

# NAG C Library Function Document

## nag\_rank\_regsn\_censored (g08rbc)

### 1 Purpose

nag\_rank\_regsn\_censored (g08rbc) 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 when some of the observations may be right-censored.

### 2 Specification

```
void nag_rank_regsn_censored (Nag_OrderType order, Integer ns, const Integer nv[],
    const double y[], Integer p, const double x[], Integer pdx,
    const Integer icen[], double gamma, Integer nmax, double tol, double parvar[],
    Integer pdpvar, 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 the regression model where the location parameters of the observations,  $\theta_i$ ,  $i = 1, 2, \dots, n$ , are related by  $\theta = X\beta$ . Here  $X$  is an  $n$  by  $p$  matrix of explanatory variables and  $\beta$  is a vector of  $p$  unknown regression parameters. The observations are replaced by their ranks and an approximation, based on Taylor's series expansion, made to the rank marginal likelihood. For details of the approximation see Pettitt (1982).

An observation is said to be right-censored if we can only observe  $Y_j^*$  with  $Y_j^* \leq Y_j$ . We rank censored and uncensored observations as follows. Suppose we can observe  $Y_j$ , for  $j = 1, 2, \dots, n$ , directly but  $Y_j^*$ , for  $j = n + 1, n + 2, \dots, q$ ;  $n \leq q$ , are censored on the right. We define the rank  $r_j$  of  $Y_j$ , for  $j = 1, 2, \dots, n$ , in the usual way;  $r_j$  equals  $i$  if and only if  $Y_j$  is the  $i$ th smallest amongst the  $Y_1, Y_2, \dots, Y_n$ . The right-censored  $Y_j^*$ , for  $j = n + 1, n + 2, \dots, q$ , has rank  $r_j$  if and only if  $Y_j^*$  lies in the interval  $[Y_{(r_j)}, Y_{(r_j+1)}]$ , with  $Y_0 = -\infty$ ,  $Y_{(n+1)} = +\infty$  and  $Y_{(1)} < \dots < Y_{(n)}$  the ordered  $Y_j$ , for  $j = 1, 2, \dots, n$ .

The distribution of the  $Y$  is assumed to be of the following form. Let  $F_L(y) = e^y / (1 + e^y)$ , the logistic distribution function, and consider the distribution function  $F_\gamma(y)$  defined by  $1 - F_\gamma = [1 - F_L(y)]^{1/\gamma}$ . This distribution function can be thought of as either the distribution function of the minimum,  $X_{1,\gamma}$ , of a random sample of size  $\gamma^{-1}$  from the logistic distribution, or as the  $F_\gamma(y - \log \gamma)$  being the distribution function of a random variable having the  $F$ -distribution with 2 and  $2\gamma^{-1}$  degrees of freedom. This family of generalized logistic distribution functions  $[F_\gamma(\cdot); 0 \leq \gamma < \infty]$  naturally links the symmetric logistic distribution ( $\gamma = 1$ ) with the skew extreme value distribution ( $\lim \gamma \rightarrow 0$ ) and with the limiting negative exponential distribution ( $\lim \gamma \rightarrow \infty$ ). For this family explicit results are available for right-censored data. See Pettitt (1983) for details.

Let  $l_R$  denote the logarithm of the rank marginal likelihood of the observations and define the  $q \times 1$  vector  $a$  by  $a = \dot{l}_R(\theta = 0)$ , and let the  $q$  by  $q$  diagonal matrix  $B$  and  $q$  by  $q$  symmetric matrix  $A$  be given by  $B - A = -\ddot{l}_R(\theta = 0)$ . Then various statistics can be found from the analysis.

- The score statistic  $X^T a$ . This statistic is used to test the hypothesis  $H_0 : \beta = 0$  (see (e)).
- The estimated variance-covariance matrix of the score statistic in (a).
- The estimate  $\hat{\beta}_R = M X^T a$ .
- The estimated variance-covariance matrix  $M = (X^T (B - A) X)^{-1}$  of the estimate  $\hat{\beta}_R$ .

- (e) The  $\chi^2$  statistic  $Q = \hat{\beta}_R M^{-1} \hat{\beta}_r = a^T X (X^T (B - A) X)^{-1} X^T a$ , 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.,  $Z_i = \hat{\beta}_{R_i} / se(\hat{\beta}_{R_i})$  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,

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

where  $\mathbf{ns}$  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

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

and

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

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

- Kalbfleisch J D and Prentice R L (1980) *The Statistical Analysis of Failure Time Data* Wiley
- Pettitt A N (1982) Inference for the linear model using a likelihood based on ranks *J. Roy. Statist. Soc. Ser. B* **44** 234–243
- Pettitt A N (1983) Approximate methods using ranks for regression with censored data *Biometrika* **70** 121–132

## 5 Parameters

- 1: **order** – Nag\_OrderType Input  
*On entry:* the **order** parameter 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 2.2.1.4 of the Essential Introduction for a more detailed explanation of the use of this parameter.  
*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  $i = 1, 2, \dots, \mathbf{ns}$ .  
*Constraint:* **nv**[ $i$ ]  $\geq$  1 for  $i = 0, 1, \dots, \mathbf{ns} - 1$ .

- 4: **y**[*dim*] – const double *Input*  
**Note:** the dimension, *dim*, of the array **y** must be at least  $\sum_{i=0}^{ns-1} \mathbf{nv}[i]$ .  
*On entry:* the observations in each sample. Specifically,  $\mathbf{y}\left(\sum_{k=1}^{i-1} \mathbf{nv}[k-1] + j\right)$  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  $\max(1, \mathbf{pdx} \times \mathbf{p})$  when **order** = **Nag\_ColMajor** and at least  $\max(1, \mathbf{pdx} \times \sum_{i=0}^{ns-1} \mathbf{nv}[i])$  when **order** = **Nag\_RowMajor**.  
 If **order** = **Nag\_ColMajor**, the (*i*, *j*)th element of the matrix *X* is stored in  $\mathbf{x}[(j-1) \times \mathbf{pdx} + i - 1]$  and if **order** = **Nag\_RowMajor**, the (*i*, *j*)th element of the matrix *X* is stored in  $\mathbf{x}[(i-1) \times \mathbf{pdx} + j - 1]$ .  
*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 observations 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 matrix row or column elements (depending on the value of **order**) in the array **x**.  
**Constraints:**  
     if **order** = **Nag\_ColMajor**,  $\mathbf{pdx} \geq \sum_{i=0}^{ns-1} \mathbf{nv}[i]$ ;  
     if **order** = **Nag\_RowMajor**,  $\mathbf{pdx} \geq \mathbf{p}$ .
- 8: **icen**[*dim*] – const Integer *Input*  
**Note:** the dimension, *dim*, of the array **icen** must be at least  $\sum_{i=0}^{ns-1} \mathbf{nv}[i]$ .  
*On entry:* defines the censoring variable for the observations in **y** as follows:  
**icen**[*i* - 1] = 0  
     If **y**[*i* - 1] is uncensored.  
**icen**[*i* - 1] = 1  
     If **y**[*i* - 1] is censored.  
**Constraint:** **icen**[*i* - 1] = 0 or 1, for  $i = 1, 2, \dots, \sum_{i=0}^{ns-1} \mathbf{nv}[i]$ .
- 9: **gamma** – double *Input*  
*On entry:* the value of the parameter defining the generalized logistic distribution. For **gamma** ≤ 0.0001, the limiting extreme value distribution is assumed.  
**Constraint:** **gamma** > 0.0.
- 10: **nmax** – Integer *Input*  
*On entry:* the value of the largest sample size.  
**Constraint:** **nmax** =  $\max_{1 \leq i \leq ns} (\mathbf{nv}[i - 1])$  and **nmax** > **p**.

- 11: **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  $|Y_i - Y_j| < \mathbf{tol}$ .  
*Constraint:* **tol** > 0.0.
- 12: **parvar**[*dim*] – double *Output*  
**Note:** the dimension, *dim*, of the array **parvar** must be at least  $\max(1, \mathbf{pdparvar} \times \mathbf{p})$  when **order** = **Nag\_ColMajor** and at least  $\max(1, \mathbf{pdparvar} \times \mathbf{p} + 1)$  when **order** = **Nag\_RowMajor**.  
Where **PARVAR**(*i*, *j*) appears in this document, it refers to the array element  
if **order** = **Nag\_ColMajor**, **parvar**[(*j* - 1) × **pdparvar** + *i* - 1];  
if **order** = **Nag\_RowMajor**, **parvar**[(*i* - 1) × **pdparvar** + *j* - 1].  
*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  $1 \leq i \leq j \leq \mathbf{p}$ , **PARVAR**(*i*, *j*) contains an estimate of the covariance between the *i*th and *j*th score statistics. For  $1 \leq j \leq i \leq \mathbf{p} - 1$ , **PARVAR**(*i* + 1, *j*) contains an estimate of the covariance between the *i*th and *j*th parameter estimates.
- 13: **pdparvar** – Integer *Input*  
*On entry:* the stride separating matrix row or column elements (depending on the value of **order**) in the array **parvar**.  
*Constraints:*  
if **order** = **Nag\_ColMajor**, **pdparvar** ≥ **p** + 1;  
if **order** = **Nag\_RowMajor**, **pdparvar** ≥ **p**.
- 14: **irank**[**nmax**] – Integer *Output*  
*On exit:* for the one sample case, **irank** contains the ranks of the observations.
- 15: **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.
- 16: **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.
- 17: **vapvec**[*dim*] – double *Output*  
**Note:** the dimension, *dim*, of the array **vapvec** must be at least  $\mathbf{nmax} \times (\mathbf{nmax} + 1)/2$ .  
*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.
- 18: **parest**[*dim*] – double *Output*  
**Note:** the dimension, *dim*, of the array **parest** must be at least  $4 \times \mathbf{p} + 1$ .  
*On exit:* the statistics calculated by the routine as follows. The first **p** components of **parest** contain the score statistics. The next **p** elements contain the parameter estimates. **parest**[2 × **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.
- 19: **fail** – NagError \* *Input/Output*  
The NAG error parameter (see the Essential Introduction).

## 6 Error Indicators and Warnings

### NE\_INT

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

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

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

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

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

Constraint: **pdx**  $> 0$ .

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

Constraint: **pdparvar**  $> 0$ .

### NE\_INT\_2

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

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

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

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

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

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

On entry, **nmax**  $\leq \mathbf{p}$ : **nmax** =  $\langle value \rangle$ , **p** =  $\langle value \rangle$ .

On entry, **pdx**  $<$  the sum of **nv**[*i*]: **pdx** =  $\langle value \rangle$ , sum **nv**[*i*] =  $\langle value \rangle$ .

### NE\_INT\_ARRAY

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

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

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

*M* elements of array **icen** are not equal to 0 or 1:  $M = \langle value \rangle$ .

*M* elements of array **nv** are less than or equal to zero:  $M = \langle value \rangle$ .

### 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**  $> 0.0$ .

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

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

### 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** =  $\langle value \rangle$ .

**NE\_ALLOC\_FAIL**

Memory allocation failed.

**NE\_BAD\_PARAM**

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

**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 consult NAG for assistance.

**7 Accuracy**

The computations are believed to be stable.

**8 Further Comments**

The time taken by the routine depends on the number of samples, the total number of observations and the number of parameters 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.

**9 Example**

A program to fit a regression model to a single sample of 40 observations using just one explanatory variable.

**9.1 Program Text**

```

/* nag_rank_regsn_censored (g08rbc) 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 gamma, tol;
    Integer exit_status, i, 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, *y=0, *zin=0;
    Integer *icen=0, *irank=0, *iwa=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

```

```

#endif

INIT_FAIL(fail);
exit_status = 0;
Vprintf("g08rbc Example Program Results\n");

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

/* Read number of samples, number of parameters to be fitted, */
/* distribution power parameter and tolerance criterion for ties. */

Vscanf("%ld%ld%lf%lf%*[\n] ", &ns, &p, &gamma, &tol);
Vprintf("\n");

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

Vprintf("Number of samples =%2ld\n", ns);
Vprintf("Number of parameters fitted =%2ld\n", p);
Vprintf("Distribution power parameter =%10.5f\n", gamma);

Vprintf("Tolerance for ties =%10.5f\n", tol);

Vprintf("\n");
/* Read the number of observations in each sample */

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

nmax = 0;
nsum = 0;
for (i = 1; i <= ns; ++i)
{
    nsum += nv[i - 1];
    nmax = MAX(nmax, nv[i - 1]);
}

/* Allocate memory */
if ( !(eta = NAG_ALLOC(nmax, double)) ||
    !(parest = NAG_ALLOC(4*p+1, double)) ||
    !(parvar = NAG_ALLOC(7 * 6, double)) ||
    !(vapvec = NAG_ALLOC(nmax*(nmax+1)/2, double)) ||
    !(x = NAG_ALLOC(nsum * p, double)) ||
    !(y = NAG_ALLOC(nsum, double)) ||
    !(zin = NAG_ALLOC(nmax, double)) ||
    !(icen = NAG_ALLOC(nsum, Integer)) ||
    !(irank = NAG_ALLOC(nmax, Integer)) ||
    !(iwa = NAG_ALLOC(400, Integer)) )
{
    Vprintf("Allocation failure\n");
    exit_status = -1;
    goto END;
}
#endif
#ifdef NAG_COLUMN_MAJOR
    pdx = nsum;
    pdparvar = p+1 ;
#else
    pdx = p;
    pdparvar = p;
#endif

/* Read in observations, design matrix and censoring variable */

for (i = 1; i <= nsum; ++i)

```

```

    {
        Vscanf("%lf", &y[i - 1]);

        for (j = 1; j <= p; ++j)
            {
                Vscanf("%lf", &X(i,j));
            }
        Vscanf("%ld", &icen[i - 1]);
    }
Vscanf("%*[\n] ");

g08rbc(order, ns, nv, y, p, x, pdx, icen, gamma,
        nmax, tol, parvar, pdparvar, irank, zin, eta, vapvec,
        parest, &fail);
if (fail.code != NE_NOERROR)
    {
        Vprintf("Error from g08rbc.\n%s\n", fail.message);
        exit_status = 1;
        goto END;
    }

Vprintf("Score statistic\n");

for (i = 1; i <= p; ++i)
    Vprintf("%9.3f\n", parest[i - 1]);
Vprintf("\n");

Vprintf("Covariance matrix of score statistic\n");
for (j = 1; j <= p; ++j)
    {
        for (i = 1; i <= j; ++i)
            Vprintf("%9.3f\n", PARVAR(i,j));
        Vprintf("\n");
    }

Vprintf("Parameter estimates\n");
for (i = 1; i <= p; ++i)
    Vprintf("%9.3f\n", parest[p + i - 1]);
Vprintf("\n");
Vprintf("Covariance matrix of parameter estimates\n");
for (i = 1; i <= p; ++i)
    {
        for (j = 1; j <= i; ++j)
            Vprintf("%9.3f\n", PARVAR(i + 1,j));
        Vprintf("\n");
    }

Vprintf("Chi-squared statistic =%9.3f with%2ld d.f.\n", parest[p * 2], p);
Vprintf("\n");

Vprintf("Standard errors of estimates and\n");
Vprintf("approximate z-statistics\n");

for (i = 1; i <= p; ++i)
    Vprintf("%9.3f%14.3f\n", parest[2*p + 1 + i - 1], parest[p * 3 + 1 + i -
1]);
END:
if (eta) NAG_FREE(eta);
if (parest) NAG_FREE(parest);
if (parvar) NAG_FREE(parvar);
if (vapvec) NAG_FREE(vapvec);
if (x) NAG_FREE(x);
if (y) NAG_FREE(y);
if (zin) NAG_FREE(zin);
if (icen) NAG_FREE(icen);
if (irank) NAG_FREE(irank);
if (iwa) NAG_FREE(iwa);
if (nv) NAG_FREE(nv);

return exit_status;
}

```



## 9.2 Program Data

g08rbc Example Program Data

1 1 0.00001 0.00001

40

143.0 0.0 0 164.0 0.0 0 188.0 0.0 0 188.0 0.0 0 190.0 0.0 0  
 192.0 0.0 0 206.0 0.0 0 209.0 0.0 0 213.0 0.0 0 216.0 0.0 0  
 220.0 0.0 0 227.0 0.0 0 230.0 0.0 0 234.0 0.0 0 246.0 0.0 0  
 265.0 0.0 0 304.0 0.0 0 216.0 0.0 1 244.0 0.0 1 142.0 1.0 0  
 156.0 1.0 0 163.0 1.0 0 198.0 1.0 0 205.0 1.0 0 232.0 1.0 0  
 232.0 1.0 0 233.0 1.0 0 233.0 1.0 0 233.0 1.0 0 233.0 1.0 0  
 239.0 1.0 0 240.0 1.0 0 261.0 1.0 0 280.0 1.0 0 280.0 1.0 0  
 296.0 1.0 0 296.0 1.0 0 323.0 1.0 0 204.0 1.0 1 344.0 1.0 1

## 9.3 Program Results

g08rbc Example Program Results

Number of samples = 1

Number of parameters fitted = 1

Distribution power parameter = 0.00001

Tolerance for ties = 0.00001

Score statistic

4.584

Covariance matrix of score statistic

7.653

Parameter estimates

0.599

Covariance matrix of parameter estimates

0.131

Chi-squared statistic = 2.746 with 1 d.f.

Standard errors of estimates and  
 approximate z-statistics

0.361 1.657

---