

NAG C Library Function Document

nag_tabulate_percentile (g11bbc)

1 Purpose

nag_tabulate_percentile (g11bbc) computes a table from a set of classification factors using a given percentile or quantile, for example the median.

2 Specification

```
void nag_tabulate_percentile (Nag_TabulateVar type, Integer n, Integer nfac,
    const Integer sf[], const Integer lfac[], const Integer factor[],
    Integer tdf, double percnt, const double y[], const double wt[],
    double table[], Integer maxt, Integer *ncells, Integer *ndim,
    Integer idim[], Integer count[], NagError *fail)
```

3 Description

A data set may include both classification variables and general variables. The classification variables, known as factors, take a small number of values known as levels. For example, the factor sex would have the levels male and female. These can be coded as 1 and 2 respectively. Given several factors, a multi-way table can be constructed such that each cell of the table represents one level from each factor. For example, the two factors sex and habitat, habitat having three levels: inner-city, suburban and rural, define the 2 by 3 contingency table:

Sex	Habitat		
	Inner-city	Suburban	Rural
Male			
Female			

For each cell statistics can be computed. If a third variable in the data set was age then for each cell the median age could be computed:

Sex	Habitat		
	Inner-city	Suburban	Rural
Male	24	31	37
Female	21.5	28.5	33

That is the median age for all observations for males living in rural areas is 37. The median being the 50% quantile. Other quantiles can also be computed: the p percent quantile or percentile, q_p , is the estimate of the value such that p percent of observations are less than q_p . This is calculated in two different ways depending on whether the tabulated variable is continuous or discrete. Let there be m values in a cell and let $y_{(1)}, y_{(2)}, \dots, y_{(m)}$ be the values for that cell sorted into ascending order. Also, associated with each value there is a weight, $w_{(1)}, w_{(2)}, \dots, w_{(m)}$, which could represent the observed frequency for that value, with $W_j = \sum_{i=1}^j w_{(i)}$ and $W'_j = \sum_{i=1}^j w_{(i)} - \frac{1}{2}w_{(j)}$. For the p percentile let $p_w = (p/100)W_m$ and $p'_w = (p/100)W'_m$ then the percentiles for the two cases are as given below.

If the variable is discrete, that is takes only a limited number of (usually integer) values then the percentile is defined as:

$$y_{(j)} \quad \text{if } W_{j-1} < p_w < W_j$$

$$\frac{y_{(j+1)} + y_{(j)}}{2} \quad \text{if } p_w = W_j$$

If the data is continuous then the quantiles are estimated by linear interpolation.

$$\begin{array}{ll} y_{(1)} & \text{if } p'_w \leq W'_1 \\ (1-f)y_{(j-1)} + fy_{(j)} & \text{if } W'_{j-1} < p'_w \leq W'_j \\ y_{(m)} & \text{if } p'_w > W'_m \end{array}$$

where $f = (p'_w - W'_{j-1}) / (W'_j - W'_{j-1})$.

4 Parameters

- 1: **type** – Nag_TabulateVar *Input*
On entry: indicates whether the variable to be tabulated is discrete or continuous.
 If **type** = **Nag_TabulateVarDiscr**, the percentiles are computed for a discrete variable.
 If **type** = **Nag_TabulateVarCont**, the percentiles are computed for a continuous variable using linear interpolation.
Constraint: **type** = **Nag_TabulateVarDiscr** or **Nag_TabulateVarCont**.
- 2: **n** – Integer *Input*
On entry: the number of observations.
Constraint: **n** \geq 2.
- 3: **nfac** – Integer *Input*
On entry: the number of classifying factors in **factor**.
Constraint: **nfac** \geq 1.
- 4: **sf[nfac]** – const Integer *Input*
On entry: indicates which factors in **factor** are to be used in the tabulation.
 If **sf**[*i* – 1] > 0 the *i*th factor in **factor** is included in the tabulation.
 Note that if **sf**[*i* – 1] \leq 0, for $i = 1, 2, \dots, \mathbf{nfac}$ then the statistic for the whole sample is calculated and returned in a 1 by 1 table.
- 5: **lfac[nfac]** – const Integer *Input*
On entry: the number of levels of the classifying factors in **factor**.
Constraint: if **sf**[*i* – 1] > 0, **lfac**[*i* – 1] \geq 2, for $i = 1, 2, \dots, \mathbf{nfac}$.
- 6: **factor[n][tdf]** – const Integer *Input*
On entry: the **nfac** coded classification factors for the **n** observations.
Constraint: if **sf**[*i* – 1] > 0, $1 \leq \mathbf{factor}[i-1][j-1] \leq \mathbf{lfac}[j-1]$, for $i = 1, 2, \dots, \mathbf{n}$; $j = 1, 2, \dots, \mathbf{nfac}$.
- 7: **tdf** – Integer *Input*
On entry: the second dimension of the array **factor** as declared in the function from which `nag_tabulate_percentile` is called.
Constraint: **tdf** \geq **nfac**.
- 8: **percent** – double *Input*
On entry: the percentile to be tabulated, *p*.
Constraint: $0.0 < \mathbf{percent} < 100.0$.

- 9: **y[n]** – const double *Input*
On entry: the variable to be tabulated.
- 10: **wt[n]** – const double *Input*
On entry: **wt** must contain the **n** weights. Otherwise **wt** must be set to null pointer (double *)0.
Constraint: $\mathbf{wt}[i - 1] \geq 0.0$, for $i = 1, 2, \dots, \mathbf{n}$.
- 11: **table[maxt]** – double *Output*
On exit: the computed table. The **ncells** cells of the table are stored so that for any two factors the index relating to the factor occurring later in **lfac** and **factor** changes faster. For further details see Section 6.
- 12: **maxt** – Integer *Input*
On entry: the maximum size of the table to be computed.
Constraint: $\mathbf{maxt} \geq$ product of the levels of the factors included in the tabulation.
- 13: **ncells** – Integer * *Output*
On exit: the number of cells in the table.
- 14: **ndim** – Integer * *Output*
On exit: the number of factors defining the table.
- 15: **idim[nfac]** – Integer *Output*
On exit: the first **ndim** elements contain the number of levels for the factors defining the table.
- 16: **count[maxt]** – Integer *Output*
On exit: a table containing the number of observations contributing to each cell of the table, stored identically to **table**.
- 17: **fail** – NagError * *Input/Output*
The NAG error parameter (see the Essential Introduction).

5 Error Indicators and Warnings

NE_INT_ARG_LT

On entry, **n** must not be less than 2: **n** = <value>.

On entry, **nfac** must not be less than 1: **nfac** = <value>.

NE_2_INT_ARG_LT

On entry, **tdf** = <value> while **nfac** = <value>.

These parameters must satisfy $\mathbf{tdf} \geq \mathbf{nfac}$.

NE_REAL

On entry, **percent** = <value>.

Constraint: $0.0 < \mathbf{percent} < 100.0$.

NE_BAD_PARAM

On entry, parameter **type** had an illegal value.

NE_REAL_ARRAY_CONS

On entry, **wt**[<value>] = <value>.
 Constraint: **wt**[*i*] ≥ 0, for *i* = 0, 1, ..., *n* - 1.

NE_2_INT_ARRAY_CONS

On entry, **sf**[<value>] = <value> while **lfac**[<value>] = <value>.
 Constraint: if **sf**[*i*] > 0, **lfac**[*i*] ≥ 2, for *i* = 0, 1, ..., **nfac** - 1.

NE_2D_INT_ARRAY_CONS

On entry, **factor**[<value>][<value>] = <value>.
 Constraint: **factor**[*i*][*j*] ≥ 1, for *i* = 0, 1, ..., *n* - 1; *j* = 0, 1, ..., **nfac** - 1.

NE_2D_1D_INT_ARRAYS_CONS

On entry, **factor**[<value>][<value>] = <value> while **lfac**[<value>] = <value>.
 Constraint: **factor**[*i*][*j*] ≤ **lfac**[*j*], for *i* = 0, 1, ..., *n* - 1; *j* = 0, 1, ..., **nfac** - 1.

NE_MAXT

The maximum size of the table to be computed, **maxt** is too small.

NE_CELL_EMPTY

At least one cell is empty.

NE_ALLOC_FAIL

Memory allocation failed.

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.

6 Further Comments

The tables created by `nag_tabulate_percentile` and stored in **table** and **count** are stored in the following way. Let there be *n* factors defining the table with factor *k* having *l_k* levels, then the cell defined by the levels *i₁*, *i₂*, ..., *i_n* of the factors is stored in *m*th cell given by:

$$m = 1 + \sum_{k=1}^n \{(i_k - 1)c_k\},$$

where $c_j = \prod_{k=j+1}^n l_k$, for $j = 1, 2, \dots, n - 1$ and $c_n = 1$.

6.1 Accuracy

Not applicable.

6.2 References

John J A and Quenouille M H (1977) *Experiments: Design and Analysis* Griffin

Kendall M G and Stuart A (1969) *The Advanced Theory of Statistics (Volume 1)* Griffin (3rd Edition)

7 See Also

None.

8 Example

The data, given by John and Quenouille (1977), are for a 3 by 6 factorial experiment in 3 blocks of 18 units. The data is input in the order: blocks, factor with 3 levels, factor with 6 levels, yield, and the 3 by 6 table of treatment medians for yield over blocks is computed and printed.

8.1 Program Text

```

/* nag_tabulate_percentile (g11bbc) Example Program.
 *
 * Copyright 2000 Numerical Algorithms Group.
 *
 * Mark 6, 2000.
 */

#include <stdio.h>
#include <nag.h>
#include <nag_stdlib.h>
#include <nagg11.h>

int main (void)
{
    char type[2], weight[2];
    double percent, *table=0, *wt=0, *wtptr, *y=0;
    Integer items, i, *count=0, *idim=0, *factor=0, ifail, *sf=0, j, k, tdf;
    Integer ltmax, *lfac=0, maxt, n, ncells, ncol, ndim, nfac, nrow;
    Integer exit_status=0;
    Nag_TabulateVar type_enum;
    NagError fail;

#define FACTOR(I,J) factor[((I)-1)*nfac +(J)-1]

    INIT_FAIL(fail);
    Vprintf("g11bbc Example Program Results\n");

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

    Vscanf(" %s %s %ld %ld %lf", type, weight, &n,
    &nfac, &percent);
    ltmax = 18;
    maxt = ltmax;
    if (!(sf = NAG_ALLOC(nfac, Integer))
        || !(lfac = NAG_ALLOC(nfac, Integer))
        || !(idim = NAG_ALLOC(nfac, Integer))
        || !(factor = NAG_ALLOC(n*nfac, Integer))
        || !(count = NAG_ALLOC(maxt, Integer))
        || !(y = NAG_ALLOC(n, double))
        || !(table = NAG_ALLOC(maxt, double))
        || !(wt = NAG_ALLOC(n, double)))
    {
        Vprintf("Allocation failure\n");
        exit_status = -1;
        goto END;
    }

    if (*weight == 'W' || *weight == 'V' )
    {

```

```

        for (i = 1; i <= n; ++i)
    {
        for (j = 1; j <= nfac; ++j)
            Vscanf("%ld", &FACTOR(i,j));
        Vscanf("%lf %lf ", &y[i - 1], &wt[i - 1]);
    }
        wtptr=wt;
    }
    else
    {
        for (i = 1; i <= n; ++i)
    {
        for (j = 1; j <= nfac; ++j)
            Vscanf("%ld", &FACTOR(i,j));
        Vscanf("%lf", &y[i - 1]);
    }
        wtptr = 0;
    }
    for (j = 1; j <= nfac; ++j)
        Vscanf("%ld", &lfac[j - 1]);
    for (j = 1; j <= nfac; ++j)
        Vscanf("%ld", &sf[j - 1]);
    tdf = nfac;
    ifail = 0;
    if (*type == 'D')
        type_enum = Nag_TabulateVarDiscr;
    else if (*type == 'C')
        type_enum = Nag_TabulateVarCont;
    else
        type_enum = (Nag_TabulateVar)-999;

    g11bbc(type_enum, n, nfac, sf, lfac, factor, tdf, percnt, y,
    wtptr, table, maxt, &ncells, &ndim, idim, count,
    &fail);
    if (fail.code != NE_NOERROR)
    {
        Vprintf("Error from g11bbc.\n%s\n", fail.message);
        exit_status = 1;
        goto END;
    }
    Vprintf("\n");
    Vprintf("%s%4.0f%s\n", " Table for ", percnt, "th percentile");
    Vprintf("\n");
    ncol = idim[ndim - 1];
    nrow = ncells / ncol;
    k = 1;
    for (i = 1; i <= nrow; ++i)
    {
        for (items = 1, j = k; j <= k + ncol - 1; ++j, items++)
    {
        Vprintf("%8.2f(%2ld)%s", table[j - 1],
        count[j - 1], items%6?" ":"\n");
    }
        k += ncol;
    }
    END:
    if (sf) NAG_FREE(sf);
    if (lfac) NAG_FREE(lfac);
    if (idim) NAG_FREE(idim);

```

```

    if (factor) NAG_FREE(factor);
    if (count) NAG_FREE(count);
    if (y) NAG_FREE(y);
    if (table) NAG_FREE(table);
    if (wt) NAG_FREE(wt);
    return exit_status;
}

```

8.2 Program Data

g11bbc Example Program Data

C U 54 3 50.0

```

1 1 1 274
1 2 1 361
1 3 1 253
1 1 2 325
1 2 2 317
1 3 2 339
1 1 3 326
1 2 3 402
1 3 3 336
1 1 4 379
1 2 4 345
1 3 4 361
1 1 5 352
1 2 5 334
1 3 5 318
1 1 6 339
1 2 6 393
1 3 6 358
2 1 1 350
2 2 1 340
2 3 1 203
2 1 2 397
2 2 2 356
2 3 2 298
2 1 3 382
2 2 3 376
2 3 3 355
2 1 4 418
2 2 4 387
2 3 4 379
2 1 5 432
2 2 5 339
2 3 5 293
2 1 6 322
2 2 6 417
2 3 6 342
3 1 1 82
3 2 1 297
3 3 1 133
3 1 2 306
3 2 2 352
3 3 2 361
3 1 3 220
3 2 3 333
3 3 3 270

```

3 1 4 388
3 2 4 379
3 3 4 274
3 1 5 336
3 2 5 307
3 3 5 266
3 1 6 389
3 2 6 333
3 3 6 353

3 3 6
0 1 1

8.3 Program Results

g11bbc Example Program Results

Table for 50th percentile

226.00(3)	320.25(3)	299.50(3)	385.75(3)	348.00(3)	334.75(3)
329.25(3)	343.25(3)	365.25(3)	370.50(3)	327.25(3)	378.00(3)
185.50(3)	328.75(3)	319.50(3)	339.25(3)	286.25(3)	350.25(3)
