iyn, conres, conchn, 'n2', n2)
```

3 Description

Given a set of simultaneous equations

$$Mt = q \quad (1)$$

(which could be nonlinear) derived, for example, from a finite difference representation of a two-dimensional elliptic partial differential equation and its boundary conditions, the function determines the values of the dependent variable t . q is a known vector of length $n_1 \times n_2$ and M is a square $(n_1 \times n_2)$ by $(n_1 \times n_2)$ matrix.

The equations must be of five-diagonal form:

$$a_{ij}t_{i,j-1} + b_{ij}t_{i-1,j} + c_{ij}t_{ij} + d_{ij}t_{i+1,j} + e_{ij}t_{i,j+1} = q_{ij}$$

for $i = 1, 2, \dots, n_1$ and $j = 1, 2, \dots, n_2$, provided $c_{ij} \neq 0.0$. Indeed, if $c_{ij} = 0.0$, then the equation is assumed to be

$$t_{ij} = q_{ij}.$$

For example, if $n_1 = 3$ and $n_2 = 2$, the equations take the form:

$$\begin{bmatrix} c_{11} & d_{11} & & e_{11} & & \\ b_{21} & c_{21} & d_{21} & & e_{21} & \\ & b_{31} & c_{31} & & & e_{31} \\ a_{12} & & & c_{12} & d_{12} & \\ & a_{22} & & b_{22} & c_{22} & d_{22} \\ & & a_{32} & b_{32} & c_{32} & \end{bmatrix} \begin{bmatrix} t_{11} \\ t_{21} \\ t_{31} \\ t_{12} \\ t_{22} \\ t_{32} \end{bmatrix} = \begin{bmatrix} q_{11} \\ q_{21} \\ q_{31} \\ q_{12} \\ q_{22} \\ q_{32} \end{bmatrix}.$$

The system is solved iteratively, from a starting approximation $t^{(1)}$, by the formulae

$$r^{(n)} = q - Mt^{(n)}$$

$$Ms^{(n)} = r^{(n)}$$

$$t^{(n+1)} = t^{(n)} + s^{(n)}.$$

Thus $r^{(n)}$ is the residual of the n th approximate solution $t^{(n)}$, and $s^{(n)}$ is the update change vector. The calling program supplies an initial approximation for the values of the dependent variable in the array \mathbf{t} , the coefficients of the five-point molecule system of equations in the arrays \mathbf{a} , \mathbf{b} , \mathbf{c} , \mathbf{d} and \mathbf{e} , and the source terms in the array \mathbf{q} . The function derives the residual of the latest approximate solution and then uses the approximate LU factorization of the Strongly Implicit Procedure with the necessary

acceleration argument adjustment by calling `nag_pde_2d_ellip_fd_iter` (d03ua) at each iteration. `nag_pde_2d_ellip_fd` (d03eb) combines the newly derived change with the old approximation to obtain the new approximate solution for t . The new solution is checked for convergence against the user-supplied convergence criteria and if these have not been achieved the iterative cycle is repeated. Convergence is based on both the maximum absolute normalized residuals (calculated with reference to the previous approximate solution as these are calculated at the commencement of each iteration) and on the maximum absolute change made to the values of t .

Problems in topologically non-rectangular regions can be solved using the function by surrounding the region by a circumscribing topological rectangle. The equations for the nodal values external to the region of interest are set to zero (i.e., $c_{ij} = t_{ij} = 0$) and the boundary conditions are incorporated into the equations for the appropriate nodes.

If there is no better initial approximation when starting the iterative cycle, an array of all zeros can be used as the initial approximation.

The function can be used to solve linear elliptic equations in which case the arrays **a**, **b**, **c**, **d**, **e** and **q** are constants and for which a single call provides the required solution. It can also be used to solve nonlinear elliptic equations in which case some or all of these arrays may require updating during the progress of the iterations as more accurate solutions are derived. The function will then have to be called repeatedly in an outer iterative cycle. Dependent on the nonlinearity, some under relaxation of the coefficients and/or source terms may be needed during their recalculation using the new estimates of the solution.

The function can also be used to solve each step of a time-dependent parabolic equation in two space dimensions. The solution at each time step can be expressed in terms of an elliptic equation if the Crank–Nicolson or other form of implicit time integration is used.

Neither diagonal dominance, nor positive-definiteness, of the matrix M formed from the arrays **a**, **b**, **c**, **d**, **e** is necessary to ensure convergence.

For problems in which the solution is not unique in the sense that an arbitrary constant can be added to the solution, for example Laplace's equation with all Neumann boundary conditions, a argument is incorporated so that the solution can be rescaled by subtracting a specified nodal value from the whole solution t after the completion of every iteration to keep rounding errors to a minimum for those cases when the convergence is slow.

4 References

Jacobs D A H (1972) The strongly implicit procedure for the numerical solution of parabolic and elliptic partial differential equations *Note RD/L/N66/72* Central Electricity Research Laboratory

Stone H L (1968) Iterative solution of implicit approximations of multi-dimensional partial differential equations *SIAM J. Numer. Anal.* **5** 530–558

5 Parameters

5.1 Compulsory Input Parameters

1: **n1** – INTEGER

The number of nodes in the first coordinate direction, n_1 .

Constraint: **n1** > 1.

2: **a(lda, n2)** – REAL (KIND=nag_wp) array

lda , the first dimension of the array, must satisfy the constraint $lda \geq \mathbf{n1}$.

a(i, j) must contain the coefficient of the ‘southerly’ term involving $t_{i,j-1}$ in the (i, j)th equation of the system (1), for $i = 1, 2, \dots, \mathbf{n1}$ and $j = 1, 2, \dots, \mathbf{n2}$. The elements of **a**, for $j = 1$, must be zero after incorporating the boundary conditions, since they involve nodal values from outside the rectangle.

- 3: **b**(*lda*, **n2**) – REAL (KIND=nag_wp) array

lda, the first dimension of the array, must satisfy the constraint $lda \geq \mathbf{n1}$.

b(*i*, *j*) must contain the coefficient of the ‘westerly’ term involving $t_{i-1,j}$ in the (*i*, *j*)th equation of the system (1), for $i = 1, 2, \dots, \mathbf{n1}$ and $j = 1, 2, \dots, \mathbf{n2}$. The elements of **b**, for $i = 1$, must be zero after incorporating the boundary conditions, since they involve nodal values from outside the rectangle.

- 4: **c**(*lda*, **n2**) – REAL (KIND=nag_wp) array

lda, the first dimension of the array, must satisfy the constraint $lda \geq \mathbf{n1}$.

c(*i*, *j*) must contain the coefficient of the ‘central’ term involving t_{ij} in the (*i*, *j*)th equation of the system (1), for $i = 1, 2, \dots, \mathbf{n1}$ and $j = 1, 2, \dots, \mathbf{n2}$. The elements of **c** are checked to ensure that they are nonzero. If any element is found to be zero, the corresponding algebraic equation is assumed to be $t_{ij} = q_{ij}$. This feature can be used to define the equations for nodes at which, for example, Dirichlet boundary conditions are applied, or for nodes external to the problem of interest, by setting **c**(*i*, *j*) = 0.0 at appropriate points, and the corresponding value of **q**(*i*, *j*) to the appropriate value, namely the prescribed value of **t**(*i*, *j*) in the Dirichlet case or zero at an external point.

- 5: **d**(*lda*, **n2**) – REAL (KIND=nag_wp) array

lda, the first dimension of the array, must satisfy the constraint $lda \geq \mathbf{n1}$.

d(*i*, *j*) must contain the coefficient of the ‘easterly’ term involving $t_{i+1,j}$ in the (*i*, *j*)th equation of the system (1), for $i = 1, 2, \dots, \mathbf{n1}$ and $j = 1, 2, \dots, \mathbf{n2}$. The elements of **d**, for $i = \mathbf{n1}$, must be zero after incorporating the boundary conditions, since they involve nodal values from outside the rectangle.

- 6: **e**(*lda*, **n2**) – REAL (KIND=nag_wp) array

lda, the first dimension of the array, must satisfy the constraint $lda \geq \mathbf{n1}$.

e(*i*, *j*) must contain the coefficient of the ‘northerly’ term involving $t_{i,j+1}$ in the (*i*, *j*)th equation of the system (1), for $i = 1, 2, \dots, \mathbf{n1}$ and $j = 1, 2, \dots, \mathbf{n2}$. The elements of **e**, for $j = \mathbf{n2}$ must be zero after incorporating the boundary conditions, since they involve nodal values from outside the rectangle.

- 7: **q**(*lda*, **n2**) – REAL (KIND=nag_wp) array

lda, the first dimension of the array, must satisfy the constraint $lda \geq \mathbf{n1}$.

q(*i*, *j*) must contain q_{ij} , for $i = 1, 2, \dots, \mathbf{n1}$ and $j = 1, 2, \dots, \mathbf{n2}$, i.e., the source term values at the nodal points for the system (1).

- 8: **t**(*lda*, **n2**) – REAL (KIND=nag_wp) array

lda, the first dimension of the array, must satisfy the constraint $lda \geq \mathbf{n1}$.

t(*i*, *j*) must contain the element t_{ij} of the approximate solution to the equations supplied by the calling program as an initial starting value, for $i = 1, 2, \dots, \mathbf{n1}$ and $j = 1, 2, \dots, \mathbf{n2}$.

If no better approximation is known, an array of zeros can be used.

- 9: **aparam** – REAL (KIND=nag_wp)

The iteration acceleration factor. A value of 1.0 is adequate for most typical problems. However, if convergence is slow, the value can be reduced, typically to 0.2 or 0.1. If divergence is obtained, the value can be increased, typically to 2.0, 5.0 or 10.0.

Constraint: $0.0 < \mathbf{aparam} \leq \left((\mathbf{n1} - 1)^2 + (\mathbf{n2} - 1)^2 \right) / 2.0$.

10: **itmax** – INTEGER

The maximum number of iterations to be used by the function in seeking the solution. A reasonable value might be 30 if $n1 = n2 = 10$ or 100 if $n1 = n2 = 50$.

11: **itcoun** – INTEGER

On the first call of `nag_pde_2d_ellip_fd` (d03eb), **itcoun** must be set to 0. On subsequent entries, its value must be unchanged from the previous call.

12: **ndir** – INTEGER

Indicates whether or not the system of equations has a unique solution. For systems which have a unique solution, **ndir** must be set to any nonzero value. For systems derived from, for example, Laplace's equation with all Neumann boundary conditions, i.e., problems in which an arbitrary constant can be added to the solution, **ndir** should be set to 0 and the values of the next two arguments must be specified. For such problems the function subtracts the value of the function derived at the node (**ixn**, **iy**) from the whole solution after each iteration to reduce the possibility of large rounding errors. You must also ensure that for such problems the appropriate consistency condition on the source terms **q** is satisfied.

13: **ixn** – INTEGER

Is ignored unless **ndir** is equal to zero, in which case it must specify the first index of the nodal point at which the solution is to be set to zero. The node should not correspond to a corner node, or to a node external to the region of interest.

14: **iy** – INTEGER

Is ignored unless **ndir** is equal to zero, in which case it must specify the second index of the nodal point at which the solution is to be set to zero. The node should not correspond to a corner node, or to a node external to the region of interest.

15: **conres** – REAL (KIND=nag_wp)

The convergence criterion to be used on the maximum absolute value of the normalized residual vector components. The latter is defined as the residual of the algebraic equation divided by the central coefficient when the latter is not equal to 0.0, and defined as the residual when the central coefficient is zero.

Clearly **conres** should not be less than a reasonable multiple of the *machine precision*.

16: **conchn** – REAL (KIND=nag_wp)

The convergence criterion to be used on the maximum absolute value of the change made at each iteration to the elements of the array **t**, namely the dependent variable. Clearly **conchn** should not be less than a reasonable multiple of the *machine precision* multiplied by the maximum value of **t** attained.

Convergence is achieved when both the convergence criteria are satisfied. You can therefore set convergence on either the residual or on the change, or (as is recommended) on a requirement that both are below prescribed limits.

5.2 Optional Input Parameters

1: **n2** – INTEGER

Default: the second dimension of the arrays **a**, **b**, **c**, **d**, **e**, **q**, **t**. (An error is raised if these dimensions are not equal.)

The number of nodes in the second coordinate direction, n_2 .

Constraint: $n2 > 1$.

5.3 Output Parameters

1: **t**(*lda*, **n2**) – REAL (KIND=nag_wp) array

The solution derived by the function.

2: **itcoun** – INTEGER

Its value is increased by the number of iterations used on this call (namely **itused**). It therefore stores the accumulated number of iterations actually used. For subsequent calls for the same problem, i.e., with the same **n1** and **n2** but possibly different coefficients and/or source terms, as occur with nonlinear systems or with time-dependent systems, **itcoun** is the accumulated number of iterations. This applies to the second and subsequent calls to nag_pde_2d_ellip_fd (d03eb). In this way a suitable cycling of the sequence of iteration arguments is obtained in the calls to nag_pde_2d_ellip_fd_iter (d03ua).

3: **itused** – INTEGER

The number of iterations actually used on that call.

4: **resids**(**itmax**) – REAL (KIND=nag_wp) array

The maximum absolute value of the residuals calculated at the *i*th iteration, for $i = 1, 2, \dots, \mathbf{itused}$. If you want to know the maximum absolute residual of the solution which is returned you must calculate this in the calling program. The sequence of values **resids** indicates the rate of convergence.

5: **chngs**(**itmax**) – REAL (KIND=nag_wp) array

The maximum absolute value of the changes made to the components of the dependent variable **t** at the *i*th iteration, for $i = 1, 2, \dots, \mathbf{itused}$. The sequence of values **chngs** indicates the rate of convergence.

6: **ifail** – INTEGER

ifail = 0 unless the function detects an error (see Section 5).

6 Error Indicators and Warnings

Errors or warnings detected by the function:

ifail = 1

On entry, **n1** < 2,
or **n2** < 2.

ifail = 2

On entry, *lda* < **n1**.

ifail = 3

On entry, **aparam** ≤ 0.0.

ifail = 4

On entry, **aparam** > $\left((\mathbf{n1} - 1)^2 + (\mathbf{n2} - 1)^2 \right) / 2.0$.

ifail = 5

Convergence was not achieved after **itmax** iterations.

ifail = -99

An unexpected error has been triggered by this routine. Please contact NAG.

ifail = -399

Your licence key may have expired or may not have been installed correctly.

ifail = -999

Dynamic memory allocation failed.

7 Accuracy

The improvement in accuracy for each iteration depends on the size of the system and on the condition of the update matrix characterised by the five-diagonal coefficient arrays. The ultimate accuracy obtainable depends on the above factors and on the *machine precision*. The rate of convergence obtained with the Strongly Implicit Procedure is not always smooth because of the cyclic use of nine acceleration arguments. The convergence may become slow with very large problems, for example when $\mathbf{n1} = \mathbf{n2} = 60$. The final accuracy may be judged approximately from the rate of convergence determined from the sequence of values returned in **chngs** and the magnitude of the maximum absolute value of the change vector on the last iteration stored in **chngs(itused)**.

8 Further Comments

The time taken per iteration is approximately proportional to $\mathbf{n1} \times \mathbf{n2}$.

Convergence may not always be obtained when the problem is very large and/or the coefficients of the equations have widely disparate values. The latter case is often associated with a near ill-conditioned matrix.

9 Example

This example solves Laplace's equation in a rectangle with a non-uniform grid spacing in the x and y coordinate directions and with Dirichlet boundary conditions specifying the function on the perimeter of the rectangle equal to

$$e^{(1.0+x)/y(n_2)} \times \cos(y/y(n_2)).$$

9.1 Program Text

```
function d03eb_example

fprintf('d03eb example results\n\n');

% This example demonstrates the solution of the Laplace Equation with
% Dirichlet boundary conditions using a finite-difference 5-point stencil.

% Initialize variables and arrays.
aparam = 1;
itmax = nag_int(100);
itcoun = nag_int(0);
ndir = nag_int(1);
ixn = nag_int(0);
iyn = nag_int(0);
conres = 1e-06;
conchn = 1e-06;
n1 = nag_int(31);
n2 = 46;
x = 0:0.5:15;
y = 0:1:45;

% Set up difference equation coefficients, source terms and
% initial conditions.
```

```

a = zeros(n1, n2);
e = a; b = a; d = a; q = a; t = a;

for j = 2:n2-1
    a(2:n1-1,j) = 2/((y(j)-y(j-1))*(y(j+1)-y(j-1)));
    e(2:n1-1,j) = 2/((y(j+1)-y(j))*(y(j+1)-y(j-1)));
end
for i = 2:n1-1
    b(i,2:n2-1) = 2/((x(i)-x(i-1))*(x(i+1)-x(i-1)));
    d(i,2:n2-1) = 2/((x(i+1)-x(i))*(x(i+1)-x(i-1)));
end
c = -a - b - d - e;

q(1,1:n2) = exp(6*(x(1) +1.0)/y(n2))*cos(7*y(1:n2)/y(n2));
q(n1,1:n2) = exp(6*(x(n1) +1.0)/y(n2))*cos(7*y(1:n2)/y(n2));
q(1:n1,1) = exp(6*(x(1:n1)+1.0)/y(n2))*cos(7*y(1) /y(n2));
q(1:n1,n2) = exp(6*(x(1:n1)+1.0)/y(n2))*cos(7*y(n2) /y(n2));

[u, itcoun, itused, resids, chngs, ifail] = ...
d03eb( ...
    n1, a, b, c, d, e, q, t, aparam, itmax, ...
    itcoun, ndir, ixn, iyn, conres, conchn);

% Output solution statistics.
fprintf('Number of iterations   %4d\n',itused);
fprintf('Maximum residual       %8.3e\n',resids(itused));

% Plot results.
fig1 = figure;
display_plot(x, y, u);

function display_plot(x, y, tOut)
    % Plot array as a mesh and contours.
    meshc(y, x, tOut);
    % Set the title.
    title({'Solution to Laplace''s Equation', ...
        'using the Strongly Implicit Procedure', ...
        'on a Five-point Molecule Discretisation'});
    % Label the axes.
    xlabel('y');
    ylabel('x');
    zlabel('U(x,y)');
    % Set the view to something nice (determined empirically).
    view(35, 20);

```

9.2 Program Results

d03eb example results

```

Number of iterations   29
Maximum residual       3.750e-08

```

