

Package ‘HPdclassifier’

January 20, 2015

Type Package

Title Distributed classifiers for Big Data

Version 1.0.0

Date 2015-01-16

Author HP Vertica Analytics Team

Maintainer HP Vertica Analytics Team <distributedRTeam@external.hp.com>

Depends R (>= 3.0.0), distributedR, randomForest (>= 4.6-10)

Description

Distributed algorithms for learning classifiers. Written using HP Vertica Distributed R package.

License GPL (>= 2) | file LICENSE

R topics documented:

HPdclassifier-package	1
hpdrandomForest	2
predictHPdRF	6
Index	9

HPdclassifier-package

Distributed algorithms for classifiers

Description

HPdclassifier provides several distributed algorithms for classifiers. It is written based on the infrastructure created in HP-Labs for distributed computing in R.

Details

Package: HPdclassifier
Type: Package
Version: 1.0.0
Date: 2015-01-16

Main Functions:

- hpdrandomForest: It is a distributed function for randomForest
- predictHPdRF: distributed predict method for applying a random forest objects on a darray or a dframe

Author(s)

HP Vertica Analytics Team <distributedRTeam@external.groups.hp.com>

References

1. Using R for Iterative and Incremental Processing. Shivaram Venkataraman, Indrajit Roy, Alvin AuYoung, Rob Schreiber. HotCloud 2012, Boston, USA.

hpdrandomForest	<i>Distributed randomForest</i>
-----------------	---------------------------------

Description

hpdrandomForest function runs randomForest function of randomForest package in a distributed fashion.

Description

hpdrandomForest calls several instances of randomForest distributed across a cluster system. Therefore, the master distributes the input data among all R-executors of the distributedR environment, and trees on different sub-sections of the forest are created simultaneously. At the end, all these trees are combined to result a single forest.

The interface of hpdrandomForest is similar to randomForest. Indeed it adds two arguments nExecutor and trace, and removes several other arguments do.trace, corr.bias, keep.inbag, and oob.prox. Its returned result is also completely compatible to the result of randomForest.

Usage

```
## S3 method for class 'formula'
hpdrandomForest(formula, data=NULL, ..., ntree=500,
                 na.action=na.fail, nExecutor, trace=FALSE,
                 completeModel=FALSE)
## Default S3 method:
hpdrandomForest(x, y=NULL, xtest=NULL, ytest=NULL, ntree=500,
                 mtry=if (!is.null(y) && !is.factor(y) && !is.dframe(y))
                 max(floor(ncol(x)/3), 1) else floor(sqrt(ncol(x))),
                 replace=TRUE, classwt=NULL, cutoff, strata,
                 sampsize = if (replace) nrow(x) else ceiling(.632*nrow(x)),
                 nodesize = if (!is.null(y) && !is.factor(y) &&
                 !is.dframe(y)) 5 else 1,
                 maxnodes=NULL, importance=FALSE, localImp=FALSE, nPerm=1,
                 proximity=FALSE, norm.votes=TRUE, keep.forest=TRUE,
                 nExecutor, trace=FALSE, completeModel=FALSE, ...)
```

Arguments

<code>data</code>	a data frame or dframe which contains samples.
<code>na.action</code>	A function to specify the action to be taken if NAs are found. (NOTE: If given, this argument must be named.)
<code>formula</code>	a formula describing the model to be fitted.
<code>x</code>	when a data frame or a matrix of predictors assigned to <code>x</code> , its size should not be bigger than 2GB. For bigger datasets, <code>darray</code> should be used. <code>darray</code> does not support categorical data. Therefore, <code>dframe</code> and the first interface should be used for classification problems of large datasets with categorical data.
<code>y</code>	a response vector. If a factor, classification is assumed, otherwise regression is assumed. If omitted, <code>hpdrandomForest</code> will run in unsupervised mode. When <code>x</code> is a <code>darray</code> , <code>y</code> should be also a <code>darray</code> with a single column.
<code>xtest</code>	a data frame or matrix (like <code>x</code>) containing predictors for the test set. When <code>x</code> is a <code>darray</code> , it should be of the same type.
<code>ytest</code>	response for the test set. Its type should be consistent with <code>y</code> . Moreover, it should have a single column.
<code>ntree</code>	Number of trees to grow. This should not be set to too small a number, to ensure that every input row gets predicted at least a few times.
<code>mtry</code>	Number of variables randomly sampled as candidates at each split. Note that the default values are different for classification (\sqrt{p}) where p is number of variables in <code>x</code>) and regression ($p/3$)
<code>replace</code>	Should sampling of cases be done with or without replacement?
<code>classwt</code>	Priors of the classes. Need not add up to one. Ignored for regression.
<code>cutoff</code>	(Classification only) A vector of length equal to number of classes. The 'winning' class for an observation is the one with the maximum ratio of proportion of votes to cutoff. Default is $1/k$ where k is the number of classes (i.e., majority vote wins).
<code>strata</code>	A (factor) variable that is used for stratified sampling.
<code>sampsiz</code>	Size(s) of sample to draw. For classification, if <code>sampsiz</code> is a vector of the length the number of strata, then sampling is stratified by strata, and the elements of <code>sampsiz</code> indicate the numbers to be drawn from the strata.
<code>nodesize</code>	Minimum size of terminal nodes. Setting this number larger causes smaller trees to be grown (and thus take less time). Note that the default values are different for classification (1) and regression (5).
<code>maxnodes</code>	Maximum number of terminal nodes trees in the forest can have. If not given, trees are grown to the maximum possible (subject to limits by <code>nodesize</code>). If set larger than maximum possible, a warning is issued.
<code>importance</code>	Should importance of predictors be assessed?
<code>localImp</code>	Should casewise importance measure be computed? (Setting this to <code>TRUE</code> will override <code>importance</code> .)
<code>nPerm</code>	Number of times the OOB data are permuted per tree for assessing variable importance. Number larger than 1 gives slightly more stable estimate, but not very effective. Currently only implemented for regression.
<code>proximity</code>	a logical value which indicates if the proximity measure among the rows should be calculated. It is <code>FALSE</code> by default because it is very memory inefficient. Moreover, it is calculated only on 'out-of-bag' data

<code>norm.votes</code>	If TRUE (default), the final result of votes are expressed as fractions. If FALSE, raw vote counts are returned (useful for combining results from different runs). Ignored for regression.
<code>keep.forest</code>	If set to FALSE, the forest will not be retained in the output object.
<code>nExecutor</code>	a positive integer number indicating the number of tasks for running the function. To have optimal performance, it is recommended to have this number smaller than the number of R-executors in the distributedR environment. It cannot be bigger than <code>ntree</code> .
<code>trace</code>	when this argument is true, intermediate steps of the progress are displayed.
<code>completeModel</code>	when it is FALSE (default), the output values that preserve information per sample are discarded. They are 'y', 'oob.times', 'votes', 'predicted', and 'test'. This feature is intended to keep the size of the output model small. \item...optional parameters to be passed to the low level function.

Value

An object of class `randomForest`. The result is similar to the result of the `combine` function in `randomForest` package and will contain the following components.

<code>call</code>	the original call to <code>hpdrandomForest</code>
<code>type</code>	one of <code>regression</code> , <code>classification</code> , or <code>unsupervised</code> .
<code>predicted</code>	(only when <code>completeModel=TRUE</code>) the predicted values of the input data based on out-of-bag samples.
<code>importance</code>	a matrix with <code>nclass + 2</code> (for classification) or two (for regression) columns. For classification, the first <code>nclass</code> columns are the class-specific measures computed as mean decrease in accuracy. The <code>nclass + 1</code> st column is the mean decrease in accuracy over all classes. The last column is the mean decrease in Gini index. For Regression, the first column is the mean decrease in accuracy and the second the mean decrease in MSE. If <code>importance=FALSE</code> , the last measure is still returned as a vector.
<code>importanceSD</code>	The "standard errors" of the permutation-based importance measure. For classification, a <code>p</code> by <code>nclass + 1</code> matrix corresponding to the first <code>nclass + 1</code> columns of the importance matrix. For regression, a length <code>p</code> vector.
<code>localImp</code>	a <code>p</code> by <code>n</code> matrix containing the casewise importance measures, the <code>[i,j]</code> element of which is the importance of <code>i</code> -th variable on the <code>j</code> -th case. NULL if <code>localImp=FALSE</code> .
<code>ntree</code>	number of trees grown.
<code>mtry</code>	number of predictors sampled for splitting at each node.
<code>forest</code>	(a list that contains the entire forest; NULL if <code>hpdrandomForest</code> is run in unsupervised mode or if <code>keep.forest=FALSE</code>).
<code>err.rate</code>	(classification only) vector error rates of the prediction on the input data, the <code>i</code> -th element being the (OOB) error rate for all trees up to the <code>i</code> -th.
<code>confusion</code>	(classification only) the confusion matrix of the prediction (based on OOB data).
<code>votes</code>	(classification only, and only when <code>completeModel=TRUE</code>) a matrix with one row for each input data point and one column for each class, giving the fraction or number of (OOB) 'votes' from the random forest.
<code>oob.times</code>	(only when <code>completeModel=TRUE</code>) number of times cases are 'out-of-bag' (and thus used in computing OOB error estimate)

<code>y</code>	(only when <code>completeModel=TRUE</code>) the response vector if it is made available in the input.
<code>proximity</code>	(only when <code>completeModel=TRUE</code>) a matrix of proximity measures among the input (based on the frequency that pairs of data points are in the same terminal nodes) when <code>proximity=TRUE</code> .
<code>mse</code>	(regression only) vector of mean square errors: sum of squared residuals divided by <code>n</code> .
<code>rsq</code>	(regression only) "pseudo R-squared": $1 - \text{mse} / \text{Var}(y)$.
<code>test</code>	(only when <code>completeModel=TRUE</code>) if test set is given (through the <code>xtest</code> or additionally <code>ytest</code> arguments), this component is a list which contains the corresponding predicted, <code>err.rate</code> , <code>confusion</code> , <code>votes</code> (for classification) or predicted, <code>mse</code> and <code>rsq</code> (for regression) for the test set.
<code>terms</code>	it contains a formula identifying response and predictors (for classification and regression types).

Note

When `ntree` is not big enough in comparison to `nExecutor`, some of the returned predicted values may become `NULL`. It is the same for values of 'votes' matrix when they are normalized (`norm.votes=TRUE`). Returned values for `err.rate`, `votes`, and `oob.times` are valid only for classification type.

Three scenarios can be imagined for the type of input data. When ordinary R types are used (matrix or `data.frame`) the behavior is similar to the `randomForest` function; however, the total size of the input data cannot be bigger than 2GB. In fact, for bigger data size distributed data types; i.e., `darray` or `dframe`, should be used. When `x` is of type `darray`, in the case of existence `y` must be of type `darray` as well. Regarding the fact that `darray` does not support categorical data, this data type cannot be used for classification mode. When `x` is of type `dframe`, no value can be assigned to `y`; indeed for this data type, the formula interface should be used for classification and the default interface for unsupervised mode.

Author(s)

HP Vertica Analytics Team

References

Breiman, L. (2001), *Random Forests*, Machine Learning 45(1),5-32.

Breiman, L (2002), "Manual On Setting Up, Using, And Understanding Random Forests V3.1", http://oz.berkeley.edu/users/breiman/Using_random_forests_V3.1.pdf.

Random Forests V4.6-10, <http://cran.r-project.org/web/packages/randomForest/randomForest.pdf>.

Examples

```
## Not run:

library(HPdclassifier)
distributedR_start()
drs <- distributedR_status()
nparts <- sum(drs$Ins)
```

```

## Classification:
##data(iris)
iris.rf <- hpdrandomForest(Species ~ ., data=iris, importance=TRUE,
                           nExecutor=nparts)
print(iris.rf)

## The 'unsupervised' case:
iris.urf <- hpdrandomForest(iris[, -5], nExecutor=nparts,
                           proximity=TRUE, completeModel=TRUE)
MDSplot(iris.urf, iris$Species)

## stratified sampling: draw 20, 30, and 20 of the species to grow each tree.
(iris.rf2 <- hpdrandomForest(iris[1:4], iris$Species,
                             sampsize=c(20, 30, 20), nExecutor=nparts))

## Regression:
## data(airquality)
ozone.rf <- hpdrandomForest(Ozone ~ ., data=airquality, mtry=3,
                            importance=TRUE, na.action=na.omit,
                            nExecutor=nparts, completeModel=TRUE)
print(ozone.rf)
## Show "importance" of variables: higher value mean more important:
round(importance(ozone.rf), 2)

## "x" can be a matrix instead of a data frame:
x <- matrix(runif(5e2), 100)
y <- gl(2, 50)
(myrf <- hpdrandomForest(x, y, nExecutor=nparts))
(predict(myrf, x))

## "complicated" formula:
(swiss.rf <- hpdrandomForest(sqrt(Fertility)~. - Catholic + I(Catholic<50),
                             data=swiss, nExecutor=nparts))
(predict(swiss.rf, swiss))
## Test use of 32-level factor as a predictor:
x <- data.frame(x1=gl(32, 5), x2=runif(160), y=rnorm(160))
(rf1 <- hpdrandomForest(x[-3], x[[3]], ntree=10, nExecutor=nparts))

## Grow no more than 4 nodes per tree:
(treesize(hpdrandomForest(Species ~ ., data=iris, maxnodes=4, ntree=30,
                           nExecutor=nparts)))

distributedR_shutdown()

## End(Not run)

```

predictHPdRF

distributed predict method for applying a random forest objects on a darray or a dframe

Description

Prediction of distributed test data using random forest.

Usage

```
predictHPdRF(object, newdata, trace=FALSE)
```

Arguments

object	an object of class <code>randomForest</code> , as that created by the function <code>randomForest</code> or <code>hpdrandomForest</code> .
newdata	a darray or a dframe containing new data. darray is highly recommended when there is no categorical data
trace	when this argument is true, intermediate steps of the progress are displayed.

Value

It returns a darray or a dframe of predicted classes. When the newdata is of type darray, the type of returned value will be also darray unless the output is categorical data. When the output is a dframe when the newdata is of type dframe.

Author(s)

HP Vertica Analytics Team

References

Breiman, L. (2001), *Random Forests*, Machine Learning 45(1), 5-32.

See Also

[hpdrandomForest](#)

Examples

```
## Not run:
# example for darray
library(HPdclassifier)
distributedR_start()
drs <- distributedR_status()
nparts <- sum(drs$Ins)

nSamples <- 100
nAttributes <- 5
nPartitions <- 2

dax <- darray(c(nSamples, nAttributes),
              c(ceiling(nSamples/nPartitions), nAttributes))
day <- darray(c(nSamples, 1), c(ceiling(nSamples/nPartitions), 1))

foreach(i, 1:npartitions(dax), function(x=splits(dax, i),
                                                  y=splits(day, i), id=i) {
  x <- matrix(runif(nrow(x)*ncol(x)), nrow(x), ncol(x))
  y <- matrix(runif(nrow(y)), nrow(y), 1)
  update(x)
  update(y)
})

(myrf1 <- hpdrandomForest(dax, day, nExecutor=nparts))
```

```

dp <- predictHPdRF(myrf1, dax)

# example for dframe
nSamples <- nrow(iris)
nAttributes <- ncol(iris)
nPartitions <- 4

df <- dframe(c(nSamples,nAttributes),
             c(ceiling(nSamples/nPartitions),nAttributes))
end = cumsum(partitionsizes(df)[,1])
start = c(0,end[-length(end)])

foreach(i, 1:npartitions(df), function(xi=splits(df,i),
  start = start[i]+1, end = end[i], iris=iris){
  xi <- iris[start:end,]
  update(xi)
})
# the following line will be redundant in the next release
colnames(df) <- colnames(iris)

(myrf2 <- hpdrandomForest(Species ~ ., df, nExecutor=nPartitions))
fp2 <- predictHPdRF(myrf2, df)

## End(Not run)

```


Index

- *Topic **Big Data Analytics**
 - HPdclassifier-package, [1](#)
 - hpdrandomForest, [2](#)
- *Topic **Distributed R**
 - HPdclassifier-package, [1](#)
- *Topic **Scalable Machine Learning algorithms**
 - HPdclassifier-package, [1](#)
- *Topic **classification**
 - predictHPdRF, [6](#)
- *Topic **distributed R**
 - hpdrandomForest, [2](#)
- *Topic **distributed random forest**
 - hpdrandomForest, [2](#)
- *Topic **regression**
 - predictHPdRF, [6](#)

HPdclassifier
(*HPdclassifier-package*), [1](#)

HPdclassifier-package, [1](#)

hpdrandomForest, [2](#), [7](#)

predictHPdRF, [6](#)