NAG Fortran Library Routine Document

G02CHF

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

G02CHF performs a multiple linear regression with no constant on a set of variables whose sums of squares and cross-products about zero and correlation-like coefficients are given.

2 Specification

```
SUBROUTINE G02CHF(N, K1, K, SSPZ, ISSPZ, RZ, IRZ, RESULT, COEFF, ICOEFF,1RZINV, IRZINV, CZ, ICZ, WKZ, IWKZ, IFAIL)INTEGERN, K1, K, ISSPZ, IRZ, ICOEFF, IRZINV, ICZ, IWKZ, IFAILrealSSPZ(ISSPZ,K1), RZ(IRZ,K1), RESULT(13), COEFF(K,3),1RZINV(IRZINV,K), CZ(ICZ,K), WKZ(IWKZ,K)
```

3 Description

The routine fits a curve of the form

$$y = b_1 x_1 + b_2 x_2 + \ldots + b_k x_k$$

to the data points

$$(x_{11}, x_{21}, \dots, x_{k1}, y_1) (x_{12}, x_{22}, \dots, x_{k2}, y_2) \vdots (x_{1n}, x_{2n}, \dots, x_{kn}, y_n)$$

such that

$$y_i = b_1 x_{1i} + b_2 x_{2i} + \ldots + b_k x_{ki} + e_i, \quad i = 1, 2, \ldots, n$$

The routine calculates the regression coefficients, b_1, b_2, \ldots, b_k , (and various other statistical quantities) by minimizing

$$\sum_{i=1}^{n} e_i^2.$$

The actual data values $(x_{1i}, x_{2i}, \ldots, x_{ki}, y_i)$ are not provided as input to the routine. Instead, input to the routine consists of:

- (i) The number of cases, n, on which the regression is based.
- (ii) The total number of variables, dependent and independent, in the regression, (k+1).
- (iii) The number of independent variables in the regression, k.
- (iv) The (k+1) by (k+1) matrix $[S_{ij}]$ of sums of squares and cross-products about zero of all the variables in the regression; the terms involving the dependent variable, y, appear in the (k+1)th row and column.
- (v) The (k+1) by (k+1) matrix $[R_{ij}]$ of correlation-like coefficients for all the variables in the regression; the correlations involving the dependent variable, y, appear in the (k+1)th row and column.

The quantities calculated are:

- (a) The inverse of the k by k partition of the matrix of correlation-like coefficients, $[\hat{R}_{ij}]$, involving only the independent variables. The inverse is obtained using an accurate method which assumes that this sub-matrix is positive-definite (see Section 8).
- (b) The modified matrix, $C = [c_{ij}]$, where

$$c_{ij} = rac{ ilde{R}_{ij} ilde{r}^{ij}}{ ilde{S}_{ij}}, \quad i,j = 1, 2, \dots, k,$$

where \tilde{r}^{ij} is the (i, j)th element of the inverse matrix of $[\tilde{R}_{ij}]$ as described in (a) above. Each element of C is thus the corresponding element of the matrix of correlation-like coefficients multiplied by the corresponding element of the inverse of this matrix, divided by the corresponding element of the matrix of sums of squares and cross-products about zero.

(c) The regression coefficients:

$$b_i = \sum_{j=1}^k c_{ij} \tilde{S}_{j(k+1)}, \quad i = 1, 2, \dots, k,$$

where $\tilde{S}_{j(k+1)}$ is the sum of cross-products about zero for the independent variable x_j and the dependent variable y.

(d) The sum of squares attributable to the regression, *SSR*, the sum of squares of deviations about the regression, *SSD*, and the total sum of squares, *SST*:

 $SST = \tilde{S}_{(k+1)(k+1)}$, the sum of squares about zero for the dependent variable, y;

$$SSR = \sum_{j=1}^{k} b_j \tilde{S}_{j(k+1)}; \quad SSD = SST - SSR.$$

(e) The degrees of freedom attributable to the regression, DFR, the degrees of freedom of deviations about the regression, DFD, and the total degrees of freedom, DFT:

DFR = k; DFD = n - k; DFT = n.

(f) The mean square attributable to the regression, MSR, and the mean square of deviations about the regression, MSD:

$$MSR = SSR/DFR;$$
 $MSD = SSD/DFD.$

(g) The *F*-value for the analysis of variance:

$$F = MSR/MSD.$$

(h) The standard error estimate:

$$s = \sqrt{MSD}.$$

(i) The coefficient of multiple correlation, R, the coefficient of multiple determination, R^2 , and the coefficient of multiple determination corrected for the degrees of freedom, \bar{R}^2 :

$$R = \sqrt{1 - \frac{SSD}{SST}}; \quad R^2 = 1 - \frac{SSD}{SST}; \quad \bar{R}^2 = 1 - \frac{SSD \times DFT}{SST \times DFD};$$

(j) The standard error of the regression coefficients:

$$se(b_i) = \sqrt{MSD \times c_{ii}}, \quad i = 1, 2, \dots, k.$$

(k) The *t*-values for the regression coefficients:

$$t(b_i) = \frac{b_i}{se(b_i)}, \quad i = 1, 2, \dots, k.$$

4 References

Draper N R and Smith H (1985) Applied Regression Analysis (2nd Edition) Wiley

5 Parameters

1: N - INTEGER

On entry: the number of cases, n, used in calculating the sums of squares and cross-products and correlation-like coefficients.

2: K1 – INTEGER

On entry: the total number of variables, independent and dependent (k + 1), in the regression. Constraint: $2 \le K1 \le N$.

3: K – INTEGER

On entry: the number of independent variables k in the regression.

Constraint: K = K1 - 1.

4: SSPZ(ISSPZ,K1) - *real* array

On entry: SSPZ(i, j) must be set to \tilde{S}_{ij} , the sum of cross-products about zero for the *i*th and *j*th variables, for i, j = 1, 2, ..., k + 1; terms involving the dependent variable appear in row k + 1 and column k + 1.

5: ISSPZ – INTEGER

On entry: the first dimension of the array SSPZ as declared in the (sub)program from which G02CHF is called.

Constraint: ISSPZ \geq K1.

6: RZ(IRZ,K1) – *real* array

On entry: RZ(i, j) must be set to \tilde{R}_{ij} , the correlation-like coefficient for the *i*th and *j*th variables, for i, j = 1, 2, ..., k + 1; coefficients involving the dependent variable appear in row k + 1 and column k + 1.

On exit: RZ is used as internal workspace prior to being restored and hence is unchanged.

7: IRZ – INTEGER

On entry: the first dimension of the array RZ as declared in the (sub)program from which G02CHF is called.

Constraint: $IRZ \ge K1$.

8: RESULT(13) – *real* array

On exit: the following information:

RESULT(1) SSR, the sum of squares attributable to the regression;

- RESULT(2) DFR, the degrees of freedom attributable to the regression;
- RESULT(3) MSR, the mean square attributable to the regression;
- RESULT(4) F, the *F*-value for the analysis of variance;
- RESULT(5) SSD, the sum of squares of deviations about the regression;
- RESULT(6) DFD, the degrees of freedom of deviations about the regression;
- RESULT(7) MSD, the mean square of deviations about the regression;
- RESULT(8) SST, the total sum of squares;
- RESULT(9) DFT, the total degrees of freedom;
- RESULT(10) s, the standard error estimate;

G02CHF

Input

Input

Input

Input

Input

Input

Output

Input/Output

9:

RESULT(11)	R, the coefficient of multiple correlation;		
RESULT(12)	R^2 , the coefficient of multiple determination;		
RESULT(13)	$ar{R}^2$, the coefficient of multiple determination corrected for the degrees of	freedom.	
COEFF(K,3) – <i>real</i> array Output		Output	
On exit: for $i = 1, 2,, k$, the following information:			

COEFF(i, 1) b_i , the regression coefficient for the *i*th variable;

COEFF(i, 2) se(b_i), the standard error of the regression coefficient for the *i*th variable;

COEFF(i, 3) $t(b_i)$, the t-value of the regression coefficient for the *i*th variable.

10: ICOEFF – INTEGER

On entry: the first dimension of the array COEFF as declared in the (sub)program from which G02CHF is called.

Constraint: ICOEFF \geq K.

11: RZINV(IRZINV,K) – *real* array

On exit: the inverse of the matrix of correlation-like coefficients for the independent variables; that is, the inverse of the matrix consisting of the first k rows and columns of RZ.

12: IRZINV – INTEGER

On entry: the first dimension of the array RZINV as declared in the (sub)program from which G02CHF is called.

Constraint: IRZINV \geq K.

13: CZ(ICZ,K) - real array

On exit: the modified inverse matrix, C, where

$$CZ(i,j) = \frac{RZ(i,j) \times RZINV(i,j)}{SSPZ(i,j)}, \quad i,j = 1, 2, \dots, k.$$

14: ICZ – INTEGER

On entry: the first dimension of the array CZ as declared in the (sub)program from which G02CHF is called.

Constraint: $ICZ \ge K$.

15: WKZ(IWKZ,K) – *real* array

16: IWKZ – INTEGER

On entry: the first dimension of the array WKZ as declared in the (sub)program from which G02CHF is called.

Constraint: $IWKZ \ge K$.

17: IFAIL – INTEGER

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 (see Section 6).

For environments where it might be inappropriate to halt program execution when an error is detected, the value -1 or 1 is recommended. If the output of error messages is undesirable, then the value 1 is recommended. Otherwise, for users not familiar with this parameter the recommended value is 0. When the value -1 or 1 is used it is essential to test the value of IFAIL on exit.

[NP3546/20A]

Workspace Input

Input

Output

Input

Output

Input

Input/Output

6 Error Indicators and Warnings

If on entry IFAIL = 0 or -1, explanatory error messages are output on the current error message unit (as defined by X04AAF).

Errors or warnings detected by the routine:

IFAIL = 1

On entry, K1 < 2.

IFAIL = 2

On entry, $K1 \neq (K+1)$.

IFAIL = 3

On entry, N < K1.

IFAIL = 4

On entry,	ISSPZ < K1,
or	IRZ < K1,
or	ICOEFF < K,
or	IRZINV < K,
or	ICZ < K,
or	IWKZ < K.

IFAIL = 5

This indicates that the k by k partition of the matrix RZ, which is to be inverted, is not positivedefinite.

IFAIL = 6

This indicates that the refinement following the actual inversion fails, indicating that the k by k partition of the matrix RZ, which is to be inverted, is ill-conditioned. The use of G02DAF, which employs a different numerical technique, may avoid the difficulty.

7 Accuracy

The accuracy of any regression routine is almost entirely dependent on the accuracy of the matrix inversion method used. In this routine, it is the matrix of correlation-like coefficients rather than that of the sums of squares and cross-products about zero that is inverted; this means that all terms in the matrix for inversion are of a similar order, and reduces the scope for computational error. For details on absolute accuracy, the relevant section of the document describing the inversion routine used, F04ABF, should be consulted. G02DAF uses a different method, based on F04AMF, and that routine may well prove more reliable numerically. It does not handle missing values, nor does it provide the same output as this routine.

If, in calculating F or any of the $t(b_i)$ (see Section 3), the numbers involved are such that the result would be outside the range of numbers which can be stored by the machine, then the answer is set to the largest quantity which can be stored as a *real* variable, by means of a call to X02ALF.

8 Further Comments

The time taken by the routine depends on k.

This routine assumes that the matrix of correlation-like coefficients for the independent variables in the regression is positive-definite; it fails if this is not the case.

This correlation matrix will in fact be positive-definite whenever the correlation-like matrix and the sums of squares and cross-products (about zero) matrix have been formed either without regard to missing values, or by eliminating **completely** any cases involving missing values for any variable. If, however,

these matrices are formed by eliminating cases with missing values from only those calculations involving the variables for which the values are missing, no such statement can be made, and the correlation-like matrix may or may not be positive-definite. Users should be aware of the possible dangers of using correlation matrices formed in this way (see the G02 Chapter Introduction), but if they nevertheless wish to carry out regressions using such matrices, this routine is capable of handling the inversion of such matrices, provided they are positive-definite.

If a matrix is positive-definite, its subsequent re-organisation by either of G02CEF or G02CFF will not affect this property and the new matrix can safely be used in this routine. Thus correlation matrices produced by any of G02BDF, G02BEF, G02BKF or G02BLF, even if subsequently modified by either G02CEF or G02CFF, can be handled by this routine.

It should be noted that the routine requires the dependent variable to be the last of the k + 1 variables whose statistics are provided as input to the routine. If this variable is not correctly positioned in the original data, the means, standard deviations, sums of squares and cross-products about zero, and correlation-like coefficients can be manipulated by using G02CEF or G02CFF to re-order the variables as necessary.

9 Example

The following program reads in the sums of squares and cross-products about zero, and correlation-like coefficients for three variables. A multiple linear regression with no constant is then performed with the third and final variable as the dependent variable. Finally the results are printed.

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.

```
GO2CHF Example Program Text
*
     Mark 14 Revised. NAG Copyright 1989.
*
      .. Parameters ..
                       K1, N, K, ISSP, ICORR, ICOEFF, IRINV, IC, IW
      INTEGER
                       (K1=3,N=5,K=K1-1,ISSP=K1,ICORR=K1,ICOEFF=K,
     PARAMETER
     +
                       IRINV=K,IC=K,IW=K)
     INTEGER
                       NIN, NOUT
      PARAMETER
                       (NIN=5,NOUT=6)
      .. Local Scalars ..
      INTEGER
                      I, IFAIL, J
      .. Local Arrays ..
*
                       C(IC,K), COEFFT(ICOEFF,3), CORR(ICORR,K1),
     real
                       RESULT(13), RINV(IRINV,K), SSP(K1,K1), W(IW,K)
     +
      .. External Subroutines ..
*
     EXTERNAL
                       G02CHF
      .. Executable Statements ..
*
      WRITE (NOUT, *) 'GO2CHF Example Program Results'
      Skip heading in data file
      READ (NIN, *)
     READ (NIN,*) ((SSP(I,J),J=1,K1),I=1,K1),
       ((CORR(I,J),J=1,K1),I=1,K1)
     WRITE (NOUT, *)
     WRITE (NOUT, *) 'Sums of squares and cross-products about zero:'
     WRITE (NOUT, 99999) (J, J=1, K1)
      WRITE (NOUT, 99998) (I, (SSP(I,J), J=1, K1), I=1, K1)
      WRITE (NOUT, *)
     WRITE (NOUT, *) 'Correlation-like coefficients:'
      WRITE (NOUT, 99999) (J, J=1, K1)
     WRITE (NOUT, 99998) (I, (CORR(I,J), J=1, K1), I=1, K1)
      WRITE (NOUT, *)
      IFAIL = 1
*
     CALL G02CHF(N,K1,K,SSP,ISSP,CORR,ICORR,RESULT,COEFFT,ICOEFF,RINV,
                  IRINV,C,IC,W,IW,IFAIL)
      IF (IFAIL.NE.O) THEN
         WRITE (NOUT, 99997) 'Routine fails, IFAIL =', IFAIL
```

```
ELSE
          WRITE (NOUT,*) 'Vble Coefft
                                                    Std err
                                                                   t-value'
          WRITE (NOUT,99996) (I,(COEFFT(I,J),J=1,3),I=1,K)
          WRITE (NOUT, *)
          WRITE (NOUT,*) 'Analysis of regression table :-'
          WRITE (NOUT, *)
          WRITE (NOUT, *)
      +′
              Source
                             Sum of squares D.F.
                                                          Mean square F-val
     +ue′
          WRITE (NOUT, *)
          WRITE (NOUT,99995) 'Due to regression', (RESULT(I),I=1,4)
          WRITE (NOUT,99995) 'About regression', (RESULT(I),I=5,7)
WRITE (NOUT,99995) 'Total ', (RESULT(I),I=8,9)
          WRITE (NOUT, *)
          WRITE (NOUT,99994) 'Standard error of estimate =', RESULT(10)
          WRITE (NOUT,99994) 'Multiple correlation (R) =', RESULT(11)
WRITE (NOUT,99994) 'Determination (R squared) =', RESULT(12)
WRITE (NOUT,99994) 'Corrected R squared =', RESULT(13)
          WRITE (NOUT, *)
          WRITE (NOUT, *)
            'Inverse of correlation matrix of independent variables:'
      +
          WRITE (NOUT, 99993) (J, J=1, K)
          WRITE (NOUT,99992) (I,(RINV(I,J),J=1,K),I=1,K)
          WRITE (NOUT, *)
          WRITE (NOUT, *) 'Modified inverse matrix:'
          WRITE (NOUT, 99993) (J, J=1, K)
          WRITE (NOUT, 99992) (I, (C(I,J), J=1,K), I=1,K)
       END IF
       STOP
99999 FORMAT (1X,3I10)
99998 FORMAT (1X,14,3F10.4)
99997 FORMAT (1X,A,I2)
99996 FORMAT (1X,I3,3F13.4)
99995 FORMAT (1X,A,F14.4,F8.0,2F14.4)
99994 FORMAT (1X,A,F8.4)
99993 FORMAT (1X,2I10)
99992 FORMAT (1X,14,2F10.4)
      END
```

9.2 Program Data

G02CHF Example Program Data245.000099.000082.000099.0000271.000052.000082.000052.000054.00001.00000.38420.71290.38421.00000.42990.71290.42991.0000

9.3 **Program Results**

GO2CHF Example Program Results

Sums of squares and cross-products about zero: 1 2 3 1 245.0000 99.0000 82.0000 99.0000271.000052.000082.000052.000054.0000 2 3 Correlation-like coefficients: 1 2 3 1 1.0000 0.3842 0.7129 0.3842 1.0000 0.4299 0.7129 0.4299 1.0000 2 3 Std err t-value Vble Coefft 0.3017 0.1998 0.0817 0.1900 1.5098 0.4299 1 2

Analysis of regression table :-Source Sum of squares D.F. Mean square F-value 28.9857 2. 25.0143 3. 14.4929 8.3381 Due to regression 1.7382 About regression 5. 54.0000 Total Standard error of estimate = 2.8876 Standard error of estimate - 2.0070 Multiple correlation (R) = 0.7326 Determination (R squared) = 0.5368 Corrected R squared = 0.2280 Inverse of correlation matrix of independent variables: 1 2 -0.4507 1.1732 1 2 -0.4507 1.1732 Modified inverse matrix: 1 2 0.0048 -0.0017 -0.0017 0.0043 1 2