

# NAG Toolbox for MATLAB

## g03ad

### 1 Purpose

g03ad performs canonical correlation analysis upon input data matrices.

### 2 Syntax

```
[e, ncv, cvx, cvy, ifail] = g03ad(weight, n, z, isz, nx, ny, wt, mcv,
tol, 'm', m)
```

### 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 summarize the relationships between them. The variables thus formed are known as canonical variates.

Let the variance-covariance matrix of the two data sets 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  $\Sigma$ . However, g03ad 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  $k_x$  rows 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  $k_y$  rows 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  $\Delta$ ,  $\delta_i$ , 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,  $\lambda_i^2$ , of the matrix  $\Sigma$  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  $\pi_i$ '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  $\chi^2$  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  $\chi^2$ -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* Butterworths

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 Parameters

### 5.1 Compulsory Input Parameters

1: **weight** – string

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** – int32 scalar

$n$ , the number of observations.

*Constraint:* **n** > **nx** + **ny**.

3: **z(ldz,m)** – double array

**ldz**, the first dimension of the array, must be at least **n**.

**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.

4: **isz(m)** – int32 array

**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  $< 0$ .

5: **nx – int32 scalar**

The number of  $x$  variables in the analysis,  $n_x$ .

*Constraint:* **nx**  $\geq 1$ .

6: **ny – int32 scalar**

The number of  $y$  variables in the analysis,  $n_y$ .

*Constraint:* **ny**  $\geq 1$ .

7: **wt(\*) – double array**

**Note:** the dimension of the array **wt** must be at least **n** if **weight** = 'W', and at least 1 otherwise.

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  $n$ .

*Constraint:* **wt(i)**  $\geq 0.0$ , for  $i = 1, 2, \dots, n$  and the sum of weights  $\geq \mathbf{nx} + \mathbf{ny} + 1$ .

8: **mcv – int32 scalar**

an upper limit to the number of canonical variates.

*Constraint:* **mcv**  $\geq \min(\mathbf{nx}, \mathbf{ny})$ .

9: **tol – double scalar**

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 nonnegative value of **tol** less than *machine precision* is entered, the square root of *machine precision* is used instead.

*Constraint:* **tol**  $\geq 0.0$ .

## 5.2 Optional Input Parameters

1: **m – int32 scalar**

*Default:* The dimension of the arrays **isz**, **z**. (An error is raised if these dimensions are not equal.)

$m$ , the total number of variables.

*Constraint:* **m**  $\geq \mathbf{nx} + \mathbf{ny}$ .

## 5.3 Input Parameters Omitted from the MATLAB Interface

ldz, lde, ldcvx, ldcvy, wk, iwk

## 5.4 Output Parameters

### 1: **e(lde,6) – double array**

The statistics of the canonical variate analysis.

**e(i,1)**

The canonical correlations,  $\delta_i$ , for  $i = 1, 2, \dots, l$ .

**e(i,2)**

The eigenvalues of  $\Sigma$ ,  $\lambda_i^2$ , 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  $\chi^2$  statistic for the  $i$ th canonical variate, for  $i = 1, 2, \dots, l$ .

**e(i,5)**

The degrees of freedom for  $\chi^2$  statistic for the  $i$ th canonical variate, for  $i = 1, 2, \dots, l$ .

**e(i,6)**

The significance level for the  $\chi^2$  statistic for the  $i$ th canonical variate, for  $i = 1, 2, \dots, l$ .

### 2: **ncv – int32 scalar**

The number of canonical correlations,  $l$ . This will be the minimum of the rank of X and the rank of Y.

### 3: **cvx(ldcvx,mcv) – double array**

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.

### 4: **cvy(ldcvy,mcv) – double array**

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.

### 5: **ifail – int32 scalar**

0 unless the function detects an error (see Section 6).

## 6 Error Indicators and Warnings

Errors or warnings detected by the function:

**ifail = 1**

On entry, **nx** < 1,  
 or **ny** < 1,  
 or **m** < **nx** + **ny**,  
 or **n** ≤ **nx** + **ny**,  
 or **mcv** < min(**nx**, **ny**),  
 or **ldz** < **n**,  
 or **ldcvx** < **nx**,  
 or **ldcvy** < **ny**,  
 or **lde** < min(**nx**, **ny**),  
 or **nx** ≥ **ny** and  
**iwk** < **n** × **nx** + **nx** + **ny** + max((5 × (**nx** – 1) + **nx** × **nx**), **n** × **ny**),

or **nx** < **ny** and  
**iwk** < **n** × **ny** + **nx** + **ny** + max((5 × (**ny** − 1) + **ny** × **ny**), **n** × **nx**),  
or **weight** ≠ '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 f02wu. 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.

## 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, g03ad should be less affected by ill-conditioned problems.

## 8 Further Comments

None.

## 9 Example

```
weight = 'U';
n = int32(9);
z = [80, 58.4, 14, 21;
     75, 59.2, 15, 27;
     78, 60.3, 15, 27;
     75, 57.4, 13, 22;
     79, 59.5, 14, 26;
     78, 58.1, 14.5, 26;
     75, 58, 12.5, 23;
     64, 55.5, 11, 22;
     80, 59.2, 12.5, 22];
isz = [int32(-1);
       int32(1);
       int32(1);
       int32(-1)];
```

```
nx = int32(2);  
ny = int32(2);  
wt = [];  
mcv = int32(2);  
tol = 1e-06;  
[e, ncv, cvx, cvy, ifail] = g03ad(weight, n, z, isz, nx, ny, wt, mcv,  
tol)
```

```
e =  
    0.9570    0.9159    0.8746   14.3914    4.0000    0.0061  
    0.3624    0.1313    0.1254    0.7744    1.0000    0.3789  
ncv =  
      2  
cvx =  
   -0.4261    1.0337  
   -0.3444   -1.1136  
cvy =  
   -0.1415    0.1504  
   -0.2384   -0.3424  
ifail =  
      0
```

---