

Cluster Analysis: nagdmc_nrgp

Purpose

nagdmc_nrgp finds the nearest group to a data record given the group centroids.

Declaration

```
#include <nagdmc.h>

void nagdmc_nrgp(long rec1, long nvar, long nrec, long dblk, double data[],
                 void (*dfun)(long, long, double [], char *, int *), char *comm,
                 long chunksize, long nxvar, long xvar[], long ng, double g[],
                 long ing[], int *info);
```

Parameters

- | | | |
|----|---|---------------------------|
| 1: | rec1 – long | <i>Input</i> |
| | <i>On entry:</i> the index in the data of the first data record used in the analysis. | |
| | <i>Constraint:</i> rec1 ≥ 0 . | |
| 2: | nvar – long | <i>Input</i> |
| | <i>On entry:</i> the number of variables in the data. | |
| | <i>Constraint:</i> nvar ≥ 1 . | |
| 3: | nrec – long | <i>Input</i> |
| | <i>On entry:</i> the number of consecutive records, beginning at rec1 , used in the analysis. | |
| | <i>Constraint:</i> nrec > 1 . | |
| 4: | dblk – long | <i>Input</i> |
| | <i>On entry:</i> the total number of records in the data block. | |
| | <i>Constraint:</i> dblk $\geq \text{rec1} + \text{nrec}$. | |
| 5: | data [dblk * nvar] – double | <i>Input</i> |
| | <i>On entry:</i> the data values for the j th variable (for $j = 0, 1, \dots, \text{nvar} - 1$) are stored in data [$i * \text{nvar} + j$], for $i = 0, 1, \dots, \text{dblk} - 1$. When the data function is used, data is not referenced. | |
| 6: | dfun – function supplied by user | <i>External Procedure</i> |
| | <i>On entry:</i> the pointer to a data function supplied by the user. | |
| | <i>Constraint:</i> if dfun is a valid pointer, data must be 0. | |

The specification of **dfun** is:

<pre>void dfun(long irec, long chunksize, double x[], char *comm, int *ierr)</pre>		
1:	irec – long	<i>Input</i>
	<i>On entry:</i> the index in the data of the first record returned.	
2:	chunksize – long	<i>Input</i>
	<i>On entry:</i> the number of consecutive records returned.	
3:	x [chunksize * nvar] – double	<i>Output</i>
	<i>On exit:</i> data values for the j th variable (for $j = 0, 1, \dots, \text{nvar} - 1$) must be returned in x [$i * \text{nvar} + j$], for $i = 0, 1, \dots, \text{chunksize} - 1$.	
4:	comm – char *	<i>Input</i>
	<i>On entry:</i> a communication parameter allowing additional information to be passed to dfun . This parameter is passed ‘as is’ through the calling function.	

- | | | |
|---|---------------------|---------------|
| 5: | ierr – int * | <i>Output</i> |
| <i>On exit:</i> if the value pointed to by ierr on return is greater than 100, the NAG DMC function will terminate immediately and info will point to this value. | | |
- 7: **comm** – char * *Input*
On entry: a communication parameter allowing additional information to be passed to **dfun**. This parameter is passed ‘as is’ through the calling function.
- 8: **chunksize** – long *Input*
On entry: if the data function is used, the function inputs no more than **chunksize** data records at a time; otherwise **chunksize** is not referenced.
Constraint: if **dfun** \neq 0, **chunksize** \geq 1.
- 9: **nxvar** – long *Input*
On entry: the number of variables in the analysis. If **nxvar** = 0, all variables in the data are used in the analysis.
Constraint: $0 \leq \mathbf{nxvar} \leq \mathbf{nvar}$.
- 10: **xvar[nxvar]** – long *Input*
On entry: the indices indicating the position in **data** in which the variables are stored. If **nxvar** = 0 then **xvar** must be 0, and the indices of variables are given by $j = 0, 1, \dots, \mathbf{nvar} - 1$.
Constraints: if **nxvar** $>$ 0, $0 \leq \mathbf{xvar}[i] < \mathbf{nvar}$, for $i = 0, 1, \dots, \mathbf{nxvar} - 1$; otherwise **xvar** must be 0.
- 11: **ng** – long *Input*
On entry: the number of groups in the clustering.
Constraint: **ng** $>$ 1.
- 12: **g[ng*nvar]** – double *Input*
On entry: **g**[$i * \mathbf{nvar} + j$] contains the mean value for the j th variable of the i th group, for $j = 0, 1, \dots, \mathbf{nvar} - 1$; for $i = 0, 1, \dots, \mathbf{ng} - 1$. Note that the value corresponding to the weights, if any, will be ignored.
- 13: **ing[nrec]** – long *Output*
On exit: **ing**[i] is the nearest group to the i th data record in the analysis, for $i = 0, 1, \dots, \mathbf{nrec} - 1$.
- 14: **info** – int * *Output*
On exit: **info** gives information on the success of the function call:
- 0: the function successfully completed its task.
 - i ; $i = 1, 2, 3, 4, 6, 8, 9, \dots, 11$: the specification of the i th formal parameter was incorrect.
 - 99: the function failed to allocate enough memory.
 - > 100 : an error occurred in a function specified by the user.

Notation

- nrec** the number of data records, n .
data the data set X .
nxvar determines the number of variables in the analysis.
ng the number of groups in the clustering.
g the vectors of group centroids c_k , for $k = 1, 2, \dots, l$.
ing the allocation, $a_i - 1$, of data records to groups, for $i = 1, 2, \dots, n$.

Description

Let X be a set of n data records x_i on p variables, for $i = 1, 2, \dots, n$, and c_k be a user-supplied vector of p elements that defines the centroid of group k . Given the centroids of a clustering containing l

groups, the Euclidean distance, d_{ik} , from the i th data record to the k th centroid is:

$$d_{ik} = \left[\sum_{j=1}^l (x_{ij} - c_{kj})^2 \right]^{\frac{1}{2}}, \quad i = 1, 2, \dots, n,$$

where x_{ij} and c_{kj} are the values of the i th data record and k th centroid on variable j , respectively.

The i th data record is allocated to the group number a_i with the minimum distance in d_{ik} , for $k = 1, 2, \dots, l$.

References and Further Reading

None.

See Also

None.
