

## NAG Toolbox

### nag\_stat\_prob\_kolmogorov2 (g01ez)

#### 1 Purpose

nag\_stat\_prob\_kolmogorov2 (g01ez) returns the probability associated with the upper tail of the Kolmogorov–Smirnov two sample distribution.

#### 2 Syntax

```
[result, ifail] = nag_stat_prob_kolmogorov2(n1, n2, d)
[result, ifail] = g01ez(n1, n2, d)
```

#### 3 Description

Let  $F_{n_1}(x)$  and  $G_{n_2}(x)$  denote the empirical cumulative distribution functions for the two samples, where  $n_1$  and  $n_2$  are the sizes of the first and second samples respectively.

The function nag\_stat\_prob\_kolmogorov2 (g01ez) computes the upper tail probability for the Kolmogorov–Smirnov two sample two-sided test statistic  $D_{n_1, n_2}$ , where

$$D_{n_1, n_2} = \sup_x |F_{n_1}(x) - G_{n_2}(x)|.$$

The probability is computed exactly if  $n_1, n_2 \leq 10000$  and  $\max(n_1, n_2) \leq 2500$  using a method given by Kim and Jenrich (1973). For the case where  $\min(n_1, n_2) \leq 10\%$  of the  $\max(n_1, n_2)$  and  $\min(n_1, n_2) \leq 80$  the Smirnov approximation is used. For all other cases the Kolmogorov approximation is used. These two approximations are discussed in Kim and Jenrich (1973).

#### 4 References

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

Feller W (1948) On the Kolmogorov–Smirnov limit theorems for empirical distributions *Ann. Math. Statist.* **19** 179–181

Kendall M G and Stuart A (1973) *The Advanced Theory of Statistics (Volume 2)* (3rd Edition) Griffin

Kim P J and Jenrich R I (1973) Tables of exact sampling distribution of the two sample Kolmogorov–Smirnov criterion  $D_{mn}(m < n)$  *Selected Tables in Mathematical Statistics* **1** 80–129 American Mathematical Society

Siegel S (1956) *Non-parametric Statistics for the Behavioral Sciences* McGraw–Hill

Smirnov N (1948) Table for estimating the goodness of fit of empirical distributions *Ann. Math. Statist.* **19** 279–281

#### 5 Parameters

##### 5.1 Compulsory Input Parameters

1: **n1** – INTEGER

The number of observations in the first sample,  $n_1$ .

*Constraint:* **n1**  $\geq$  1.

2: **n2** – INTEGER

The number of observations in the second sample,  $n_2$ .

*Constraint:*  $\mathbf{n2} \geq 1$ .

3: **d** – REAL (KIND=nag\_wp)

The test statistic  $D_{n_1, n_2}$ , for the two sample Kolmogorov–Smirnov goodness-of-fit test, that is the maximum difference between the empirical cumulative distribution functions (CDFs) of the two samples.

*Constraint:*  $0.0 \leq \mathbf{d} \leq 1.0$ .

## 5.2 Optional Input Parameters

None.

## 5.3 Output Parameters

1: **result**

The result of the function.

2: **ifail** – INTEGER

**ifail** = 0 unless the function detects an error (see Section 5).

## 6 Error Indicators and Warnings

Errors or warnings detected by the function:

**ifail** = 1

On entry,  $\mathbf{n1} < 1$ ,  
or  $\mathbf{n2} < 1$ .

**ifail** = 2

On entry,  $\mathbf{d} < 0.0$ ,  
or  $\mathbf{d} > 1.0$ .

**ifail** = 3

The approximation solution did not converge in 500 iterations. A tail probability of 1.0 is returned by nag\_stat\_prob\_kolmogorov2 (g01ez).

**ifail** = -99

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

**ifail** = -399

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

**ifail** = -999

Dynamic memory allocation failed.

## 7 Accuracy

The large sample distributions used as approximations to the exact distribution should have a relative error of less than 5% for most cases.

## 8 Further Comments

The upper tail probability for the one-sided statistics,  $D_{n_1, n_2}^+$  or  $D_{n_1, n_2}^-$ , can be approximated by halving the two-sided upper tail probability returned by `nag_stat_prob_kolmogorov2` (`g01ez`), that is  $p/2$ . This approximation to the upper tail probability for either  $D_{n_1, n_2}^+$  or  $D_{n_1, n_2}^-$  is good for small probabilities, (e.g.,  $p \leq 0.10$ ) but becomes poor for larger probabilities.

The time taken by the function increases with  $n_1$  and  $n_2$ , until  $n_1 n_2 > 10000$  or  $\max(n_1, n_2) \geq 2500$ . At this point one of the approximations is used and the time decreases significantly. The time then increases again modestly with  $n_1$  and  $n_2$ .

## 9 Example

The following example reads in 10 different sample sizes and values for the test statistic  $D_{n_1, n_2}$ . The upper tail probability is computed and printed for each case.

### 9.1 Program Text

```
function g01ez_example

fprintf('g01ez example results\n\n');

% Upper tail probabilities for 2-sample Kolmogorov--Smirnov distribution.
n1 = [nag_int( 5); 10; 20; 20; 400; 200; 1000; 200; 15; 100];
n2 = [nag_int(10); 10; 10; 15; 200; 20; 20; 50; 200; 100];
d = [0.5; 0.5; 0.5; 0.4833; 0.1412; 0.2861; 0.2113; 0.1796; 0.18; 0.18];

fprintf('      d      n1  n2  two-sided probability\n');
for j = 1:numel(d)

    [p, ifail] = g01ez( ...
                    n1(j), n2(j), d(j));

fprintf('%8.4f%6d%6d%17.4f\n', d(j), n1(j), n2(j), p);
end
```

### 9.2 Program Results

```
g01ez example results
```

d	n1	n2	two-sided probability
0.5000	5	10	0.3506
0.5000	10	10	0.1678
0.5000	20	10	0.0623
0.4833	20	15	0.0261
0.1412	400	200	0.0083
0.2861	200	20	0.0789
0.2113	1000	20	0.2941
0.1796	200	50	0.1392
0.1800	15	200	0.6926
0.1800	100	100	0.0782

---