

## NAG Library Function Document

### nag\_rank\_regsn (g08rac)

## 1 Purpose

nag\_rank\_regsn (g08rac) calculates the parameter estimates, score statistics and their variance-covariance matrices for the linear model using a likelihood based on the ranks of the observations.

## 2 Specification

```
#include <nag.h>
#include <nagg08.h>
void nag_rank_regsn (Nag_OrderType order, Integer ns, const Integer nv[],
                      const double y[], Integer p, const double x[], Integer pdx,
                      Integer idist, Integer nmax, double tol, double prvr[],
                      Integer pdparvar, Integer irank[], double zin[], double eta[],
                      double vapvec[], double parest[], NagError *fail)
```

## 3 Description

Analysis of data can be made by replacing observations by their ranks. The analysis produces inference for regression arguments arising from the following model.

For random variables  $Y_1, Y_2, \dots, Y_n$  we assume that, after an arbitrary monotone increasing differentiable transformation,  $h(\cdot)$ , the model

$$h(Y_i) = x_i^T \beta + \epsilon_i \quad (1)$$

holds, where  $x_i$  is a known vector of explanatory variables and  $\beta$  is a vector of  $p$  unknown regression coefficients. The  $\epsilon_i$  are random variables assumed to be independent and identically distributed with a completely known distribution which can be one of the following: Normal, logistic, extreme value or double-exponential. In Pettitt (1982) an estimate for  $\beta$  is proposed as  $\hat{\beta} = MX^T a$  with estimated variance-covariance matrix  $M$ . The statistics  $a$  and  $M$  depend on the ranks  $r_i$  of the observations  $Y_i$  and the density chosen for  $\epsilon_i$ .

The matrix  $X$  is the  $n$  by  $p$  matrix of explanatory variables. It is assumed that  $X$  is of rank  $p$  and that a column or a linear combination of columns of  $X$  is not equal to the column vector of 1 or a multiple of it. This means that a constant term cannot be included in the model (1). The statistics  $a$  and  $M$  are found as follows. Let  $\epsilon_i$  have pdf  $f(\epsilon)$  and let  $g = -f'/f$ . Let  $W_1, W_2, \dots, W_n$  be order statistics for a random sample of size  $n$  with the density  $f(\cdot)$ . Define  $Z_i = g(W_i)$ , then  $a_i = E(Z_{r_i})$ . To define  $M$  we need  $M^{-1} = X^T(B - A)X$ , where  $B$  is an  $n$  by  $n$  diagonal matrix with  $B_{ii} = E(g'(W_{r_i}))$  and  $A$  is a symmetric matrix with  $A_{ij} = \text{cov}(Z_{r_i}, Z_{r_j})$ . In the case of the Normal distribution, the  $Z_1 < \dots < Z_n$  are standard Normal order statistics and  $E(g'(W_i)) = 1$ , for  $i = 1, 2, \dots, n$ .

The analysis can also deal with ties in the data. Two observations are adjudged to be tied if  $|Y_i - Y_j| < \text{tol}$ , where  $\text{tol}$  is a user-supplied tolerance level.

Various statistics can be found from the analysis:

- (a) The score statistic  $X^T a$ . This statistic is used to test the hypothesis  $H_0 : \beta = 0$ , see (e).
- (b) The estimated variance-covariance matrix  $X^T(B - A)X$  of the score statistic in (a).
- (c) The estimate  $\hat{\beta} = MX^T a$ .
- (d) The estimated variance-covariance matrix  $M = (X^T(B - A)X)^{-1}$  of the estimate  $\hat{\beta}$ .
- (e) The  $\chi^2$  statistic  $Q = \hat{\beta}^T M^{-1} \hat{\beta} = a^T X(X^T(B - A)X)^{-1} X^T a$  used to test  $H_0 : \beta = 0$ . Under  $H_0$ ,  $Q$  has an approximate  $\chi^2$ -distribution with  $p$  degrees of freedom.

- (f) The standard errors  $M_{ii}^{1/2}$  of the estimates given in (c).
- (g) Approximate  $z$ -statistics, i.e.,  $Z_i = \hat{\beta}_i / se(\hat{\beta}_i)$  for testing  $H_0 : \beta_i = 0$ . For  $i = 1, 2, \dots, n$ ,  $Z_i$  has an approximate  $N(0, 1)$  distribution.

In many situations, more than one sample of observations will be available. In this case we assume the model

$$h_k(Y_k) = X_k^T \beta + e_k, \quad k = 1, 2, \dots, \text{ns},$$

where **ns** is the number of samples. In an obvious manner,  $Y_k$  and  $X_k$  are the vector of observations and the design matrix for the  $k$ th sample respectively. Note that the arbitrary transformation  $h_k$  can be assumed different for each sample since observations are ranked within the sample.

The earlier analysis can be extended to give a combined estimate of  $\beta$  as  $\hat{\beta} = Dd$ , where

$$D^{-1} = \sum_{k=1}^{\text{ns}} X_k^T (B_k - A_k) X_k$$

and

$$d = \sum_{k=1}^{\text{ns}} X_k^T a_k,$$

with  $a_k$ ,  $B_k$  and  $A_k$  defined as  $a$ ,  $B$  and  $A$  above but for the  $k$ th sample.

The remaining statistics are calculated as for the one sample case.

## 4 References

Pettitt A N (1982) Inference for the linear model using a likelihood based on ranks *J. Roy. Statist. Soc. Ser: B* **44** 234–243

## 5 Arguments

1: **order** – Nag\_OrderType *Input*

*On entry:* the **order** argument specifies the two-dimensional storage scheme being used, i.e., row-major ordering or column-major ordering. C language defined storage is specified by **order** = Nag\_RowMajor. See Section 3.2.1.3 in the Essential Introduction for a more detailed explanation of the use of this argument.

*Constraint:* **order** = Nag\_RowMajor or Nag\_ColMajor.

2: **ns** – Integer *Input*

*On entry:* the number of samples.

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

3: **nv[ns]** – const Integer *Input*

*On entry:* the number of observations in the  $i$ th sample, for  $i = 1, 2, \dots, \text{ns}$ .

*Constraint:* **nv[i]**  $\geq 1$ , for  $i = 0, 1, \dots, \text{ns} - 1$ .

4: **y[dim]** – const double *Input*

**Note:** the dimension, *dim*, of the array **y** must be at least  $\left( \sum_{i=1}^{\text{ns}} \mathbf{nv}[i-1] \right)$ .

*On entry:* the observations in each sample. Specifically,  $\mathbf{y}[\sum_{k=1}^{i-1} \mathbf{nv}[k-1] + j - 1]$  must contain the  $j$ th observation in the  $i$ th sample.

5: **p** – Integer *Input*

*On entry:* the number of parameters to be fitted.

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

6: **x[dim]** – const double *Input*

**Note:** the dimension,  $dim$ , of the array **x** must be at least

$$\begin{aligned} & \max(1, \mathbf{pdx} \times \mathbf{p}) \text{ when } \mathbf{order} = \text{Nag\_ColMajor}; \\ & \max\left(1, \left(\sum_{i=1}^{\mathbf{ns}} \mathbf{nv}[i-1]\right) \times \mathbf{pdx}\right) \text{ when } \mathbf{order} = \text{Nag\_RowMajor}. \end{aligned}$$

Where  $\mathbf{X}(i, j)$  appears in this document, it refers to the array element

$$\begin{aligned} & \mathbf{x}[(j-1) \times \mathbf{pdx} + i - 1] \text{ when } \mathbf{order} = \text{Nag\_ColMajor}; \\ & \mathbf{x}[(i-1) \times \mathbf{pdx} + j - 1] \text{ when } \mathbf{order} = \text{Nag\_RowMajor}. \end{aligned}$$

*On entry:* the design matrices for each sample. Specifically,  $\mathbf{X}\left(\sum_{k=1}^{i-1} \mathbf{nv}[k-1] + j, l\right)$  must contain the value of the  $l$ th explanatory variable for the  $j$ th observation in the  $i$ th sample.

*Constraint:* **x** must not contain a column with all elements equal.

7: **pdx** – Integer *Input*

*On entry:* the stride separating row or column elements (depending on the value of **order**) in the array **x**.

*Constraints:*

$$\begin{aligned} & \text{if } \mathbf{order} = \text{Nag\_ColMajor}, \mathbf{pdx} \geq \left(\sum_{i=1}^{\mathbf{ns}} \mathbf{nv}[i-1]\right); \\ & \text{if } \mathbf{order} = \text{Nag\_RowMajor}, \mathbf{pdx} \geq \mathbf{p}. \end{aligned}$$

8: **idist** – Integer *Input*

*On entry:* the error distribution to be used in the analysis.

**idist** = 1  
Normal.

**idist** = 2  
Logistic.

**idist** = 3  
Extreme value.

**idist** = 4  
Double-exponential.

*Constraint:*  $1 \leq \mathbf{idist} \leq 4$ .

9: **nmax** – Integer *Input*

*On entry:* the value of the largest sample size.

*Constraint:*  $\mathbf{nmax} = \max_{1 \leq i \leq \mathbf{ns}} (\mathbf{nv}[i-1])$  and  $\mathbf{nmax} > \mathbf{p}$ .

10: **tol** – double *Input*

*On entry:* the tolerance for judging whether two observations are tied. Thus, observations  $Y_i$  and  $Y_j$  are adjudged to be tied if  $|Y_i - Y_j| < \text{tol}$ .

*Constraint:*  $\text{tol} > 0.0$ .

11: **prvr**[*dim*] – double *Output*

**Note:** the dimension, *dim*, of the array **prvr** must be at least

$\max(1, \text{pdparvar} \times \mathbf{p})$  when **order** = Nag\_ColMajor;  
 $\max(1, \mathbf{p} + 1 \times \text{pdparvar})$  when **order** = Nag\_RowMajor.

Where **PRVR**(*i,j*) appears in this document, it refers to the array element

**prvr**[ $((j-1) \times \text{pdparvar} + i - 1)$ ] when **order** = Nag\_ColMajor;  
**prvr**[ $((i-1) \times \text{pdparvar} + j - 1)$ ] when **order** = Nag\_RowMajor.

*On exit:* the variance-covariance matrices of the score statistics and the parameter estimates, the former being stored in the upper triangle and the latter in the lower triangle. Thus for  $1 \leq i \leq j \leq \mathbf{p}$ , **PRVR**(*i,j*) contains an estimate of the covariance between the *i*th and *j*th score statistics. For  $1 \leq j \leq i \leq \mathbf{p} - 1$ , **PRVR**(*i+1,j*) contains an estimate of the covariance between the *i*th and *j*th parameter estimates.

12: **pdparvar** – Integer *Input*

*On entry:* the stride separating row or column elements (depending on the value of **order**) in the array **prvr**.

*Constraints:*

if **order** = Nag\_ColMajor, **pdparvar**  $\geq \mathbf{p} + 1$ ;  
if **order** = Nag\_RowMajor, **pdparvar**  $\geq \mathbf{p}$ .

13: **irank**[*nmax*] – Integer *Output*

*On exit:* for the one sample case, **irank** contains the ranks of the observations.

14: **zin**[*nmax*] – double *Output*

*On exit:* for the one sample case, **zin** contains the expected values of the function  $g(\cdot)$  of the order statistics.

15: **eta**[*nmax*] – double *Output*

*On exit:* for the one sample case, **eta** contains the expected values of the function  $g'(\cdot)$  of the order statistics.

16: **vapvec**[*nmax*  $\times$  (*nmax* + 1)/2] – double *Output*

*On exit:* for the one sample case, **vapvec** contains the upper triangle of the variance-covariance matrix of the function  $g(\cdot)$  of the order statistics stored column-wise.

17: **parest**[4  $\times$   $\mathbf{p} + 1$ ] – double *Output*

*On exit:* the statistics calculated by the function.

The first **p** components of **parest** contain the score statistics.

The next **p** elements contain the parameter estimates.

**parest**[2  $\times$  **p**] contains the value of the  $\chi^2$  statistic.

The next **p** elements of **parest** contain the standard errors of the parameter estimates.

Finally, the remaining **p** elements of **parest** contain the *z*-statistics.

18: **fail** – NagError \*

*Input/Output*

The NAG error argument (see Section 3.6 in the Essential Introduction).

## 6 Error Indicators and Warnings

### NE\_ALLOC\_FAIL

Dynamic memory allocation failed.

### NE\_BAD\_PARAM

On entry, argument  $\langle value \rangle$  had an illegal value.

### NE\_INT

On entry, **idist** is outside the range 1 to 4: **idist** =  $\langle value \rangle$ .

On entry, **ns** =  $\langle value \rangle$ .

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

On entry, **p** =  $\langle value \rangle$ .

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

On entry, **pdpvar** =  $\langle value \rangle$ .

Constraint: **pdpvar** > 0.

On entry, **pdx** =  $\langle value \rangle$ .

Constraint: **pdx** > 0.

### NE\_INT\_2

On entry, **nmax** =  $\langle value \rangle$  and **p** =  $\langle value \rangle$ .

Constraint: **nmax** > **p**.

On entry, **pdpvar** =  $\langle value \rangle$  and **p** =  $\langle value \rangle$ .

Constraint: **pdpvar**  $\geq p$ .

On entry, **pdpvar** =  $\langle value \rangle$  and **p** =  $\langle value \rangle$ .

Constraint: **pdpvar**  $\geq p + 1$ .

On entry, **pdx** =  $\langle value \rangle$  and **p** =  $\langle value \rangle$ .

Constraint: **pdx**  $\geq p$ .

On entry, **pdx** =  $\langle value \rangle$  and sum **nv**[*i* – 1] =  $\langle value \rangle$ .

Constraint: **pdx**  $\geq$  the sum of **nv**[*i* – 1].

### NE\_INT\_ARRAY

On entry, **nv**[*i*] =  $\langle value \rangle$  and **ns** =  $\langle value \rangle$ .

Constraint: **nv**[*i*]  $\geq 1$ , for *i* = 0, 1, …, **ns** – 1.

### NE\_INT\_ARRAY\_ELEM\_CONS

On entry **M** =  $\langle value \rangle$ .

Constraint: **M** elements of array **nv** > 0.

### NE\_INTERNAL\_ERROR

An internal error has occurred in this function. Check the function call and any array sizes. If the call is correct then please contact NAG for assistance.

### NE\_MAT\_ILL\_DEFINED

The matrix  $X^T(B - A)X$  is either singular or non positive definite.

**NE\_OBSERVATIONS**

All the observations were adjudged to be tied.

**NE\_REAL**

On entry, **tol** =  $\langle value \rangle$ .

Constraint: **tol** > 0.0.

**NE\_REAL\_ARRAY\_ELEM\_CONS**

On entry, all elements in column  $\langle value \rangle$  of **x** are equal to  $\langle value \rangle$ .

**NE\_SAMPLE**

The largest sample size is  $\langle value \rangle$  which is not equal to **nmax**, **nmax** =  $\langle value \rangle$ .

**7 Accuracy**

The computations are believed to be stable.

**8 Parallelism and Performance**

`nag_rank_regsn` (g08rac) is threaded by NAG for parallel execution in multithreaded implementations of the NAG Library.

`nag_rank_regsn` (g08rac) 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 Users' Note for your implementation for any additional implementation-specific information.

**9 Further Comments**

The time taken by `nag_rank_regsn` (g08rac) depends on the number of samples, the total number of observations and the number of arguments fitted.

In extreme cases the parameter estimates for certain models can be infinite, although this is unlikely to occur in practice. See Pettitt (1982) for further details.

**10 Example**

A program to fit a regression model to a single sample of 20 observations using two explanatory variables. The error distribution will be taken to be logistic.

**10.1 Program Text**

```
/* nag_rank_regsn (g08rac) Example Program.
*
* Copyright 2001 Numerical Algorithms Group.
*
* Mark 7, 2001.
*/
#include <stdio.h>
#include <nag.h>
#include <nag_stdlib.h>
#include <nagg08.h>

int main(void)
{
    /* Scalars */
    /*
```

```

double      tol;
Integer     exit_status, i, idist, p, j, nmax, ns, nsum;
Integer     pdx, pdparvar;
NagError    fail;
Nag_OrderType order;

/* Arrays */
double      *eta = 0, *parest = 0, *parvar = 0, *vapvec = 0, *x = 0;
double      *y = 0, *zin = 0;
Integer     *irank = 0, *nv = 0;

#ifdef NAG_COLUMN_MAJOR
#define X(I, J)      x[(J-1)*pdx + I - 1]
#define PARVAR(I, J) parvar[(J-1)*pdparvar + I - 1]
    order = Nag_ColMajor;
#else
#define X(I, J)      x[(I-1)*pdx + J - 1]
#define PARVAR(I, J) parvar[(I-1)*pdparvar + J - 1]
    order = Nag_RowMajor;
#endif

INIT_FAIL(fail);

exit_status = 0;
printf("nag_rank_regsn (g08rac) Example Program Results\n");

/* Skip heading in data file */
scanf("%*[^\n] ");

/* Read number of samples, number of parameters to be fitted,
 * error distribution parameter and tolerance criterion for ties.
 */
scanf("%ld%ld%lf%*[^\n] ", &ns, &p, &idist,
      &tol);

/* Allocate memory to nv only */
if (!(nv = NAG_ALLOC(ns, Integer)))
{
    printf("Allocation failure\n");
    exit_status = -1;
    goto END;
}

printf("\n");
printf("Number of samples =%2ld\n", ns);
printf("Number of parameters fitted =%2ld\n", p);
printf("Distribution =%2ld\n", idist);
printf("Tolerance for ties =%8.5f\n", tol);

/* Read the number of observations in each sample. */

for (i = 1; i <= ns; ++i)
    scanf("%ld", &nv[i - 1]);
scanf("%*[^\n] ");

nmax = 0;
nsum = 0;
for (i = 1; i <= ns; ++i)
{
    nsum += nv[i - 1];
    nmax = MAX(nmax, nv[i - 1]);
}
if (nmax > 0 && nmax <= 100 && nsum > 0 && nsum <= 100)
{
    /* Allocate memory */
    if (!(eta = NAG_ALLOC(nmax, double)) ||
        !(parest = NAG_ALLOC(4*p+1, double)) ||
        !(parvar = NAG_ALLOC((p+1)*p, double)) ||
        !(vapvec = NAG_ALLOC(nmax*(nmax+1)/2, double)) ||
        !(x = NAG_ALLOC(nsum * p, double)) ||
        !(y = NAG_ALLOC(nsum, double)) ||
        !(zin = NAG_ALLOC(nsum, double)) ||
        !(irank = NAG_ALLOC(ns, Integer)) ||
        !(fail = NAG_ALLOC(1, NagError))) ||
        !(order = NAG_ALLOC(1, Nag_OrderType))) ||
        !(exit_status = NAG_ALLOC(1, Integer)))
        exit_status = -1;
}

```

```

    !(y = NAG_ALLOC(nsum, double)) ||
    !(zin = NAG_ALLOC(nmax, double)) ||
    !(irank = NAG_ALLOC(nmax, Integer)))
{
    printf("Allocation failure\n");
    exit_status = -1;
    goto END;
}
#endif NAG_COLUMN_MAJOR
    pdx = nsum;
    pdparvar = p+1;
#else
    pdx = p;
    pdparvar = p;
#endif

/* Read in observations and design matrices for each sample.*/
for (i = 1; i <= nsum; ++i)
{
    scanf("%lf", &y[i - 1]);
    for (j = 1; j <= p; ++j)
        scanf("%lf", &x(i, j));
}
scanf("%*[^\n] ");

/* nag_rank_regsn (g08rac).
 * Regression using ranks, uncensored data
 */
nag_rank_regsn(order, ns, nv, y, p, x, pdx, idist, nmax, tol,
                 parvar, pdparvar, irank, zin, eta, vapvec, parest, &fail);
if (fail.code != NE_NOERROR)
{
    printf("Error from nag_rank_regsn (g08rac).\n%s\n",
           fail.message);
    exit_status = 1;
    goto END;
}

printf("\n");
printf("Score statistic\n");
for (i = 1; i <= p; ++i)
    printf("%9.3f%s", parest[i - 1], i%2 == 0 || i == p?"\n":" ");
printf("\n");

printf("Covariance matrix of score statistic\n");
for (j = 1; j <= p; ++j)
{
    for (i = 1; i <= j; ++i)
        printf("%9.3f%s", PARVAR(i, j),
               i%2 == 0 || i == j?"\n":" ");
}
printf("\n");

printf("Parameter estimates\n");
for (i = 1; i <= p; ++i)
    printf("%9.3f%s", parest[p + i - 1],
           i%2 == 0 || i == p?"\n":" ");
printf("\n");

printf("Covariance matrix of parameter estimates\n");
for (i = 1; i <= p; ++i)

{
    printf(" ");
    for (j = 1; j <= i; ++j)
        printf("%9.3f%s", PARVAR(i + 1, j),
               j%2 == 0 || j == i?"\n":" ");
}

```

```

    printf("\n");

    printf("Chi-squared statistic =%9.3f with%2ld d.f.\n",
           parest[p * 2], p);
    printf("\n");
    printf("Standard errors of estimates and\n");
    printf("approximate z-statistics\n");
    for (i = 1; i <= p; ++i)
        printf("%9.3f%14.3f\n", parest[2*p + 1 + i - 1],
               parest[p * 3 + 1 + i - 1]);
    printf("\n");
}
END:
NAG_FREE(eta);
NAG_FREE(parest);
NAG_FREE(parvar);
NAG_FREE(vapvec);
NAG_FREE(x);
NAG_FREE(y);
NAG_FREE(zin);
NAG_FREE(irank);
NAG_FREE(nv);

return exit_status;
}

```

## 10.2 Program Data

```
nag_rank_regsn (g08rac) Example Program Data
1 2 2 0.00001
20
1.0 1.0 23.0
1.0 1.0 32.0
3.0 1.0 37.0
4.0 1.0 41.0
2.0 1.0 41.0
4.0 1.0 48.0
1.0 1.0 48.0
5.0 1.0 55.0
4.0 1.0 55.0
4.0 0.0 56.0
4.0 1.0 57.0
4.0 1.0 57.0
4.0 1.0 57.0
1.0 0.0 58.0
4.0 1.0 59.0
5.0 0.0 59.0
5.0 0.0 60.0
4.0 1.0 61.0
4.0 1.0 62.0
3.0 1.0 62.0
```

## 10.3 Program Results

```
nag_rank_regsn (g08rac) Example Program Results

Number of samples = 1
Number of parameters fitted = 2
Distribution = 2
Tolerance for ties = 0.00001

Score statistic
-1.048      64.333

Covariance matrix of score statistic
  0.673
-4.159      533.670

Parameter estimates
-0.852      0.114
```

Covariance matrix of parameter estimates  
1.560  
0.012      0.002

Chi-squared statistic =    8.221 with 2 d.f.

Standard errors of estimates and  
approximate z-statistics  
1.249      -0.682  
0.044      2.567

---