

# NAG Fortran Library Routine Document

## G01DHF

**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

G01DHF computes the ranks, Normal scores, an approximation to the Normal scores or the exponential scores as requested by the user.

### 2 Specification

```
SUBROUTINE G01DHF(SCORES, TIES, N, X, R, IWRK, IFAIL)
INTEGER          N, IWRK(N), IFAIL
real           X(N), R(N)
CHARACTER*1     SCORES, TIES
```

### 3 Description

G01DHF computes one of the following scores for a sample of observations,  $x_1, x_2, \dots, x_n$ .

#### 1. Rank Scores

The ranks are assigned to the data in ascending order, that is the  $i$ th observation has score  $s_i = k$  if it is the  $k$ th smallest observation in the sample.

#### 2. Normal Scores

The Normal scores are the expected values of the Normal order statistics from a sample of size  $n$ . If  $x_i$  is the  $k$ th smallest observation in the sample, then the score for that observation,  $s_i$ , is  $E(Z_k)$  where  $Z_k$  is the  $k$ th order statistic in a sample of size  $n$  from a standard Normal distribution and  $E$  is the expectation operator.

#### 3. Blom, Tukey and van der Waerden scores

These scores are approximations to the Normal scores. The scores are obtained by evaluating the inverse cumulative Normal distribution function,  $\Phi^{-1}(\cdot)$ , at the values of the ranks scaled into the interval (0,1) using different scaling transformations.

The Blom scores use the scaling transformation  $(r_i - 3/8)/(n + 1/4)$  for the rank  $r_i$ , for  $i = 1, 2, \dots, n$ . Thus the Blom score corresponding to the observation  $x_i$  is

$$s_i = \Phi^{-1}\left(\frac{r_i - 3/8}{n + 1/4}\right).$$

The Tukey scores use the scaling transformation  $(r_i - 1/3)/(n + 1/3)$ ; the Tukey score corresponding to the observation  $x_i$  is

$$s_i = \Phi^{-1}\left(\frac{r_i - 1/3}{n + 1/3}\right).$$

The van der Waerden scores use the scaling transformation  $r_i/(n + 1)$ ; the van der Waerden score corresponding to the observation  $x_i$  is

$$s_i = \Phi^{-1}\left(\frac{r_i}{n + 1}\right).$$

The van der Waerden scores may be used to carry out the van der Waerden test for testing for differences between several population distributions, see Conover (1980).

#### 4. Savage scores

The Savage scores are the expected values of the exponential order statistics from a sample of size  $n$ . They may be used in a test discussed by Savage (1956) and Lehmann (1975). If  $x_i$  is the  $k$ th smallest observation in the sample, then the score for that observation is

$$s_i = E(Y_k) = \frac{1}{n} + \frac{1}{n-1} + \cdots + \frac{1}{n-k+1}$$

where  $Y_k$  is the  $k$ th order statistic in a sample of size  $n$  from a standard exponential distribution and  $E$  is the expectation operator.

Ties may be handled in one of five ways. Let  $x_{t(i)}$ , for  $i = 1, 2, \dots, m$  denote  $m$  tied observations, that is  $x_{t(1)} = x_{t(2)} = \cdots = x_{t(m)}$  with  $t(1) < t(2) < \cdots < t(m)$ . If the rank of  $x_{t(1)}$  is  $k$ , then if ties are ignored the rank of  $x_{t(j)}$  will be  $k + j - 1$ . Let the scores ignoring ties be  $s_{t(1)}^*, s_{t(2)}^*, \dots, s_{t(m)}^*$ . Then the scores,  $s_{t(i)}$ , for  $i = 1, 2, \dots, m$  may be calculated as follows.

If averages are used, then  $s_{t(i)} = \sum_{j=1}^m s_{t(j)}^* / m$ .

If the lowest score is used, then  $s_{t(i)} = s_{t(1)}^*$ .

If the highest score is used, then  $s_{t(i)} = s_{t(m)}^*$ .

If ties are to be broken randomly, then  $s_{t(i)} = s_{t(I)}^*$  where  $I \in \{ \text{random permutation of } 1, 2, \dots, m \}$ .

If ties are to be ignored, then  $s_{t(i)} = s_{t(i)}^*$ .

## 4 References

Blom G (1958) *Statistical Estimates and Transformed Beta-variables* Wiley

Conover W J (1980) *Practical Nonparametric Statistics* Wiley

Lehmann E L (1975) *Nonparametrics: Statistical Methods Based on Ranks* Holden-Day

Savage I R (1956) Contributions to the theory of rank order statistics – the two-sample case *Ann. Math. Statist.* **27** 590–615

Tukey J W (1962) The future of data analysis *Ann. Math. Statist.* **33** 1–67

## 5 Parameters

1: SCORES – CHARACTER\*1 array *Input*

*On entry:* indicates which of the following scores are required;

If SCORES = 'R', the ranks,

If SCORES = 'N', the Normal scores, that is the expected value of the Normal order statistics,

If SCORES = 'B', the Blom version of the Normal scores,

If SCORES = 'T', the Tukey version of the Normal scores,

If SCORES = 'V', the van der Waerden version of the Normal scores,

If SCORES = 'S', the Savage scores, that is the expected value of the exponential order statistics.

*Constraint:* SCORES = 'R', 'N', 'B', 'T', 'V', or 'S'.

2: TIES – CHARACTER\*1 array *Input*

*On entry:* indicates which of the following methods is to be used to assign scores to tied observations;

If TIES = 'A', the average of the scores for tied observations is used,

If TIES = 'L', the lowest score in the group of ties is used,

If TIES = 'H', the highest score in the group of ties is used,

If TIES = 'R', the random number generator is used to randomly untie any group of tied observations,

If TIES = 'I', any ties are ignored, that is the scores are assigned to tied observations in the order that they appear in the data.

*Constraint:* TIES = 'A', 'L', 'H', 'R' or 'I'.

3: N – INTEGER *Input*

*On entry:* the number,  $n$ , of observations.

*Constraint:*  $N \geq 1$ .

4: X(N) – *real* array *Input*

*On entry:* the sample of observations,  $x_i$ , for  $i = 1, 2, \dots, n$ .

5: R(N) – *real* array *Output*

*On exit:* contains the scores,  $s_i$ , for  $i = 1, 2, \dots, n$ , as specified by SCORES.

6: IWRK(N) – INTEGER array *Workspace*

7: 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 (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.**

## 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, SCORES is an invalid character,  
or TIES is an invalid character,  
or  $N < 1$ .

## 7 Accuracy

For SCORES='R', the results should be accurate to *machine precision*. For SCORES='S', the results should be accurate to a small multiple of *machine precision*. For SCORES='N', the results should have a relative accuracy of at least  $\max(100 \times \epsilon, 10^{-8})$  where  $\epsilon$  is the *machine precision*. For SCORES='B', 'T' or 'V', the results should have a relative accuracy of at least  $\max(10 \times \epsilon, 10^{-12})$ .

## 8 Further Comments

If more accurate Normal scores are required G01DAF should be used with appropriate settings for the input parameter ETOL.

Note that when ties are resolved randomly the routine G05EHF is used which calls the NAG random number generator G05CAF. If the user does not initialise the generator then the default seed will be used. If the routine is called at different times using the same data and using either the default seed or a fixed seed, by calling G05CBF, then the same permutation will arise and ties will thus be resolved in the same way. If the user wishes ties to be resolved differently then the generator should be initialised to a non-repeatable number using G05CCF.

## 9 Example

This example program computes and prints the Savage scores for a sample of 5 observations. The average of the scores of any tied observations is used.

### 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.

```
*      G01DHF Example Program Text
*      Mark 15 Release. NAG Copyright 1991.
*      .. Parameters ..
      INTEGER          NIN, NOUT
      PARAMETER       (NIN=5,NOUT=6)
      INTEGER          NMAX
      PARAMETER       (NMAX=20)
*      .. Local Scalars ..
      INTEGER          I, IFAIL, N
*      .. Local Arrays ..
      real            R(NMAX), X(NMAX)
      INTEGER          IWRK(NMAX)
*      .. External Subroutines ..
      EXTERNAL        G01DHF
*      .. Executable Statements ..
      WRITE (NOUT,*) 'G01DHF Example Program Results'
      WRITE (NOUT,*)
*      Skip heading in data file
      READ (NIN,*)
      READ (NIN,*) N
      IF (N.LE.NMAX) THEN
         READ (NIN,*) (X(I),I=1,N)
         IFAIL = 0
*
         CALL G01DHF('Savage','Average',N,X,R,IWRK,IFAIL)
*
         WRITE (NOUT,*) 'The Savage Scores : '
         WRITE (NOUT,*)
+        ' (Average scores are used for tied observations)'
         WRITE (NOUT,*)
         WRITE (NOUT,99999) (R(I),I=1,N)
      ELSE
         WRITE (NOUT,*) 'N is larger than NMAX'
      END IF
      STOP
*
99999 FORMAT (1X,F10.4)
      END
```

### 9.2 Program Data

```
G01DHF Example Program Data
5
2 0 2 2 0
```

### 9.3 Program Results

G01DHF Example Program Results

The Savage Scores :  
(Average scores are used for tied observations)

1.4500  
0.3250  
1.4500  
1.4500  
0.3250

---