

NAG C Library Function Document

nag_sum_sqs (g02buc)

1 Purpose

nag_sum_sqs (g02buc) calculates the sample means and sums of squares and cross-products, or sums of squares and cross-products of deviations from the mean, in a single pass for a set of data. The data may be weighted.

2 Specification

```
#include <nag.h>
#include <nagg02.h>

void nag_sum_sqs (Nag_OrderType order, Nag_SumSquare mean, Integer n, Integer m,
                  const double x[], Integer pdx, const double wt[], double *sw,
                  double wmean[], double c[], NagError *fail)
```

3 Description

nag_sum_sqs (g02buc) is an adaptation of West's WV2 algorithm; see West (1979). This function calculates the (optionally weighted) sample means and (optionally weighted) sums of squares and cross-products or sums of squares and cross-products of deviations from the (weighted) mean for a sample of n observations on m variables X_j , for $j = 1, 2, \dots, m$. The algorithm makes a single pass through the data.

For the first $i - 1$ observations let the mean of the j th variable be $\bar{x}_j(i - 1)$, the cross-product about the mean for the j th and k th variables be $c_{jk}(i - 1)$ and the sum of weights be W_{i-1} . These are updated by the i th observation, x_{ij} , for $j = 1, 2, \dots, m$, with weight w_i as follows:

$$\begin{aligned} W_i &= W_{i-1} + w_i \\ \bar{x}_j(i) &= \bar{x}_j(i-1) + \frac{w_i}{W_i}(x_{ij} - \bar{x}_j(i-1)), \quad j = 1, 2, \dots, m \end{aligned}$$

and

$$c_{jk}(i) = c_{jk}(i-1) + \frac{w_i}{W_i}(x_{ij} - \bar{x}_j(i-1))(x_{ik} - \bar{x}_k(i-1))W_{i-1}, \quad j = 1, 2, \dots, m; k = j, j+1, \dots, m.$$

The algorithm is initialized by taking $\bar{x}_j(1) = x_{1j}$, the first observation, and $c_{ij}(1) = 0.0$.

For the unweighted case $w_i = 1$ and $W_i = i$ for all i .

Note that only the upper triangle of the matrix is calculated and returned packed by column.

4 References

Chan T F, Golub G H and Leveque R J (1982) *Updating Formulae and a Pairwise Algorithm for Computing Sample Variances* Compstat, Physica-Verlag

West D H D (1979) Updating mean and variance estimates: An improved method *Comm. ACM* **22** 532–555

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 2.2.1.4 of the Essential Introduction for a more detailed explanation of the use of this argument.

Constraint: **order = Nag_RowMajor or Nag_ColMajor.**

2: **mean** – Nag_SumSquare *Input*

On entry: indicates whether nag_sum_sqs (g02buc) is to calculate sums of squares and cross-products, or sums of squares and cross-products of deviations about the mean.

mean = Nag_AboutMean

The sums of squares and cross-products of deviations about the mean are calculated.

mean = Nag_AboutZero

The sums of squares and cross-products are calculated.

Constraint: **mean = Nag_AboutMean or Nag_AboutZero.**

3: **n** – Integer *Input*

On entry: n , the number of observations in the data set.

Constraint: $n \geq 1$.

4: **m** – Integer *Input*

On entry: m , the number of variables.

Constraint: $m \geq 1$.

5: **x[dim]** – const double *Input*

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

$\max(1, pdx \times m)$ when **order = Nag_ColMajor**;
 $\max(1, n \times pdx)$ when **order = Nag_RowMajor**.

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

if **order = Nag_ColMajor**, $x[(j-1) \times pdx + i - 1]$;
if **order = Nag_RowMajor**, $x[(i-1) \times pdx + j - 1]$.

On entry: $X(i,j)$ must contain the i th observation on the j th variable, for $i = 1, 2, \dots, n$;
 $j = 1, 2, \dots, m$.

6: **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**, $pdx \geq n$;
if **order = Nag_RowMajor**, $pdx \geq m$.

7: **wt[dim]** – const double *Input*

Note: the dimension, dim , of the array **wt** must be at least **n**.

On entry: the optional weights of each observation. If weights are not provided then **wt** must be set to the null pointer, i.e., `(double *)0`, otherwise **wt[i]** must contain the weight for the $i - 1$ th observation.

Constraint: $wt[i] \geq 0.0$ if **wt** is not **NULL**, for $i = 0, 1, \dots, n - 1$.

8: **sw** – double * *Output*

On exit: the sum of weights.

If **wt** is **NULL**, **sw** contains the number of observations, n .

9: **wmean[m]** – double *Output*

On exit: the sample means. **wmean**[$j - 1$] contains the mean for the j th variable.

10: **c[dim]** – double *Output*

Note: the dimension, dim , of the array **c** must be at least $(m \times m + m)/2$.

On exit: the cross-products.

If **mean** = **Nag_AboutMean**, **c** contains the upper triangular part of the matrix of (weighted) sums of squares and cross-products of deviations about the mean.

If **mean** = **Nag_AboutZero**, **c** contains the upper triangular part of the matrix of (weighted) sums of squares and cross-products.

These are stored packed by columns, i.e., the cross-product between the j th and k th variable, $k \geq j$, is stored in **c**($k \times (k - 1)/2 + j$).

11: **fail** – NagError * *Input/Output*

The NAG error argument (see Section 2.6 of 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, **m** = $\langle value \rangle$.

Constraint: **m** ≥ 1 .

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

Constraint: **n** ≥ 1 .

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

Constraint: **pdx** > 0 .

NE_INT_2

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

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

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

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

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.

NE_REAL_ARRAY_ELEM CONS

On entry, **wt**[$\langle value \rangle$] < 0.0.

7 Accuracy

For a detailed discussion of the accuracy of this algorithm see Chan *et al.* (1982) or West (1979).

8 Further Comments

`nag_cov_to_corr` (`g02bw`) may be used to calculate the correlation coefficients from the cross-products of deviations about the mean. The cross-products of deviations about the mean may be scaled to give a variance-covariance matrix.

The means and cross-products produced by `nag_sum_sq` (`g02buc`) may be updated by adding or removing observations using `nag_sum_sq_update` (`g02btc`).

9 Example

A program to calculate the means, the required sums of squares and cross-products matrix, and the variance matrix for a set of 3 observations of 3 variables.

9.1 Program Text

```
/* nag_sum_sq (g02buc) Example Program.
*
* Copyright 2002 Numerical Algorithms Group.
*
* Mark 7, 2002.
*/
#include <stdio.h>
#include <string.h>
#include <nag.h>
#include <nag_stdlib.h>
#include <nagf06.h>
#include <nagg02.h>
#include <nagx04.h>

int main(void)
{
    /* Scalars */
    double alpha, sw;
    Integer exit_status, j, k, m, mm, n, pdx;
    NagError fail;
    Nag_SumSquare mean_enum;

    /* Arrays */
    char mean[2], weight[2];
    double *c=0, *v=0, *wmean=0, *wt=0, *x=0;
    double *wptr=0;

    Nag_OrderType order;

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

    INIT_FAIL(fail);
    exit_status = 0;
    Vprintf("nag_sum_sq (g02buc) Example Program Results\n");

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

    while (scanf(" %1s ' ' %1s '%ld%ld%*[^\n]", mean, weight,
                &m, &n) != EOF)
    {
        /* Allocate memory */
        if ( !(c = NAG_ALLOC((m*m+m)/2, double)) ||
            !(v = NAG_ALLOC((m*m+m)/2, double)) ||
            !(wmean = NAG_ALLOC(m, double)) ||
            !(wptr = NAG_ALLOC(m, double)) )
            fail.code = 'E';
    }
}
```

```

    !(wt = NAG_ALLOC(n, double)) ||
    !(x = NAG_ALLOC(n * m, double)) )
{
    Vprintf("Allocation failure\n");
    exit_status = -1;
    goto END;
}

#ifndef NAG_COLUMN_MAJOR
    pdx = n;
#else
    pdx = m;
#endif

for (j = 1; j <= n; ++j)
    Vscanf("%lf", &wt[j-1]);

Vscanf("%*[^\n] ");
for (j = 1; j <= n; ++j)
{
    for (k = 1; k <= m; ++k)
        Vscanf("%lf", &x(j,k));
}
Vscanf("%*[^\n] ");

if (mean[0] == 'M')
    mean_enum = Nag_AboutMean;
else if (mean[0] == 'Z')
    mean_enum = Nag_AboutZero;
else
{
    Vprintf("Incorrect value for mean\n");
    exit_status = -1;
    goto END;
}
if (weight[0] == 'W')
    wptr = wt;

/* Calculate sums of squares and cross-products matrix */
/* nag_sum_sq (g02buc).
 * Computes a weighted sum of squares matrix
 */
nag_sum_sq(order, mean_enum, n, m, x, pdx, wptr, &sw, wmean, c, &fail);
if (fail.code != NE_NOERROR)
{
    Vprintf("Error from nag_sum_sq (g02buc).\n%s\n", fail.message);
    exit_status = 1;
    goto END;
}

Vprintf("\n");
Vprintf("Means\n");
for (j = 1; j <= m; ++j)
    Vprintf("%14.4f%s", wmean[j-1], j%6 == 0 || j == m ? "\n:" " ");
if (wptr)
{
    Vprintf("\n");
    Vprintf("Weights\n");
    for (j = 1; j <= n; ++j)
        Vprintf("%14.4f%s", wptr[j-1], j%6 == 0 || j == n ? "\n:" " ");
    Vprintf("\n");
}

/* Print the sums of squares and cross products matrix */
/* nag_pack_real_mat_print (x04ccc).
 * Print real packed triangular matrix (easy-to-use)
 */
nag_pack_real_mat_print(Nag_ColMajor, Nag_Upper, Nag_NonUnitDiag, m, c,
                        "Sums of squares and cross-products", 0, &fail);
if (fail.code != NE_NOERROR)
{

```

```

Vprintf("Error from nag_pack_real_mat_print (x04ccc).\n%s\n",
        fail.message);
exit_status = 1;
goto END;
}
if (sw > 1.0)
{
    /* Calculate the variance matrix */
    alpha = 1.0 / (sw - 1.0);
    mm = m * (m + 1) / 2;
    f06fdc(mm, alpha, c, 1, v, 1);
    /* Print the variance matrix */
    Vprintf("\n");
    /* nag_pack_real_mat_print (x04ccc), see above. */
    nag_pack_real_mat_print(Nag_ColMajor, Nag_Upper, Nag_NonUnitDiag, m,
                            v, "Variance matrix", 0, &fail);
    if (fail.code != NE_NOERROR)
    {
        Vprintf("Error from nag_pack_real_mat_print (x04ccc).\n%s\n",
                fail.message);
        exit_status = 1;
        goto END;
    }
}
if (c) NAG_FREE(c);
if (v) NAG_FREE(v);
if (wmean) NAG_FREE(wmean);
if (wt) NAG_FREE(wt);
if (x) NAG_FREE(x);
}

END:
if (c) NAG_FREE(c);
if (v) NAG_FREE(v);
if (wmean) NAG_FREE(wmean);
if (wt) NAG_FREE(wt);
if (x) NAG_FREE(x);

return exit_status;
}

```

9.2 Program Data

```
nag_sum_sq (g02buc) Example Program Data
'M'      'W'      3      3
0.1300   1.3070   0.3700
9.1231   3.7011   4.5230
0.9310   0.0900   0.8870
0.0009   0.0099   0.0999
```

9.3 Program Results

```
nag_sum_sq (g02buc) Example Program Results

Means
     1.3299          0.3334          0.9874

Weights
     0.1300          1.3070          0.3700

Sums of squares and cross-products

          1          2          3
1     8.7569      3.6978      4.0707
2           1.5905      1.6861
3           1.9297

Variance matrix
```

1	10.8512	2	5.0443
2		1.9709	2.0893
3			2.3912
