

Package ‘ranger’

January 10, 2020

Type Package

Title A Fast Implementation of Random Forests

Version 0.12.1

Date 2020-01-10

Author Marvin N. Wright [aut, cre], Stefan Wager [ctb], Philipp Probst [ctb]

Maintainer Marvin N. Wright <cran@wrig.de>

Description A fast implementation of Random Forests, particularly suited for high dimensional data. Ensembles of classification, regression, survival and probability prediction trees are supported. Data from genome-wide association studies can be analyzed efficiently. In addition to data frames, datasets of class 'gwa.data' (R package 'GenABEL') and 'dgCMatrix' (R package 'Matrix') can be directly analyzed.

License GPL-3

Imports Rcpp (>= 0.11.2), Matrix

LinkingTo Rcpp, RcppEigen

Depends R (>= 3.1)

Suggests survival, testthat

Encoding UTF-8

RoxygenNote 7.0.2

URL <https://github.com/imbs-hl/ranger>

BugReports <https://github.com/imbs-hl/ranger/issues>

NeedsCompilation yes

Repository CRAN

Date/Publication 2020-01-10 17:50:05 UTC

R topics documented:

csrf	2
getTerminalNodeIDs	3

holdoutRF	4
importance.ranger	5
importance_pvalues	5
parse.formula	7
predict.ranger	8
predict.ranger.forest	10
predictions.ranger	12
predictions.ranger.prediction	13
print.ranger	13
print.ranger.forest	14
print.ranger.prediction	14
ranger	15
timepoints.ranger	22
timepoints.ranger.prediction	23
treeInfo	23

Index	25
--------------	-----------

csrf	<i>Case-specific random forests.</i>
------	--------------------------------------

Description

In case-specific random forests (CSRF), random forests are built specific to the cases of interest. Instead of using equal probabilities, the cases are weighted according to their difference to the case of interest.

Usage

```
csrf(formula, training_data, test_data, params1 = list(), params2 = list())
```

Arguments

formula	Object of class formula or character describing the model to fit.
training_data	Training data of class data.frame.
test_data	Test data of class data.frame.
params1	Parameters for the proximity random forest grown in the first step.
params2	Parameters for the prediction random forests grown in the second step.

Details

The algorithm consists of 3 steps:

1. Grow a random forest on the training data
2. For each observation of interest (test data), the weights of all training observations are computed by counting the number of trees in which both observations are in the same terminal node.

3. For each test observation, grow a weighted random forest on the training data, using the weights obtained in step 2. Predict the outcome of the test observation as usual.

In total, $n+1$ random forests are grown, where n is the number observations in the test dataset. For details, see Xu et al. (2014).

Value

Predictions for the test dataset.

Author(s)

Marvin N. Wright

References

Xu, R., Nettleton, D. & Nordman, D.J. (2014). Case-specific random forests. *J Comp Graph Stat* 25:49-65. <https://doi.org/10.1080/10618600.2014.983641>.

Examples

```
## Split in training and test data
train.idx <- sample(nrow(iris), 2/3 * nrow(iris))
iris.train <- iris[train.idx, ]
iris.test <- iris[-train.idx, ]

## Run case-specific RF
csrf(Species ~ ., training_data = iris.train, test_data = iris.test,
     params1 = list(num.trees = 50, mtry = 4),
     params2 = list(num.trees = 5))
```

getTerminalNodeIDs *Get terminal node IDs (deprecated)*

Description

This function is deprecated. Please use `predict()` with `type = "terminalNodes"` instead. This function calls `predict()` now.

Usage

```
getTerminalNodeIDs(rf, dat)
```

Arguments

rf	ranger object.
dat	New dataset. Terminal node IDs for this dataset are obtained.

Value

Matrix with terminal nodeIDs for all observations in dataset and trees.

Examples

```
rf <- ranger(Species ~ ., data = iris, num.trees = 5, write.forest = TRUE)
getTerminalNodeIDs(rf, iris)
```

holdoutRF	<i>Hold-out random forests</i>
-----------	--------------------------------

Description

Grow two random forests on two cross-validation folds. Instead of out-of-bag data, the other fold is used to compute permutation importance. Related to the novel permutation variable importance by Janitza et al. (2015).

Usage

```
holdoutRF(...)
```

Arguments

... All arguments are passed to `ranger()` (except importance, case.weights, replace and holdout.).

Value

Hold-out random forests with variable importance.

Author(s)

Marvin N. Wright

References

Janitza, S., Celik, E. & Boulesteix, A.-L., (2015). A computationally fast variable importance test for random forests for high-dimensional data. *Adv Data Anal Classif* <https://doi.org/10.1007/s11634-016-0276-4>.

See Also

[ranger](#)

importance.ranger	<i>ranger variable importance</i>
-------------------	-----------------------------------

Description

Extract variable importance of ranger object.

Usage

```
## S3 method for class 'ranger'  
importance(x, ...)
```

Arguments

x	ranger object.
...	Further arguments passed to or from other methods.

Value

Variable importance measures.

Author(s)

Marvin N. Wright

See Also

[ranger](#)

importance_pvalues	<i>ranger variable importance p-values</i>
--------------------	--

Description

Compute variable importance with p-values. For high dimensional data, the fast method of Janitza et al. (2016) can be used. The permutation approach of Altmann et al. (2010) is computationally intensive but can be used with all kinds of data. See below for details.

Usage

```
importance_pvalues(  
  x,  
  method = c("janitza", "altmann"),  
  num.permutations = 100,  
  formula = NULL,  
  data = NULL,  
  ...  
)
```

Arguments

x	ranger or holdoutRF object.
method	Method to compute p-values. Use "janitza" for the method by Janitza et al. (2016) or "altmann" for the non-parametric method by Altmann et al. (2010).
num.permutations	Number of permutations. Used in the "altmann" method only.
formula	Object of class formula or character describing the model to fit. Used in the "altmann" method only.
data	Training data of class data.frame or matrix. Used in the "altmann" method only.
...	Further arguments passed to ranger(). Used in the "altmann" method only.

Details

The method of Janitza et al. (2016) uses a clever trick: With an unbiased variable importance measure, the importance values of non-associated variables vary randomly around zero. Thus, all non-positive importance values are assumed to correspond to these non-associated variables and they are used to construct a distribution of the importance under the null hypothesis of no association to the response. Since only the non-positive values of this distribution can be observed, the positive values are created by mirroring the negative distribution. See Janitza et al. (2016) for details.

The method of Altmann et al. (2010) uses a simple permutation test: The distribution of the importance under the null hypothesis of no association to the response is created by several replications of permuting the response, growing an RF and computing the variable importance. The authors recommend 50-100 permutations. However, much larger numbers have to be used to estimate more precise p-values. We add 1 to the numerator and denominator to avoid zero p-values.

Value

Variable importance and p-value for each variable.

Author(s)

Marvin N. Wright

References

Janitza, S., Celik, E. & Boulesteix, A.-L., (2016). A computationally fast variable importance test for random forests for high-dimensional data. *Adv Data Anal Classif* <https://doi.org/10.1007/s11634-016-0276-4>.

Altmann, A., Tolosi, L., Sander, O. & Lengauer, T. (2010). Permutation importance: a corrected feature importance measure, *Bioinformatics* 26:1340-1347.

See Also

[ranger](#)

Examples

```
## Janitza's p-values with corrected Gini importance
n <- 50
p <- 400
dat <- data.frame(y = factor(rbinom(n, 1, .5)), replicate(p, runif(n)))
rf.sim <- ranger(y ~ ., dat, importance = "impurity_corrected")
importance_pvalues(rf.sim, method = "janitza")

## Permutation p-values
## Not run:
rf.iris <- ranger(Species ~ ., data = iris, importance = 'permutation')
importance_pvalues(rf.iris, method = "altmann", formula = Species ~ ., data = iris)

## End(Not run)
```

parse.formula	<i>Parse formula</i>
---------------	----------------------

Description

Parse formula and return dataset containing selected columns. Interactions are supported for numerical columns only. An interaction column is the product of all interacting columns.

Usage

```
parse.formula(formula, data, env = parent.frame())
```

Arguments

formula	Object of class formula or character describing the model to fit.
data	Training data of class data.frame.
env	The environment in which the left hand side of formula is evaluated.

Value

Dataset including selected columns and interactions.

predict.ranger *Ranger prediction*

Description

Prediction with new data and a saved forest from Ranger.

Usage

```
## S3 method for class 'ranger'
predict(
  object,
  data = NULL,
  predict.all = FALSE,
  num.trees = object$num.trees,
  type = "response",
  se.method = "infjack",
  quantiles = c(0.1, 0.5, 0.9),
  what = NULL,
  seed = NULL,
  num.threads = NULL,
  verbose = TRUE,
  ...
)
```

Arguments

object	Ranger ranger object.
data	New test data of class <code>data.frame</code> or <code>gwaal.data</code> (GenABEL).
predict.all	Return individual predictions for each tree instead of aggregated predictions for all trees. Return a matrix (sample x tree) for classification and regression, a 3d array for probability estimation (sample x class x tree) and survival (sample x time x tree).
num.trees	Number of trees used for prediction. The first <code>num.trees</code> in the forest are used.
type	Type of prediction. One of 'response', 'se', 'terminalNodes', 'quantiles' with default 'response'. See below for details.
se.method	Method to compute standard errors. One of 'jack', 'infjack' with default 'infjack'. Only applicable if <code>type = 'se'</code> . See below for details.
quantiles	Vector of quantiles for quantile prediction. Set <code>type = 'quantiles'</code> to use.
what	User specified function for quantile prediction used instead of quantile. Must return numeric vector, see examples.
seed	Random seed. Default is NULL, which generates the seed from R. Set to 0 to ignore the R seed. The seed is used in case of ties in classification mode.
num.threads	Number of threads. Default is number of CPUs available.
verbose	Verbose output on or off.
...	further arguments passed to or from other methods.

Details

For type = 'response' (the default), the predicted classes (classification), predicted numeric values (regression), predicted probabilities (probability estimation) or survival probabilities (survival) are returned. For type = 'se', the standard error of the predictions are returned (regression only). The jackknife-after-bootstrap or infinitesimal jackknife for bagging is used to estimate the standard errors based on out-of-bag predictions. See Wager et al. (2014) for details. For type = 'terminalNodes', the IDs of the terminal node in each tree for each observation in the given dataset are returned. For type = 'quantiles', the selected quantiles for each observation are estimated. See Meinshausen (2006) for details.

If type = 'se' is selected, the method to estimate the variances can be chosen with `se.method`. Set `se.method = 'jack'` for jackknife-after-bootstrap and `se.method = 'infjack'` for the infinitesimal jackknife for bagging.

For classification and `predict.all = TRUE`, a factor levels are returned as numerics. To retrieve the corresponding factor levels, use `rf$forest$levels`, if `rf` is the ranger object.

Value

Object of class `ranger.prediction` with elements

<code>predictions</code>	Predicted classes/values (only for classification and regression)
<code>unique.death.times</code>	Unique death times (only for survival).
<code>chf</code>	Estimated cumulative hazard function for each sample (only for survival).
<code>survival</code>	Estimated survival function for each sample (only for survival).
<code>num.trees</code>	Number of trees.
<code>num.independent.variables</code>	Number of independent variables.
<code>treetype</code>	Type of forest/tree. Classification, regression or survival.
<code>num.samples</code>	Number of samples.

Author(s)

Marvin N. Wright

References

- Wright, M. N. & Ziegler, A. (2017). `ranger`: A Fast Implementation of Random Forests for High Dimensional Data in C++ and R. *J Stat Softw* 77:1-17. <https://doi.org/10.18637/jss.v077.i01>.
- Wager, S., Hastie T., & Efron, B. (2014). Confidence Intervals for Random Forests: The Jackknife and the Infinitesimal Jackknife. *J Mach Learn Res* 15:1625-1651. <http://jmlr.org/papers/v15/wager14a.html>.
- Meinshausen (2006). Quantile Regression Forests. *J Mach Learn Res* 7:983-999. <http://www.jmlr.org/papers/v7/meinshausen06a.html>.

See Also

[ranger](#)

Examples

```
## Classification forest
ranger(Species ~ ., data = iris)
train.idx <- sample(nrow(iris), 2/3 * nrow(iris))
iris.train <- iris[train.idx, ]
iris.test <- iris[-train.idx, ]
rg.iris <- ranger(Species ~ ., data = iris.train)
pred.iris <- predict(rg.iris, data = iris.test)
table(iris.test$Species, pred.iris$predictions)

## Quantile regression forest
rf <- ranger(mpg ~ ., mtcars[1:26, ], quantreg = TRUE)
pred <- predict(rf, mtcars[27:32, ], type = "quantiles", quantiles = c(0.1, 0.5, 0.9))
pred$predictions

## Quantile regression forest with user-specified function
rf <- ranger(mpg ~ ., mtcars[1:26, ], quantreg = TRUE)
pred <- predict(rf, mtcars[27:32, ], type = "quantiles",
               what = function(x) sample(x, 10, replace = TRUE))
pred$predictions
```

predict.ranger.forest *Ranger prediction*

Description

Prediction with new data and a saved forest from Ranger.

Usage

```
## S3 method for class 'ranger.forest'
predict(
  object,
  data,
  predict.all = FALSE,
  num.trees = object$num.trees,
  type = "response",
  se.method = "infjack",
  seed = NULL,
  num.threads = NULL,
  verbose = TRUE,
  inbag.counts = NULL,
  ...
)
```

Arguments

object	Ranger ranger.forest object.
data	New test data of class data.frame or gwaa.data (GenABEL).
predict.all	Return individual predictions for each tree instead of aggregated predictions for all trees. Return a matrix (sample x tree) for classification and regression, a 3d array for probability estimation (sample x class x tree) and survival (sample x time x tree).
num.trees	Number of trees used for prediction. The first num.trees in the forest are used.
type	Type of prediction. One of 'response', 'se', 'terminalNodes', 'quantiles' with default 'response'. See below for details.
se.method	Method to compute standard errors. One of 'jack', 'infjack' with default 'infjack'. Only applicable if type = 'se'. See below for details.
seed	Random seed. Default is NULL, which generates the seed from R. Set to 0 to ignore the R seed. The seed is used in case of ties in classification mode.
num.threads	Number of threads. Default is number of CPUs available.
verbose	Verbose output on or off.
inbag.counts	Number of times the observations are in-bag in the trees.
...	further arguments passed to or from other methods.

Details

For type = 'response' (the default), the predicted classes (classification), predicted numeric values (regression), predicted probabilities (probability estimation) or survival probabilities (survival) are returned. For type = 'se', the standard error of the predictions are returned (regression only). The jackknife-after-bootstrap or infinitesimal jackknife for bagging is used to estimate the standard errors based on out-of-bag predictions. See Wager et al. (2014) for details. For type = 'terminalNodes', the IDs of the terminal node in each tree for each observation in the given dataset are returned.

If type = 'se' is selected, the method to estimate the variances can be chosen with se.method. Set se.method = 'jack' for jackknife after bootstrap and se.method = 'infjack' for the infinitesimal jackknife for bagging.

For classification and predict.all = TRUE, a factor levels are returned as numerics. To retrieve the corresponding factor levels, use rf\$forest\$levels, if rf is the ranger object.

Value

Object of class ranger.prediction with elements

predictions	Predicted classes/values (only for classification and regression)
unique.death.times	Unique death times (only for survival).
chf	Estimated cumulative hazard function for each sample (only for survival).
survival	Estimated survival function for each sample (only for survival).
num.trees	Number of trees.
num.independent.variables	Number of independent variables.
treetype	Type of forest/tree. Classification, regression or survival.
num.samples	Number of samples.

Author(s)

Marvin N. Wright

References

- Wright, M. N. & Ziegler, A. (2017). ranger: A Fast Implementation of Random Forests for High Dimensional Data in C++ and R. J Stat Softw 77:1-17. <https://doi.org/10.18637/jss.v077.i01>.
- Wager, S., Hastie T., & Efron, B. (2014). Confidence Intervals for Random Forests: The Jackknife and the Infinitesimal Jackknife. J Mach Learn Res 15:1625-1651. <http://jmlr.org/papers/v15/wager14a.html>.

See Also

[ranger](#)

predictions.ranger *Ranger predictions*

Description

Extract training data predictions of Ranger object.

Usage

```
## S3 method for class 'ranger'  
predictions(x, ...)
```

Arguments

x Ranger object.
... Further arguments passed to or from other methods.

Value

Predictions: Classes for Classification forests, Numerical values for Regressions forests and the estimated survival functions for all individuals for Survival forests.

Author(s)

Marvin N. Wright

See Also

[ranger](#)

`predictions.ranger.prediction`
Ranger predictions

Description

Extract predictions of Ranger prediction object.

Usage

```
## S3 method for class 'ranger.prediction'
predictions(x, ...)
```

Arguments

`x` Ranger prediction object.
`...` Further arguments passed to or from other methods.

Value

Predictions: Classes for Classification forests, Numerical values for Regressions forests and the estimated survival functions for all individuals for Survival forests.

Author(s)

Marvin N. Wright

See Also

[ranger](#)

`print.ranger` *Print Ranger*

Description

Print contents of Ranger object.

Usage

```
## S3 method for class 'ranger'
print(x, ...)
```

Arguments

`x` Object of class 'ranger'.
`...` Further arguments passed to or from other methods.

Author(s)

Marvin N. Wright

See Also

[ranger](#)

`print.ranger.forest` *Print Ranger forest*

Description

Print contents of Ranger forest object.

Usage

```
## S3 method for class 'ranger.forest'  
print(x, ...)
```

Arguments

`x` Object of class 'ranger.forest'.
`...` further arguments passed to or from other methods.

Author(s)

Marvin N. Wright

`print.ranger.prediction`
 Print Ranger prediction

Description

Print contents of Ranger prediction object.

Usage

```
## S3 method for class 'ranger.prediction'  
print(x, ...)
```

Arguments

`x` Object of class 'ranger.prediction'.
`...` further arguments passed to or from other methods.

Author(s)

Marvin N. Wright

*ranger**Ranger*

Description

Ranger is a fast implementation of random forests (Breiman 2001) or recursive partitioning, particularly suited for high dimensional data. Classification, regression, and survival forests are supported. Classification and regression forests are implemented as in the original Random Forest (Breiman 2001), survival forests as in Random Survival Forests (Ishwaran et al. 2008). Includes implementations of extremely randomized trees (Geurts et al. 2006) and quantile regression forests (Meinshausen 2006).

Usage

```
ranger(  
  formula = NULL,  
  data = NULL,  
  num.trees = 500,  
  mtry = NULL,  
  importance = "none",  
  write.forest = TRUE,  
  probability = FALSE,  
  min.node.size = NULL,  
  max.depth = NULL,  
  replace = TRUE,  
  sample.fraction = ifelse(replace, 1, 0.632),  
  case.weights = NULL,  
  class.weights = NULL,  
  splitrule = NULL,  
  num.random.splits = 1,  
  alpha = 0.5,  
  minprop = 0.1,  
  split.select.weights = NULL,  
  always.split.variables = NULL,  
  respect.unordered.factors = NULL,  
  scale.permutation.importance = FALSE,  
  local.importance = FALSE,  
  regularization.factor = 1,  
  regularization.usedepth = FALSE,  
  keep.inbag = FALSE,  
  inbag = NULL,  
  holdout = FALSE,  
  quantreg = FALSE,
```

```

oob.error = TRUE,
num.threads = NULL,
save.memory = FALSE,
verbose = TRUE,
seed = NULL,
dependent.variable.name = NULL,
status.variable.name = NULL,
classification = NULL,
x = NULL,
y = NULL
)

```

Arguments

<code>formula</code>	Object of class <code>formula</code> or character describing the model to fit. Interaction terms supported only for numerical variables.
<code>data</code>	Training data of class <code>data.frame</code> , <code>matrix</code> , <code>dgCMatrix</code> (<code>Matrix</code>) or <code>gwaa.data</code> (<code>GenABEL</code>).
<code>num.trees</code>	Number of trees.
<code>mtry</code>	Number of variables to possibly split at in each node. Default is the (rounded down) square root of the number variables. Alternatively, a single argument function returning an integer, given the number of independent variables.
<code>importance</code>	Variable importance mode, one of 'none', 'impurity', 'impurity_corrected', 'permutation'. The 'impurity' measure is the Gini index for classification, the variance of the responses for regression and the sum of test statistics (see <code>splitrule</code>) for survival.
<code>write.forest</code>	Save <code>ranger.forest</code> object, required for prediction. Set to <code>FALSE</code> to reduce memory usage if no prediction intended.
<code>probability</code>	Grow a probability forest as in Malley et al. (2012).
<code>min.node.size</code>	Minimal node size. Default 1 for classification, 5 for regression, 3 for survival, and 10 for probability.
<code>max.depth</code>	Maximal tree depth. A value of <code>NULL</code> or 0 (the default) corresponds to unlimited depth, 1 to tree stumps (1 split per tree).
<code>replace</code>	Sample with replacement.
<code>sample.fraction</code>	Fraction of observations to sample. Default is 1 for sampling with replacement and 0.632 for sampling without replacement. For classification, this can be a vector of class-specific values.
<code>case.weights</code>	Weights for sampling of training observations. Observations with larger weights will be selected with higher probability in the bootstrap (or subsampled) samples for the trees.
<code>class.weights</code>	Weights for the outcome classes (in order of the factor levels) in the splitting rule (cost sensitive learning). Classification and probability prediction only. For classification the weights are also applied in the majority vote in terminal nodes.

<code>splitrule</code>	Splitting rule. For classification and probability estimation "gini", "extratrees" or "hellinger" with default "gini". For regression "variance", "extratrees", "maxstat" or "beta" with default "variance". For survival "logrank", "extratrees", "C" or "maxstat" with default "logrank".
<code>num.random.splits</code>	For "extratrees" <code>splitrule</code> .: Number of random splits to consider for each candidate splitting variable.
<code>alpha</code>	For "maxstat" <code>splitrule</code> : Significance threshold to allow splitting.
<code>minprop</code>	For "maxstat" <code>splitrule</code> : Lower quantile of covariate distribution to be considered for splitting.
<code>split.select.weights</code>	Numeric vector with weights between 0 and 1, representing the probability to select variables for splitting. Alternatively, a list of size <code>num.trees</code> , containing split select weight vectors for each tree can be used.
<code>always.split.variables</code>	Character vector with variable names to be always selected in addition to the <code>mtry</code> variables tried for splitting.
<code>respect.unordered.factors</code>	Handling of unordered factor covariates. One of 'ignore', 'order' and 'partition'. For the "extratrees" <code>splitrule</code> the default is "partition" for all other <code>splitrules</code> 'ignore'. Alternatively TRUE (= 'order') or FALSE (= 'ignore') can be used. See below for details.
<code>scale.permutation.importance</code>	Scale permutation importance by standard error as in (Breiman 2001). Only applicable if permutation variable importance mode selected.
<code>local.importance</code>	Calculate and return local importance values as in (Breiman 2001). Only applicable if <code>importance</code> is set to 'permutation'.
<code>regularization.factor</code>	Regularization factor (gain penalization), either a vector of length <code>p</code> or one value for all variables.
<code>regularization.usedepth</code>	Consider the depth in regularization.
<code>keep.inbag</code>	Save how often observations are in-bag in each tree.
<code>inbag</code>	Manually set observations per tree. List of size <code>num.trees</code> , containing <code>inbag</code> counts for each observation. Can be used for stratified sampling.
<code>holdout</code>	Hold-out mode. Hold-out all samples with case weight 0 and use these for variable importance and prediction error.
<code>quantreg</code>	Prepare quantile prediction as in quantile regression forests (Meinshausen 2006). Regression only. Set <code>keep.inbag = TRUE</code> to prepare out-of-bag quantile prediction.
<code>oob.error</code>	Compute OOB prediction error. Set to FALSE to save computation time, e.g. for large survival forests.
<code>num.threads</code>	Number of threads. Default is number of CPUs available.

<code>save.memory</code>	Use memory saving (but slower) splitting mode. No effect for survival and GWAS data. Warning: This option slows down the tree growing, use only if you encounter memory problems.
<code>verbose</code>	Show computation status and estimated runtime.
<code>seed</code>	Random seed. Default is NULL, which generates the seed from R. Set to 0 to ignore the R seed.
<code>dependent.variable.name</code>	Name of dependent variable, needed if no formula given. For survival forests this is the time variable.
<code>status.variable.name</code>	Name of status variable, only applicable to survival data and needed if no formula given. Use 1 for event and 0 for censoring.
<code>classification</code>	Set to TRUE to grow a classification forest. Only needed if the data is a matrix or the response numeric.
<code>x</code>	Predictor data (independent variables), alternative interface to data with formula or <code>dependent.variable.name</code> .
<code>y</code>	Response vector (dependent variable), alternative interface to data with formula or <code>dependent.variable.name</code> . For survival use a <code>Surv()</code> object or a matrix with time and status.

Details

The tree type is determined by the type of the dependent variable. For factors classification trees are grown, for numeric values regression trees and for survival objects survival trees. The Gini index is used as default splitting rule for classification. For regression, the estimated response variances or maximally selected rank statistics (Wright et al. 2016) can be used. For Survival the log-rank test, a C-index based splitting rule (Schmid et al. 2015) and maximally selected rank statistics (Wright et al. 2016) are available. For all tree types, forests of extremely randomized trees (Geurts et al. 2006) can be grown.

With the `probability` option and factor dependent variable a probability forest is grown. Here, the node impurity is used for splitting, as in classification forests. Predictions are class probabilities for each sample. In contrast to other implementations, each tree returns a probability estimate and these estimates are averaged for the forest probability estimate. For details see Malley et al. (2012).

Note that for classification and regression nodes with size smaller than `min.node.size` can occur, as in original Random Forests. For survival all nodes contain at `min.node.size` samples. Variables selected with `always.split.variables` are tried additionally to the `mtry` variables randomly selected. In `split.select.weights` variables weighted with 0 are never selected and variables with 1 are always selected. Weights do not need to sum up to 1, they will be normalized later. The weights are assigned to the variables in the order they appear in the formula or in the data if no formula is used. Names of the `split.select.weights` vector are ignored. The usage of `split.select.weights` can increase the computation times for large forests.

Unordered factor covariates can be handled in 3 different ways by using `respect.unordered.factors`: For 'ignore' all factors are regarded ordered, for 'partition' all possible 2-partitions are considered for splitting. For 'order' and 2-class classification the factor levels are ordered by their proportion falling in the second class, for regression by their mean response, as described in Hastie et al. (2009), chapter 9.2.4. For multiclass classification the factor levels are ordered by the first principal

component of the weighted covariance matrix of the contingency table (Coppersmith et al. 1999), for survival by the median survival (or the largest available quantile if the median is not available). The use of 'order' is recommended, as it computationally fast and can handle an unlimited number of factor levels. Note that the factors are only reordered once and not again in each split.

The 'impurity_corrected' importance measure is unbiased in terms of the number of categories and category frequencies and is almost as fast as the standard impurity importance. It is a modified version of the method by Sandri & Zuccolotto (2008), which is faster and more memory efficient. See Nembrini et al. (2018) for details. This importance measure can be combined with the methods to estimate p-values in `importance_pvalues`.

Regularization works by penalizing new variables by multiplying the splitting criterion by a factor, see Deng & Runger (2012) for details. If `regularization.usedepth=TRUE`, f^d is used, where f is the regularization factor and d the depth of the node. If regularization is used, multithreading is deactivated because all trees need access to the list of variables that are already included in the model.

For a large number of variables and data frames as input data the formula interface can be slow or impossible to use. Alternatively `dependent.variable.name` (and `status.variable.name` for survival) or `x` and `y` can be used. Use `x` and `y` with a matrix for `x` to avoid conversions and save memory. Consider setting `save.memory = TRUE` if you encounter memory problems for very large datasets, but be aware that this option slows down the tree growing.

For GWAS data consider combining ranger with the GenABEL package. See the Examples section below for a demonstration using PLink data. All SNPs in the GenABEL object will be used for splitting. To use only the SNPs without sex or other covariates from the phenotype file, use `0` on the right hand side of the formula. Note that missing values are treated as an extra category while splitting.

See <https://github.com/imbs-hl/ranger> for the development version.

With recent R versions, multithreading on Windows platforms should just work. If you compile yourself, the new RTools toolchain is required.

Value

Object of class `ranger` with elements

<code>forest</code>	Saved forest (If <code>write.forest</code> set to <code>TRUE</code>). Note that the variable IDs in the <code>split.varIDs</code> object do not necessarily represent the column number in R.
<code>predictions</code>	Predicted classes/values, based on out of bag samples (classification and regression only).
<code>variable.importance</code>	Variable importance for each independent variable.
<code>variable.importance.local</code>	Variable importance for each independent variable and each sample, if <code>local.importance</code> is set to <code>TRUE</code> and <code>importance</code> is set to 'permutation'.
<code>prediction.error</code>	Overall out of bag prediction error. For classification this is the fraction of misclassified samples, for probability estimation the Brier score, for regression the mean squared error and for survival one minus Harrell's C-index.

<code>r.squared</code>	R squared. Also called explained variance or coefficient of determination (regression only). Computed on out of bag data.
<code>confusion.matrix</code>	Contingency table for classes and predictions based on out of bag samples (classification only).
<code>unique.death.times</code>	Unique death times (survival only).
<code>chf</code>	Estimated cumulative hazard function for each sample (survival only).
<code>survival</code>	Estimated survival function for each sample (survival only).
<code>call</code>	Function call.
<code>num.trees</code>	Number of trees.
<code>num.independent.variables</code>	Number of independent variables.
<code>mtry</code>	Value of mtry used.
<code>min.node.size</code>	Value of minimal node size used.
<code>treetype</code>	Type of forest/tree. classification, regression or survival.
<code>importance.mode</code>	Importance mode used.
<code>num.samples</code>	Number of samples.
<code>inbag.counts</code>	Number of times the observations are in-bag in the trees.

Author(s)

Marvin N. Wright

References

- Wright, M. N. & Ziegler, A. (2017). ranger: A fast implementation of random forests for high dimensional data in C++ and R. *J Stat Softw* 77:1-17. <https://doi.org/10.18637/jss.v077.i01>.
- Schmid, M., Wright, M. N. & Ziegler, A. (2016). On the use of Harrell's C for clinical risk prediction via random survival forests. *Expert Syst Appl* 63:450-459. <https://doi.org/10.1016/j.eswa.2016.07.018>.
- Wright, M. N., Dankowski, T. & Ziegler, A. (2017). Unbiased split variable selection for random survival forests using maximally selected rank statistics. *Stat Med* 36:1272-1284. <https://doi.org/10.1002/sim.7212>.
- Nembrini, S., Koenig, I. R. & Wright, M. N. (2018). The revival of the Gini Importance? *Bioinformatics*. <https://doi.org/10.1093/bioinformatics/bty373>.
- Breiman, L. (2001). Random forests. *Mach Learn*, 45:5-32. <https://doi.org/10.1023/A:1010933404324>.
- Ishwaran, H., Kogalur, U. B., Blackstone, E. H., & Lauer, M. S. (2008). Random survival forests. *Ann Appl Stat* 2:841-860. <https://doi.org/10.1097/JTO.0b013e318233d835>.

- Malley, J. D., Kruppa, J., Dasgupta, A., Malley, K. G., & Ziegler, A. (2012). Probability machines: consistent probability estimation using nonparametric learning machines. *Methods Inf Med* 51:74-81. <https://doi.org/10.3414/ME00-01-0052>.
- Hastie, T., Tibshirani, R., Friedman, J. (2009). *The Elements of Statistical Learning*. Springer, New York. 2nd edition.
- Geurts, P., Ernst, D., Wehenkel, L. (2006). Extremely randomized trees. *Mach Learn* 63:3-42. <https://doi.org/10.1007/s10994-006-6226-1>.
- Meinshausen (2006). Quantile Regression Forests. *J Mach Learn Res* 7:983-999. <http://www.jmlr.org/papers/v7/meinshausen06a.html>.
- Sandri, M. & Zuccolotto, P. (2008). A bias correction algorithm for the Gini variable importance measure in classification trees. *J Comput Graph Stat*, 17:611-628. <https://doi.org/10.1198/106186008X344522>.
- Coppersmith D., Hong S. J., Hosking J. R. (1999). Partitioning nominal attributes in decision trees. *Data Min Knowl Discov* 3:197-217. <https://doi.org/10.1023/A:1009869804967>.
- Deng & Runger (2012). Feature selection via regularized trees. The 2012 International Joint Conference on Neural Networks (IJCNN), Brisbane, Australia. <https://doi.org/10.1109/IJCNN.2012.6252640>.

See Also

[predict.ranger](#)

Examples

```
## Classification forest with default settings
ranger(Species ~ ., data = iris)

## Prediction
train.idx <- sample(nrow(iris), 2/3 * nrow(iris))
iris.train <- iris[train.idx, ]
iris.test <- iris[-train.idx, ]
rg.iris <- ranger(Species ~ ., data = iris.train)
pred.iris <- predict(rg.iris, data = iris.test)
table(iris.test$Species, pred.iris$predictions)

## Quantile regression forest
rf <- ranger(mpg ~ ., mtcars[1:26, ], quantreg = TRUE)
pred <- predict(rf, mtcars[27:32, ], type = "quantiles")
pred$predictions

## Variable importance
rg.iris <- ranger(Species ~ ., data = iris, importance = "impurity")
rg.iris$variable.importance

## Survival forest
require(survival)
rg.veteran <- ranger(Surv(time, status) ~ ., data = veteran)
plot(rg.veteran$unique.death.times, rg.veteran$survival[1,])
```

```
## Alternative interfaces (same results)
ranger(dependent.variable.name = "Species", data = iris)
ranger(y = iris[, 5], x = iris[, -5])

## Not run:
## Use GenABEL interface to read Plink data into R and grow a classification forest
## The ped and map files are not included
library(GenABEL)
convert.snp.ped("data.ped", "data.map", "data.raw")
dat.gwaa <- load.gwaa.data("data.pheno", "data.raw")
phdata(dat.gwaa)$trait <- factor(phdata(dat.gwaa)$trait)
ranger(trait ~ ., data = dat.gwaa)

## End(Not run)
```

timepoints.ranger *Ranger timepoints*

Description

Extract unique death times of Ranger Survival forest

Usage

```
## S3 method for class 'ranger'
timepoints(x, ...)
```

Arguments

x Ranger Survival forest object.
... Further arguments passed to or from other methods.

Value

Unique death times

Author(s)

Marvin N. Wright

See Also

[ranger](#)

timepoints.ranger.prediction
Ranger timepoints

Description

Extract unique death times of Ranger Survival prediction object.

Usage

```
## S3 method for class 'ranger.prediction'  
timepoints(x, ...)
```

Arguments

x Ranger Survival prediction object.
... Further arguments passed to or from other methods.

Value

Unique death times

Author(s)

Marvin N. Wright

See Also

[ranger](#)

treeInfo *Tree information in human readable format*

Description

Extract tree information of a ranger object.

Usage

```
treeInfo(object, tree = 1)
```

Arguments

object ranger object.
tree Number of the tree of interest.

Details

Node and variable ID's are 0-indexed, i.e., node 0 is the root node. If the formula interface is used in the `ranger` call, the variable ID's are usually different to the original data used to grow the tree. Refer to the variable name instead to be sure.

Splitting at unordered factors (nominal variables) depends on the option `respect.unordered.factors` in the `ranger` call. For the "ignore" and "order" approaches, all values smaller or equal the `splitval` value go to the left and all values larger go to the right, as usual. However, with "order" the values correspond to the order in `object$forest$covariate.levels` instead of the original order (usually alphabetical). In the "partition" mode, the `splitval` values for unordered factor are comma separated lists of values, representing the factor levels (in the original order) going to the right.

Value

A `data.frame` with the columns

<code>nodeID</code>	The <code>nodeID</code> , 0-indexed.
<code>leftChild</code>	ID of the left child node, 0-indexed.
<code>rightChild</code>	ID of the right child node, 0-indexed.
<code>splitvarID</code>	ID of the splitting variable, 0-indexed. Caution, the variable order changes if the formula interface is used.
<code>splitvarName</code>	Name of the splitting variable.
<code>splitval</code>	The splitting value. For numeric or ordinal variables, all values smaller or equal go to the left, larger values to the right.
<code>terminal</code>	Logical, TRUE for terminal nodes.
<code>prediction</code>	One column with the predicted class (factor) for classification and the predicted numerical value for regression.

Author(s)

Marvin N. Wright

See Also

[ranger](#)

Examples

```
rf <- ranger(Species ~ ., data = iris)
treeInfo(rf, 1)
```


Index

csrf, [2](#)

getTerminalNodeIDs, [3](#)

holdoutRF, [4](#)

importance (importance.ranger), [5](#)
importance.ranger, [5](#)
importance_pvalues, [5](#), [19](#)

parse.formula, [7](#)
predict.ranger, [8](#), [21](#)
predict.ranger.forest, [10](#)
predictions
 (predictions.ranger.prediction),
 [13](#)
predictions.ranger, [12](#)
predictions.ranger.prediction, [13](#)
print.ranger, [13](#)
print.ranger.forest, [14](#)
print.ranger.prediction, [14](#)

ranger, [4–6](#), [9](#), [12–14](#), [15](#), [22–24](#)

timepoints (timepoints.ranger), [22](#)
timepoints.ranger, [22](#)
timepoints.ranger.prediction, [23](#)
treeInfo, [23](#)