

NAG Library Routine Document

G03ADF

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

G03ADF performs canonical correlation analysis upon input data matrices.

2 Specification

```

SUBROUTINE G03ADF (WEIGHT, N, M, Z, LDZ, ISZ, NX, NY, WT, E, LDE, NCV,      &
                  CVX, LDCVX, MCV, CVY, LDCVY, TOL, WK, IWK, IFAIL)
INTEGER           N, M, LDZ, ISZ(M), NX, NY, LDE, NCV, LDCVX, MCV,      &
                  LDCVY, IWK, IFAIL
REAL (KIND=nag_wp) Z(LDZ,M), WT(*), E(LDE,6), CVX(LDCVX,MCV),      &
                  CVY(LDCVY,MCV), TOL, WK(IWK)
CHARACTER(1)     WEIGHT

```

3 Description

Let there be two sets of variables, x and y . For a sample of n observations on n_x variables in a data matrix X and n_y variables in a data matrix Y , canonical correlation analysis seeks to find a small number of linear combinations of each set of variables in order to explain or summarise the relationships between them. The variables thus formed are known as canonical variates.

Let the variance-covariance matrix of the two datasets be

$$\begin{pmatrix} S_{xx} & S_{xy} \\ S_{yx} & S_{yy} \end{pmatrix}$$

and let

$$\Sigma = S_{yy}^{-1} S_{yx} S_{xx}^{-1} S_{xy}$$

then the canonical correlations can be calculated from the eigenvalues of the matrix Σ . However, G03ADF calculates the canonical correlations by means of a singular value decomposition (SVD) of a matrix V . If the rank of the data matrix X is k_x and the rank of the data matrix Y is k_y , and both X and Y have had variable (column) means subtracted then the k_x by k_y matrix V is given by:

$$V = Q_x^T Q_y,$$

where Q_x is the first k_x columns of the orthogonal matrix Q either from the QR decomposition of X if X is of full column rank, i.e., $k_x = n_x$:

$$X = Q_x R_x$$

or from the SVD of X if $k_x < n_x$:

$$X = Q_x D_x P_x^T.$$

Similarly Q_y is the first k_y columns of the orthogonal matrix Q either from the QR decomposition of Y if Y is of full column rank, i.e., $k_y = n_y$:

$$Y = Q_y R_y$$

or from the SVD of Y if $k_y < n_y$:

$$Y = Q_y D_y P_y^T.$$

Let the SVD of V be:

$$V = U_x \Delta U_y^T$$

then the nonzero elements of the diagonal matrix Δ , δ_i , for $i = 1, 2, \dots, l$, are the l canonical correlations associated with the l canonical variates, where $l = \min(k_x, k_y)$.

The eigenvalues, λ_i^2 , of the matrix Σ are given by:

$$\lambda_i^2 = \delta_i^2.$$

The value of $\pi_i = \lambda_i^2 / \sum \lambda_i^2$ gives the proportion of variation explained by the i th canonical variate. The values of the π_i 's give an indication as to how many canonical variates are needed to adequately describe the data, i.e., the dimensionality of the problem.

To test for a significant dimensionality greater than i the χ^2 statistic:

$$\left(n - \frac{1}{2}(k_x + k_y + 3)\right) \sum_{j=i+1}^l \log(1 - \delta_j^2)$$

can be used. This is asymptotically distributed as a χ^2 -distribution with $(k_x - i)(k_y - i)$ degrees of freedom. If the test for $i = k_{\min}$ is not significant, then the remaining tests for $i > k_{\min}$ should be ignored.

The loadings for the canonical variates are calculated from the matrices U_x and U_y respectively. These matrices are scaled so that the canonical variates have unit variance.

4 References

Hastings N A J and Peacock J B (1975) *Statistical Distributions* Butterworth

Kendall M G and Stuart A (1976) *The Advanced Theory of Statistics (Volume 3)* (3rd Edition) Griffin

Morrison D F (1967) *Multivariate Statistical Methods* McGraw-Hill

5 Arguments

- 1: WEIGHT – CHARACTER(1) *Input*
On entry: indicates if weights are to be used.
 WEIGHT = 'U'
 No weights are used.
 WEIGHT = 'W'
 Weights are used and must be supplied in WT.
Constraint: WEIGHT = 'U' or 'W'.
- 2: N – INTEGER *Input*
On entry: n , the number of observations.
Constraint: $N > NX + NY$.
- 3: M – INTEGER *Input*
On entry: m , the total number of variables.
Constraint: $M \geq NX + NY$.
- 4: Z(LDZ, M) – REAL (KIND=nag_wp) array *Input*
On entry: $Z(i, j)$ must contain the i th observation for the j th variable, for $i = 1, 2, \dots, n$ and $j = 1, 2, \dots, m$.

Both x and y variables are to be included in Z , the indicator array, ISZ, being used to assign the variables in Z to the x or y sets as appropriate.

- 5: LDZ – INTEGER *Input*
On entry: the first dimension of the array Z as declared in the (sub)program from which G03ADF is called.
Constraint: $LDZ \geq N$.
- 6: ISZ(M) – INTEGER array *Input*
On entry: ISZ(j) indicates whether or not the j th variable is included in the analysis and to which set of variables it belongs.
 ISZ(j) > 0
 The variable contained in the j th column of Z is included as an x variable in the analysis.
 ISZ(j) < 0
 The variable contained in the j th column of Z is included as a y variable in the analysis.
 ISZ(j) = 0
 The variable contained in the j th column of Z is not included in the analysis.
Constraint: only NX elements of ISZ can be > 0 and only NY elements of ISZ can be < 0.
- 7: NX – INTEGER *Input*
On entry: the number of x variables in the analysis, n_x .
Constraint: $NX \geq 1$.
- 8: NY – INTEGER *Input*
On entry: the number of y variables in the analysis, n_y .
Constraint: $NY \geq 1$.
- 9: WT(*) – REAL (KIND=nag_wp) array *Input*
Note: the dimension of the array WT must be at least N if WEIGHT = 'W', and at least 1 otherwise.
On entry: if WEIGHT = 'W', the first n elements of WT must contain the weights to be used in the analysis.
 If WT(i) = 0.0, the i th observation is not included in the analysis. The effective number of observations is the sum of weights.
 If WEIGHT = 'U', WT is not referenced and the effective number of observations is n .
Constraints:
 $WT(i) \geq 0.0$, for $i = 1, 2, \dots, n$;
 the sum of weights $\geq NX + NY + 1$.
- 10: E(LDE,6) – REAL (KIND=nag_wp) array *Output*
On exit: the statistics of the canonical variate analysis.
 E(i , 1)
 The canonical correlations, δ_i , for $i = 1, 2, \dots, l$.
 E(i , 2)
 The eigenvalues of Σ , λ_i^2 , for $i = 1, 2, \dots, l$.
 E(i , 3)
 The proportion of variation explained by the i th canonical variate, for $i = 1, 2, \dots, l$.

- $E(i, 4)$
The χ^2 statistic for the i th canonical variate, for $i = 1, 2, \dots, l$.
- $E(i, 5)$
The degrees of freedom for χ^2 statistic for the i th canonical variate, for $i = 1, 2, \dots, l$.
- $E(i, 6)$
The significance level for the χ^2 statistic for the i th canonical variate, for $i = 1, 2, \dots, l$.
- 11: LDE – INTEGER *Input*
On entry: the first dimension of the array E as declared in the (sub)program from which G03ADF is called.
Constraint: $LDE \geq \min(NX, NY)$.
- 12: NCV – INTEGER *Output*
On exit: the number of canonical correlations, l . This will be the minimum of the rank of X and the rank of Y.
- 13: CVX(LDCVX, MCV) – REAL (KIND=nag_wp) array *Output*
On exit: the canonical variate loadings for the x variables. $CVX(i, j)$ contains the loading coefficient for the i th x variable on the j th canonical variate.
- 14: LDCVX – INTEGER *Input*
On entry: the first dimension of the array CVX as declared in the (sub)program from which G03ADF is called.
Constraint: $LDCVX \geq NX$.
- 15: MCV – INTEGER *Input*
On entry: an upper limit to the number of canonical variates.
Constraint: $MCV \geq \min(NX, NY)$.
- 16: CVY(LDCVY, MCV) – REAL (KIND=nag_wp) array *Output*
On exit: the canonical variate loadings for the y variables. $CVY(i, j)$ contains the loading coefficient for the i th y variable on the j th canonical variate.
- 17: LDCVY – INTEGER *Input*
On entry: the first dimension of the array CVY as declared in the (sub)program from which G03ADF is called.
Constraint: $LDCVY \geq NY$.
- 18: TOL – REAL (KIND=nag_wp) *Input*
On entry: the value of TOL is used to decide if the variables are of full rank and, if not, what is the rank of the variables. The smaller the value of TOL the stricter the criterion for selecting the singular value decomposition. If a non-negative value of TOL less than *machine precision* is entered, the square root of *machine precision* is used instead.
Constraint: $TOL \geq 0.0$.
- 19: WK(IWK) – REAL (KIND=nag_wp) array *Workspace*
20: IWK – INTEGER *Input*
On entry: the dimension of the array WK as declared in the (sub)program from which G03ADF is called.

Constraints:

if $NX \geq NY$,
 $IWK \geq N \times NX + NX + NY + \max((5 \times (NX - 1) + NX \times NX), N \times NY) + 1$;
 if $NX < NY$,
 $IWK \geq N \times NY + NX + NY + \max((5 \times (NY - 1) + NY \times NY), N \times NX) + 1$.

21: IFAIL – INTEGER

Input/Output

On entry: IFAIL must be set to 0, -1 or 1. If you are unfamiliar with this argument you should refer to Section 3.4 in How to Use the NAG Library and its Documentation for details.

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, if you are not familiar with this argument, the recommended value is 0. **When the value -1 or 1 is used it is essential to test the value of IFAIL on exit.**

On exit: IFAIL = 0 unless the routine detects an error or a warning has been flagged (see Section 6).

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, $NX < 1$,
 or $NY < 1$,
 or $M < NX + NY$,
 or $N \leq NX + NY$,
 or $MCV < \min(NX, NY)$,
 or $LDZ < N$,
 or $LDCVX < NX$,
 or $LDCVY < NY$,
 or $LDE < \min(NX, NY)$,
 or $NX \geq NY$ and
 $IWK < N \times NX + NX + NY + \max((5 \times (NX - 1) + NX \times NX), N \times NY)$,
 or $NX < NY$ and
 $IWK < N \times NY + NX + NY + \max((5 \times (NY - 1) + NY \times NY), N \times NX)$,
 or $WEIGHT \neq 'U'$ or $'W'$,
 or $TOL < 0.0$.

IFAIL = 2

On entry, a WEIGHT = 'W' and value of WT < 0.0.

IFAIL = 3

On entry, the number of x variables to be included in the analysis as indicated by ISZ is not equal to NX.
 or the number of y variables to be included in the analysis as indicated by ISZ is not equal to NY.

IFAIL = 4

On entry, the effective number of observations is less than $NX + NY + 1$.

IFAIL = 5

A singular value decomposition has failed to converge. See F02WUF. This is an unlikely error exit.

IFAIL = 6

A canonical correlation is equal to 1. This will happen if the x and y variables are perfectly correlated.

IFAIL = 7

On entry, the rank of the X matrix or the rank of the Y matrix is 0. This will happen if all the x or y variables are constants.

IFAIL = -99

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

See Section 3.9 in How to Use the NAG Library and its Documentation for further information.

IFAIL = -399

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

See Section 3.8 in How to Use the NAG Library and its Documentation for further information.

IFAIL = -999

Dynamic memory allocation failed.

See Section 3.7 in How to Use the NAG Library and its Documentation for further information.

7 Accuracy

As the computation involves the use of orthogonal matrices and a singular value decomposition rather than the traditional computing of a sum of squares matrix and the use of an eigenvalue decomposition, G03ADF should be less affected by ill-conditioned problems.

8 Parallelism and Performance

G03ADF is threaded by NAG for parallel execution in multithreaded implementations of the NAG Library.

G03ADF makes calls to BLAS and/or LAPACK routines, which may be threaded within the vendor library used by this implementation. Consult the documentation for the vendor library for further information.

Please consult the X06 Chapter Introduction for information on how to control and interrogate the OpenMP environment used within this routine. Please also consult the Users' Note for your implementation for any additional implementation-specific information.

9 Further Comments

None.

10 Example

This example has nine observations and two variables in each set of the four variables read in, the second and third are x variables while the first and last are y variables. Canonical variate analysis is performed and the results printed.

10.1 Program Text

```

Program g03adfe

!      G03ADF Example Program Text

!      Mark 26 Release. NAG Copyright 2016.

!      .. Use Statements ..
Use nag_library, Only: g03adf, nag_wp, x04caf
!      .. Implicit None Statement ..
Implicit None
!      .. Parameters ..
Integer, Parameter          :: nin = 5, nout = 6
!      .. Local Scalars ..
Real (Kind=nag_wp)         :: tol
Integer                    :: i, ifail, iwk, ldcvx, ldcvy, lde,      &
                          ldz, lwt, m, mcv, n, ncv, nx, ny
Character (1)              :: weight
!      .. Local Arrays ..
Real (Kind=nag_wp), Allocatable :: cvx(:,,:), cvy(:,,:), e(:,,:), wk(:),      &
                          wt(:), z(:,:)
Integer, Allocatable       :: isz(:)
!      .. Intrinsic Procedures ..
Intrinsic                  :: max
!      .. Executable Statements ..
Write (nout,*) 'G03ADF Example Program Results'
Write (nout,*)

!      Skip heading in data file
Read (nin,*)

!      Read in problem size
Read (nin,*) n, m, nx, ny, weight

If (weight=='W' .Or. weight=='w') Then
  lwt = n
Else
  lwt = 0
End If
If (nx>=ny) Then
  iwk = n*nx + nx*ny + max(5*(nx-1)+nx*nx,n*ny) + 1
  lde = ny
  mcv = ny
Else
  iwk = n*ny + nx*ny + max(5*(ny-1)+ny*ny,n*nx) + 1
  lde = nx
  mcv = nx
End If
ldz = n
ldcvx = nx
ldcvy = ny
Allocate (z(ldz,m),isz(m),wt(lwt),e(lde,6),cvx(ldcvx,mcv),      &
          cvy(ldcvy,mcv),wk(iwk))

!      Read in data
If (lwt>0) Then
  Read (nin,*)(z(i,1:m),wt(i),i=1,n)
Else
  Read (nin,*)(z(i,1:m),i=1,n)
End If

!      Read in variable inclusion flags
Read (nin,*) isz(1:m)

!      Use default tolerance
tol = 0.0E0_nag_wp

!      Perform canonical correlation analysis
ifail = 0
Call g03adf(weight,n,m,z,ldz,isz,nx,ny,wt,e,lde,ncv,cvx,ldcvx,mcv,cvy,      &

```

```

        ldcvy,tol,wk,iwk,ifail)

!      Display results
      Write (nout,99999) 'Rank of X = ', nx, ' Rank of Y = ', ny
      Write (nout,*)
      Write (nout,*)
      'Canonical      Eigenvalues Percentage      Chisq      DF      Sig'
      Write (nout,*) 'correlations      variation'
      Write (nout,99998)(e(i,1:6),i=1,ncv)
      Write (nout,*)
      Flush (nout)
      ifail = 0
      Call x04caf('General',' ',nx,ncv,cvx,ldcvx,
        'Canonical Coefficients for X',ifail)
      Write (nout,*)
      Flush (nout)
      ifail = 0
      Call x04caf('General',' ',ny,ncv,cvy,ldcvy,
        'Canonical Coefficients for Y',ifail)

99999 Format (1X,A,I0,A,I0)
99998 Format (1X,2F12.4,F11.4,F10.4,F8.1,F8.4)
      End Program g03adfe

```

10.2 Program Data

```

G03ADF Example Program Data
 9 4 2 2 'U'
80.0 58.4 14.0 21.0
75.0 59.2 15.0 27.0
78.0 60.3 15.0 27.0
75.0 57.4 13.0 22.0
79.0 59.5 14.0 26.0
78.0 58.1 14.5 26.0
75.0 58.0 12.5 23.0
64.0 55.5 11.0 22.0
80.0 59.2 12.5 22.0
-1   1   1   -1

```

10.3 Program Results

G03ADF Example Program Results

Rank of X = 2 Rank of Y = 2

Canonical correlations	Eigenvalues	Percentage variation	Chisq	DF	Sig
0.9570	0.9159	0.8746	14.3914	4.0	0.0061
0.3624	0.1313	0.1254	0.7744	1.0	0.3789

Canonical Coefficients for X

	1	2
1	-0.4261	1.0337
2	-0.3444	-1.1136

Canonical Coefficients for Y

	1	2
1	-0.1415	0.1504
2	-0.2384	-0.3424