

Decision Tree: nagdmc_predict_entropy_tree

Purpose

nagdmc_predict_entropy_tree predicts values for new data given a decision tree computed by **nagdmc_entropy_tree**.

Declaration

```
#include <nagdmc.h>
void nagdmc_predict_entropy_tree(long rec1, long nvar, long nrec, long dblk,
                                double data[], long iproot, int optrand,
                                long izeed, long res[], double prob[],
                                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 rec1 + nrec . | |
| 5: | data [dblk * nvar] – double | <i>Input</i> |
| | <i>On entry:</i> the data values for the j th variable (for $j = 0, 1, \dots, \mathbf{nvar} - 1$) are stored in data [$i * \mathbf{nvar} + j$], for $i = 0, 1, \dots, \mathbf{dblk} - 1$. | |
| 6: | iproot – long | <i>Input</i> |
| | <i>On entry:</i> the integer value of the root node of a decision tree as returned by nagdmc_entropy_tree . | |
| 7: | optrand – int | <i>Input</i> |
| | <i>On entry:</i> 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 data records may be unclassified, i.e., will be classified as -1. | |
| | <i>Constraint:</i> optrand $\in \{0, 1\}$. | |
| 8: | izeed – long | <i>Input</i> |
| | <i>On entry:</i> 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 izeed is not referenced. | |
| 9: | res [nrec] – long | <i>Output</i> |
| | <i>On exit:</i> res [i] contains the decision tree classification for $i = 0, 1, \dots, \mathbf{nrec} - 1$. | |
| 10: | prob [nrec] – double | <i>Output</i> |
| | <i>On exit:</i> prob [i] contains the probability of res [i] given the training data, for $i = 0, 1, \dots, \mathbf{nrec} - 1$. | |
| 11: | info – int * | <i>Output</i> |
| | <i>On exit:</i> info gives information on the success of the function call: | |
| | 0: the function successfully completed its task. | |
| | i ; $i = 1, 2, 3, 4, 7$: the specification of the i th formal parameter was incorrect. | |
| | 99: the function failed to allocate enough memory. | |
| | 100: an internal error occurred during the execution of the function. | |

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 classifications y_i , for $i = 1, 2, \dots, n$.
prob accuracy of classifications a_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 decision tree calculated by using an entropy criterion each leaf node stores the modal class 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 class distribution of data records at leaf node l_i , for $i = 1, 2, \dots, n$. Suppose that r_i of m_i data records associated with l_i (and used to fit T) belong to the modal class, then a measure of the accuracy a_i of the classification is given by,

$$a_i = \frac{r_i}{m_i}, \quad i = 1, 2, \dots, n.$$

References and Further Reading

None.

See Also

- | | |
|---|---|
| nagdmc_entropy_tree | computes an decision tree by using an entropy-based criterion. |
| nagdmc_free_entropy_tree | returns to the operating system memory used by an entropy tree. |
| nagdmc_load_entropy_tree | loads an entropy tree from a file. |
| nagdmc_prune_entropy_tree | prunes an entropy tree using pessimistic error pruning. |
| nagdmc_save_entropy_tree | writes an entropy tree to a file. |
| entropy_tree_ex.c | the example calling program. |