## E02GCF – NAG Fortran Library Routine Document

**Note.** Before using this routine, please read the Users' Note for your implementation to check the interpretation of bold italicised terms and other implementation-dependent details.

## 1 Purpose

E02GCF calculates an $l_\infty$ solution to an over-determined system of linear equations.

## 2 Specification

```
    SUBROUTINE E02GCF(M, N, MDIM, NDIM, A, B, TOL, RELERR, X, RESMAX,
   1                  IRANK, ITER, IFAIL)
    INTEGER          M, N, MDIM, NDIM, IRANK, ITER, IFAIL
    real             A(NDIM,MDIM), B(M), TOL, RELERR, X(N), RESMAX
```

## 3 Description

Given a matrix $A$ with $m$ rows and $n$ columns $(m \geq n)$ and a vector $b$ with $m$ elements, the routine calculates an $l_\infty$ solution to the over-determined system of equations

$$Ax = b.$$

That is to say, it calculates a vector $x$, with $n$ elements, which minimizes the $l_\infty$-norm of the residuals (the absolutely largest residual)

$$r(x) = \max_{1 \leq i \leq m} |r_i|$$

where the residuals $r_i$ are given by

$$r_i = b_i - \sum_{j=1}^{n} a_{ij} x_j, \quad i = 1, 2, \ldots, m.$$

Here $a_{ij}$ is the element in row $i$ and column $j$ of $A$, $b_i$ is the $i$th element of $b$ and $x_j$ the $j$th element of $x$. The matrix $A$ need not be of full rank. The solution is not unique in this case, and may not be unique even if $A$ is of full rank.

Alternatively, in applications where a complete minimization of the $l_\infty$-norm is not necessary, the user may obtain an approximate solution, usually in shorter time, by giving an appropriate value to the parameter RELERR.

Typically in applications to data fitting, data consisting of $m$ points with co-ordinates $(t_i, y_i)$ is to be approximated in the $l_\infty$-norm by a linear combination of known functions $\phi_j(t)$,

$$\alpha_1 \phi_1(t) + \alpha_2 \phi_2(t) + \ldots + \alpha_n \phi_n(t).$$

This is equivalent to finding an $l_\infty$ solution to the over-determined system of equations

$$\sum_{j=1}^{n} \phi_j(t_i)\alpha_j = y_i, \quad i = 1, 2, \ldots, m.$$

Thus if, for each value of $i$ and $j$ the element $a_{ij}$ of the matrix $A$ above is set equal to the value of $\phi_j(t_i)$ and $b_i$ is set equal to $y_i$, the solution vector $x$ will contain the required values of the $\alpha_j$. Note that the independent variable $t$ above can, instead, be a vector of several independent variables (this includes the case where each $\phi_i$ is a function of a different variable, or set of variables).

The algorithm is a modification of the simplex method of linear programming applied to the dual formation of the $l_\infty$ problem (see Barrodale and Phillips [1] and [2]). The modifications are designed to improve the efficiency and stability of the simplex method for this particular application.

# 4 References

[**1**] Barrodale I and Phillips C (1974) An improved algorithm for discrete Chebyshev linear approximation *Proc. 4th Manitoba Conf. Numerical Mathematics* University of Manitoba, Canada 177–190

[**2**] Barrodale I and Phillips C (1975) Solution of an overdetermined system of linear equations in the Chebyshev norm [F4] (Algorithm 495) *ACM Trans. Math. Software* **1 (3)** 264–270

# 5 Parameters

**1:** M — INTEGER *Input*

*On entry:* the number of equations, $m$ (the number of rows of the matrix $A$).

*Constraint:* M $\geq$ N.

**2:** N — INTEGER *Input*

*On entry:* the number of unknowns, $n$ (the number of columns of the matrix $A$).

*Constraint:* N $\geq 1$.

**3:** MDIM — INTEGER *Input*

*On entry:* the second dimension of the array A as declared in the (sub)program from which E02GCF is called.

*Constraint:* MDIM $\geq$ M + 1.

**4:** NDIM — INTEGER *Input*

*On entry:* the first dimension of the array A as declared in the (sub)program from which E02GCF is called.

*Constraint:* NDIM $\geq$ N + 3.

**5:** A(NDIM,MDIM) — ***real*** array *Input/Output*

*On entry:* A$(j,i)$ must contain $a_{ij}$, element in the $i$th row and $j$th column of the matrix $A$ for, $i = 1, 2, \ldots, m$; $j = 1, 2, \ldots, n$ (that is, the **transpose** of the matrix). The remaining elements need not be set. Preferably, the columns of the matrix $A$ (rows of the parameter A) should be scaled before entry: see Section 7.

*On exit:* A contains the last simplex tableau.

**6:** B(M) — ***real*** array *Input/Output*

*On entry:* $b_i$, the $i$th element of the vector $b$, for $i = 1, 2, \ldots, m$.

*On exit:* the $i$th residual $r_i$ corresponding to the solution vector $x$, for $i = 1, 2, \ldots, m$. Note however that these residuals may contain few significant figures, especially when RESMAX is within one or two orders of magnitude of TOL. Indeed if RESMAX $\leq$ TOL, the elements B$(i)$ may all be set to zero. It is therefore often advisable to compute the residuals directly.

**7:** TOL — ***real*** *Input*

*On entry:* a threshold below which numbers are regarded as zero. The recommended threshold value is $10.0 \times \epsilon$, where $\epsilon$ is the ***machine precision***. If TOL $\leq 0.0$ on entry, the recommended value is used within the routine. If premature termination occurs, a larger value for TOL may result in a valid solution.

*Suggested value:* 0.0.

**8:**    RELERR — **real**                                                                    *Input/Output*

*On entry:* RELERR must be set to a bound on the relative error acceptable in the maximum residual at the solution.

If RELERR $\leq$ 0.0, then the $l_\infty$ solution is computed, and RELERR is set to 0.0 on exit.

If RELERR $>$ 0.0, then the routine obtains instead an approximate solution for which the largest residual is less than $1.0 +$ RELERR times that of the $l_\infty$ solution; on exit, RELERR contains a smaller value such that the above bound still applies. (The usual result of this option, say with RELERR $= 0.1$, is a saving in the number of simplex iterations).

*On exit:* RELERR is altered as described above.

**9:**    X(N) — **real** array                                                                    *Output*

*On exit:* if IFAIL $= 0$ or 1, X($j$) contains the $j$th element of the solution vector $x$, for $j = 1, 2, \ldots, n$. Whether this is an $l_\infty$ solution or an approximation to one, depends on the value of RELERR on entry.

**10:**    RESMAX — **real**                                                                    *Output*

*On exit:* if IFAIL $= 0$ or 1, RESMAX contains the absolute value of the largest residual(s) for the solution vector $x$. (See B above.)

**11:**    IRANK — INTEGER                                                                    *Output*

*On exit:* if IFAIL $= 0$ or 1, IRANK contains the computed rank of the matrix $A$.

**12:**    ITER — INTEGER                                                                    *Output*

*On exit:* if IFAIL $= 0$ or 1, ITER contains the number of iterations taken by the simplex method.

**13:**    IFAIL — INTEGER                                                                    *Input/Output*

*On entry:* IFAIL must be set to 0, $-1$ or 1. Users who are unfamiliar with this parameter should refer to Chapter P01 for details.

*On exit:* IFAIL $= 0$ unless the routine detects an error or gives a warning (see Section 6).

**For this routine**, because the values of output parameters may be useful even if IFAIL $\neq 0$ on exit, users are recommended to set IFAIL to $-1$ before entry. **It is then essential to test the value of IFAIL on exit**. To suppress the output of an error message when soft failure occurs, set IFAIL to 1.

# 6    Error Indicators and Warnings

Errors or warnings specified by the routine:

IFAIL $= 1$

An optimal solution has been obtained but this may not be unique (perhaps simply because the matrix $A$ is not of full rank, i.e., IRANK $<$ N).

IFAIL $= 2$

The calculations have terminated prematurely due to rounding errors. Experiment with larger values of TOL or try rescaling the columns of the matrix (see Section 8).

IFAIL $= 3$

On entry,   NDIM $<$ N $+ 3$

            or   MDIM $<$ M $+ 1$

            or   M $<$ N

            or   N $< 1$.

## 7   Accuracy

Experience suggests that the computational accuracy of the solution $x$ is comparable with the accuracy that could be obtained by applying Gaussian elimination with partial pivoting to the $n+1$ equations which have residuals of largest absolute value. The accuracy therefore varies with the conditioning of the problem, but has been found generally very satisfactory in practice.

## 8   Further Comments

The effects of $m$ and $n$ on the time and on the number of iterations in the simplex method vary from problem to problem, but typically the number of iterations is a small multiple of $n$ and the total time is approximately proportional to $mn^2$.

It is recommended that, before the routine is entered, the columns of the matrix $A$ are scaled so that the largest element in each column is of the order of unity. This should improve the conditioning of the matrix, and also enable the parameter TOL to perform its correct function. The solution $x$ obtained will then, of course, relate to the scaled form of the matrix. Thus if the scaling is such that, for each $j = 1, 2, \ldots, n$, the elements of the $j$th column are multiplied by the constant $k_j$, the element $x_j$ of the solution vector $x$ must be multiplied by $k_j$ if it is desired to recover the solution corresponding to the original matrix $A$.

## 9   Example

Suppose we wish to approximate a set of data by a curve of the form

$$y = Ke^t + Le^{-t} + M$$

where $K$, $L$ and $M$ are unknown. Given values $y_i$ at 5 points $t_i$ we may form the over-determined set of equations for $K$, $L$ and $M$

$$e^{t_i}K + e_i^{-t}L + M = y_i, \;\; i = 1, 2, \ldots, 5.$$

E02GCF is used to solve these in the $l_\infty$ sense.

### 9.1   Program Text

**Note.** The listing of the example program presented below uses bold italicised terms to denote precision-dependent details. Please read the Users' Note for your implementation to check the interpretation of these terms. As explained in the Essential Introduction to this manual, the results produced may not be identical for all implementations.

```
*     E02GCF Example Program Text
*     Mark 14 Revised.  NAG Copyright 1989.
*     .. Parameters ..
      INTEGER          N, MMAX, NDIM, MDIM
      PARAMETER        (N=3,MMAX=5,NDIM=N+3,MDIM=MMAX+1)
      INTEGER          NIN, NOUT
      PARAMETER        (NIN=5,NOUT=6)
*     .. Local Scalars ..
      real             RELERR, RESMAX, T, TOL
      INTEGER          I, IFAIL, IRANK, ITER, M
*     .. Local Arrays ..
      real             A(NDIM,MDIM), B(MMAX), X(N)
*     .. External Subroutines ..
      EXTERNAL         E02GCF
*     .. Intrinsic Functions ..
      INTRINSIC        EXP
*     .. Executable Statements ..
      WRITE (NOUT,*) 'E02GCF Example Program Results'
*     Skip heading in data file
      READ (NIN,*)
      READ (NIN,*) M
```

```
        IF (M.GT.0 .AND. M.LE.MMAX) THEN
           DO 20 I = 1, M
              READ (NIN,*) T, B(I)
              A(1,I) = EXP(T)
              A(2,I) = EXP(-T)
              A(3,I) = 1.0e0
   20      CONTINUE
           TOL = 0.0e0
           RELERR = 0.0e0
           IFAIL = 1
*
           CALL E02GCF(M,N,MDIM,NDIM,A,B,TOL,RELERR,X,RESMAX,IRANK,ITER,
     +                 IFAIL)
*
           WRITE (NOUT,*)
           IF (IFAIL.LE.1) THEN
              WRITE (NOUT,99999) 'RESMAX = ', RESMAX, '  Rank = ', IRANK,
     +           ' Iterations = ', ITER, '  IFAIL =', IFAIL
              WRITE (NOUT,*)
              WRITE (NOUT,*) 'Solution'
              WRITE (NOUT,99998) (X(I),I=1,N)
           ELSE
              WRITE (NOUT,99997) 'E02GCF fails with error', IFAIL
           END IF
        END IF
        STOP
*
99999 FORMAT (1X,A,e10.2,A,I5,A,I5,A,I5)
99998 FORMAT (1X,6F10.4)
99997 FORMAT (1X,A,I2)
        END
```

## 9.2 Program Data

```
E02GCF Example Program Data
  5
   0.0 4.501
   0.2 4.360
   0.4 4.333
   0.6 4.418
   0.8 4.625
```

## 9.3 Program Results

```
E02GCF Example Program Results

RESMAX =   0.10E-02  Rank =     3  Iterations =     4  IFAIL =    0

Solution
     1.0049    2.0149    1.4822
```