

## Decision Tree: nagdmc\_predict\_reg\_tree

### Purpose

**nagdmc\_predict\_reg\_tree** uses a decision tree computed by **nagdmc\_reg\_tree** to predict values of data records.

### Declaration

```
#include <nagdmc.h>

void nagdmc_predict_reg_tree(long rec1, long nvar, long nrec, long dblk,
                             double data[], long bcat[], long iproot, int optrand,
                             long iseed, double res[], double acc[], int *info);
```

### Parameters

- 1: **rec1** – long *Input*  
*On entry:* the index in the data of the first data record used in the analysis.  
*Constraint:* **rec1**  $\geq 0$ .
- 2: **nvar** – long *Input*  
*On entry:* the number of variables in the data.  
*Constraint:* **nvar**  $> 1$ .
- 3: **nrec** – long *Input*  
*On entry:* the number of consecutive records, beginning at **rec1**, used in the analysis.  
*Constraint:* **nrec**  $> 1$ .
- 4: **dblk** – long *Input*  
*On entry:* the total number of records in the data block.  
*Constraint:* **dblk**  $\geq \text{rec1} + \text{nrec}$ .
- 5: **data**[**dblk** \* **nvar**] – double *Input*  
*On entry:* 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  $i = 0, 1, \dots, \text{dblk} - 1$ .
- 6: **bcat**[**nvar**] – long *Input*  
*On entry:* **bcat**[ $i$ ] contains the base level value for the **ncat**[ $i$ ] categories on the  $i$ th variable. If **ncat**[ $i$ ]  $> 0$ , for  $i = 0, 1, \dots, \text{nvar} - 1$ , the categorical values on the  $i$ th variable are given by **bcat**[ $i$ ] +  $j$ , for  $j = 0, 1, \dots, \text{ncat}[i] - 1$ ; otherwise **bcat**[ $i$ ] is not referenced. If the base level for each categorical variable is zero, **bcat** can be 0.
- 7: **iproot** – long *Input*  
*On entry:*
- 8: **optrand** – int *Input*  
*On entry:* if the value of **optrand** is set equal to 1, a random number will be used to resolve dichotomies in the decision tree; otherwise **optrand** must be set equal to 0 and some predictions may be set equal to the dummy value  $-1$ .  
*Constraint:* **optrand**  $\in \{0, 1\}$ .
- 9: **iseed** – long *Input*  
*On entry:* if **optrand** = 1, the initial values used to set the seed of the random number generator used to resolve any dichotomies in the tree; otherwise **optrand** is not referenced.
- 10: **res**[**nrec**] – double *Output*  
*On exit:* **res**[ $i$ ] contains the decision tree prediction for the (**rec1** +  $i$ )th data record, for  $i = 0, 1, \dots, \text{nrec} - 1$ .

- 11: **acc[nrec]** – double *Output*  
*On exit:* **acc**[*i*] contains the variance about the mean value, based on the training data, at the leaf node giving the *i*th prediction, for  $i = 0, 1, \dots, \mathbf{nrec} - 1$ .
- 12: **info** – int \* *Output*  
*On exit:* **info** gives information on the success of the function call:
- 0: the function successfully completed its task.
  - 32: a path down the decision tree could not be found for at least one data record, consequently not all data records have been classified; this warning can be avoided by setting **optrand** equal to one.
  - i*;  $i = 1, 2, 3, 4, 8$ : 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 used to predict values,  $n$ .  
**data** data records  $x_i$ , for  $i = 1, 2, \dots, n$ .  
**res** decision tree predictions  $y_i$ , for  $i = 1, 2, \dots, n$ .  
**acc** accuracy of predictions  $v_i$ , for  $i = 1, 2, \dots, n$ .

## Description

Let  $x_i$ , for  $i = 1, 2, \dots, n$  be a set of  $n$  data records not used to fit a decision tree,  $T$ . The *i*th prediction for the dependent variable in the data is found by using the outcome of a series of tests at the root node and internal nodes in  $T$  to associate  $x_i$  with leaf node  $l_i$ , for  $i = 1, 2, \dots, n$ . The value of the dependent variable stored at  $l_i$  is then used as the predicted value  $y_i$ , for  $i = 1, 2, \dots, n$ . In a regression decision tree each leaf node stores the mean of the dependent variable over a subset of the data records.

The outcome of each test depends on the type of variable used to partition data records at the node. Let a test at a node  $k$  be on variable  $j$  in the data and  $x_{ij}$  be the value of the *i*th data record on variable  $j$ .

If  $j$  is continuous,  $x_i$  is sent to the left child node of node  $k$  if  $x_{ij} \leq t$ , where  $t$  is the value of the continuous test as stored in node  $k$ ; otherwise  $x_i$  is sent to the right child node of node  $k$ .

If  $j$  is categorical,  $x_i$  is sent to the node associated with the category value  $x_{ij}$ . However, when the decision was fitted there may not have been a category value  $x_{ij}$  at node  $k$  and, therefore, either the *i*th data record can be assigned an unclassified value or a child node can be chosen at random from those available to node  $k$ .

This process of evaluating tests continues until  $x_i$  reaches a leaf node, say  $l_i$ , in  $T$ .

A measure of the accuracy of the *i*th prediction can be obtained by considering the variance,  $v_i$ , of the mean value of the dependent variable for data records at leaf node  $l_i$  (and used to fit  $T$ ), for  $i = 1, 2, \dots, n$ .

## References and Further Reading

None.

## See Also

- [nagdmc\\_srs](#) sets the seed for the random number generator.  
[reg\\_tree\\_ex.c](#) the example calling program.