

Package ‘optRF’

April 1, 2025

Type Package

Title Optimising Random Forest Stability by Determining the Optimal Number of Trees

Version 1.2.0

Maintainer Thomas Martin Lange <thomas.lange@uni-goettingen.de>

Description Calculating the stability of random forest with certain numbers of trees. The non-linear relationship between stability and numbers of trees is described using a logistic regression model and used to estimate the optimal number of trees.

BugReports <https://github.com/tmlange/optRF/issues>

URL <https://github.com/tmlange/optRF>

License GPL (>= 2)

Encoding UTF-8

LazyData true

RoxygenNote 7.3.2

Depends R (>= 3.6.0)

Imports minpack.lm (>= 1.2-2), ranger (>= 0.12.0), irr (>= 0.82), graphics, methods, stats

Suggests covr, knitr, rmarkdown, spelling, testthat

VignetteBuilder knitr

NeedsCompilation no

Author Thomas Martin Lange [cre, aut]
(<<https://orcid.org/0000-0003-4351-7950>>),
Felix Heinrich [ctb] (<<https://orcid.org/0000-0002-6093-8522>>)

Repository CRAN

Date/Publication 2025-04-01 11:20:02 UTC

Contents

estimate_numtrees	2
estimate_stability	3
measure_stability	4
opt_importance	5
opt_prediction	7
plot_stability	8
SNPdata	10
Index	11

estimate_numtrees	<i>Estimate the required number of trees</i>
-------------------	--

Description

Estimate the number of trees required to achieve certain stability of random forest

Usage

```
estimate_numtrees(
  optRF_object,
  measure = c("selection", "importance", "prediction"),
  for_stability = 0.95
)
```

Arguments

optRF_object	An optRF_object, either the result from the opt_importance or the opt_prediction function.
measure	A character string indicating which stability measure is to be analysed. One of "selection" (default, analyses selection stability), "prediction" (analyses prediction stability) or "importance" (analyses variable importance stability).
for_stability	Either a single stability value or a vector containing multiple stability values for which the number of trees should be estimated.

Value

A data frame summarising the estimated stability and run time in seconds for the given num.trees values.

Examples

```
## Not run:
data(SNPdata)
set.seed(123)
result_optpred = opt_prediction(y = SNPdata[,1], X=SNPdata[,-1]) # optimise random forest
estimate_numtrees(result_optpred, measure="prediction", for_stability=0.95)

## End(Not run)
```

estimate_stability *Estimate the stability of random forest*

Description

Estimate the stability of random forest with certain numbers of trees

Usage

```
estimate_stability(
  optRF_object,
  with_num.trees = c(1000, 5000, 10000, 50000, 1e+05)
)
```

Arguments

`optRF_object` An `optRF_object`, either the result from the [opt_importance](#) or the [opt_prediction](#) function.

`with_num.trees` Either a single `num.trees` value or a vector containing multiple `num.trees` values for which the stability should be estimated.

Value

A data frame summarising the estimated stability and run time in seconds for the given `num.trees` values.

Examples

```
## Not run:
data(SNPdata)
set.seed(123)
result_optpred = opt_prediction(y = SNPdata[,1], X=SNPdata[,-1]) # optimise random forest
estimate_stability(result_optpred, with_num.trees=c(1000, 5000, 10000, 50000, 100000))

## End(Not run)
```

measure_stability *Measure the stability of random forest*

Description

Measure the stability of random forest for a certain data set with a certain number of trees

Usage

```
measure_stability(
  y,
  X,
  num.trees = 500,
  method = c("prediction", "importance"),
  X_Test = NULL,
  alpha = NULL,
  select_for = c("high", "low", "zero"),
  importance = c("permutation", "impurity", "impurity_corrected"),
  number_repetitions = 10,
  verbose = TRUE,
  ...
)
```

Arguments

y	A vector containing the response variable in the training data set.
X	A data frame containing the explanatory variables in the training data set. The number of rows must be equal to the number of elements in y.
num.trees	Either a single value or a vector containing the numbers of trees for which the stability should be analysed (default = 500).
method	Either "prediction" (default) or "importance" specifying if random forest should be used for prediction or to estimate the variable importance.
X_Test	If method is "prediction", a data frame containing the explanatory variables of the test data set. If not entered, the out of bag data will be used.
alpha	If method is "prediction", the number of best individuals to be selected in the test data set (default = 0.15), if method is "importance", the number of most important variables to be selected (default = 0.05).
select_for	If method is "prediction", what should be selected? In random forest classification, this must be set to a vector containing the values of the desired classes. In random forest regression, this can be set as "high" (default) to select the individuals with the highest predicted value, "low" to select the individuals with the lowest predicted value, or "zero" to select the individuals which predicted value is closest to zero.
importance	If method is "importance", the variable importance mode, one of "permutation" (default), "impurity" or "impurity_corrected".

number_repetitions	Number of repetitions of random forest to estimate the stability. It needs to be at least 2. Default is 10.
verbose	Show computation status.
...	Any other argument from the ranger function.

Value

A data frame summarising the estimated stability for the given num.trees values.

Examples

```
## Not run:
data(SNPdata)
set.seed(123)
stability_result = measure_stability(y = SNPdata[,1], X=SNPdata[,-1], num.trees=500)
stability_result # Stability of random forest with 500 trees

## End(Not run)
```

opt_importance	<i>Optimise random forest for estimation of variable importance</i>
----------------	---

Description

Optimising random forest for estimating the importance of variables by calculating the variable importance stability with certain numbers of trees

Usage

```
opt_importance(
  y,
  X,
  number_repetitions = 10,
  alpha = 0.05,
  num.trees_values = c(250, 500, 750, 1000, 2000),
  importance = c("permutation", "impurity", "impurity_corrected"),
  visualisation = c("none", "importance", "selection"),
  recommendation = c("importance", "selection", "none"),
  rec_thresh = 1e-06,
  round_recommendation = c("thousand", "hundred", "ten", "none"),
  verbose = TRUE,
  ...
)
```

Arguments

y	A vector containing the response variable.
X	A data frame containing the explanatory variables. The number of rows must be equal to the number of elements in y.
number_repetitions	Number of repetitions of random forest to estimate the stability. It needs to be at least 2. Default is 10.
alpha	The amount of most important variables to be selected based on their estimated variable importance. If < 1 , alpha will be considered the relative amount of variables in the data set.
num.trees_values	A vector containing the numbers of trees to be analysed. If not specified, 250, 500, 750, 1000, and 2000 trees will be analysed.
importance	Variable importance mode, one of "permutation" (default), "impurity" or "impurity_corrected". The "impurity" measure is the Gini index for classification and the variance of the responses for regression.
visualisation	Can be set to "importance" to draw a plot of the variable importance stability or to "selection" to draw a plot of the selection stability for the numbers of trees to be analysed.
recommendation	If set to "importance" (default) or "selection", a recommendation will be given based on optimised variable importance or selection stability. If set to be "none", the function will analyse the stability of random forest with the inserted numbers of trees without giving a recommendation.
rec_thresh	If the number of trees leads to an increase of stability smaller or equal to the value specified, this number of trees will be recommended. Default is 1e-6.
round_recommendation	Setting to what number the recommended number of trees should be rounded to. Options: "none", "ten", "hundred", "thousand" (default).
verbose	Show computation status
...	Any other argument from the ranger function.

Value

An `opt_importance_object` containing the recommended number of trees, based on which measure the recommendation was given (importance or selection), a matrix summarising the estimated stability and computation time of a random forest with the recommended numbers of trees, a matrix containing the calculated stability and computation time for the analysed numbers of trees, and the parameters used to model the relationship between stability and numbers of trees.

Examples

```
## Not run:
data(SNPdata)
set.seed(123)
result_optimp = opt_importance(y = SNPdata[,1], X=SNPdata[,-1]) # optimise random forest
summary(result_optimp)
```

```
## End(Not run)
```

opt_prediction	<i>Optimise random forest for prediction</i>
----------------	--

Description

Optimising random forest predictions by calculating the prediction stability with certain numbers of trees

Usage

```
opt_prediction(
  y,
  X,
  X_Test = NULL,
  number_repetitions = 10,
  alpha = 0.15,
  num.trees_values = c(250, 500, 750, 1000, 2000),
  visualisation = c("none", "prediction", "selection"),
  select_for = c("high", "low", "zero"),
  recommendation = c("prediction", "selection", "none"),
  rec_thresh = 1e-06,
  round_recommendation = c("thousand", "hundred", "ten", "none"),
  verbose = TRUE,
  ...
)
```

Arguments

y	A vector containing the response variable in the training data set.
X	A data frame containing the explanatory variables in the training data set. The number of rows must be equal to the number of elements in y.
X_Test	A data frame containing the explanatory variables of the test data set. If not entered, the out of bag data will be used.
number_repetitions	Number of repetitions of random forest to estimate the stability. It needs to be at least 2. Default is 10.
alpha	The number of best individuals to be selected in the test data set based on their predicted response values. If < 1, alpha will be considered to be the relative amount of individuals in the test data set.
num.trees_values	A vector containing the numbers of trees to be analysed. If not specified, 250, 500, 750, 1000, and 2000 trees will be analysed.

visualisation	Can be set to "prediction" to draw a plot of the prediction stability or "selection" to draw a plot of the selection stability for the numbers of trees to be analysed.
select_for	What should be selected? In random forest classification, this must be set to a vector containing the values of the desired classes. In random forest regression, this can be set as "high" (default) to select the individuals with the highest predicted value, "low" to select the individuals with the lowest predicted value, or "zero" to select the individuals which predicted value is closest to zero.
recommendation	If set to "prediction" (default) or "selection", a recommendation will be given based on optimised prediction or selection stability. If set to be "none", the function will analyse the stability of random forest with the inserted numbers of trees without giving a recommendation.
rec_thresh	If the number of trees leads to an increase of stability smaller or equal to the value specified, this number of trees will be recommended. Default is 1e-6.
round_recommendation	Setting to what number the recommended number of trees should be rounded to. Options: "none", "ten", "hundred", "thousand" (default).
verbose	Show computation status
...	Any other argument from the ranger function.

Value

An `opt_prediction_object` containing the recommended number of trees, based on which measure the recommendation was given (prediction or selection), a matrix summarising the estimated stability and computation time of a random forest with the recommended numbers of trees, a matrix containing the calculated stability and computation time for the analysed numbers of trees, and the parameters used to model the relationship between stability and numbers of trees.

Examples

```
## Not run:
data(SNPdata)
set.seed(123)
result_optpred = opt_prediction(y = SNPdata[,1], X=SNPdata[,-1]) # optimise random forest
summary(result_optpred)

## End(Not run)
```

plot_stability

Plot random forest stability

Description

Plot the estimated stability of random forest against certain numbers of trees

Usage

```
plot_stability(
  optRF_object,
  measure = c("selection", "importance", "prediction"),
  from = 0,
  to = 1e+05,
  add_recommendation = TRUE,
  add = FALSE,
  ...
)
```

Arguments

optRF_object	An optRF_object, either the result from the opt_importance or the opt_prediction function.
measure	A character string indicating which stability measure is to be plotted. One of "selection" (default, visualises selection stability), "prediction" (visualises prediction stability) or "importance" (visualises variable importance stability).
from	Smallest num.trees value to be plotted.
to	Greatest num.trees value to be plotted.
add_recommendation	When set as TRUE, if a recommendation was stated within the opt_prediction or opt_importance function, the recommended num.trees value as well as the expected random forest stability will be highlighted in the graph
add	If FALSE, a new plot will be created, if TRUE, the graph will be added to an existing plot.
...	Any other arguments from the plot function.

Value

A plot showing the estimated stability of random forest for the given num.trees values.

Examples

```
## Not run:
data(SNPdata)
set.seed(123)
result_optpred = opt_prediction(y = SNPdata[,1], X=SNPdata[,-1]) # optimise random forest
plot_stability(result_optpred, measure = "prediction", add_recommendation = TRUE, add=FALSE)
plot_stability(result_optpred, measure = "selection", add_recommendation = FALSE, add=TRUE)

## End(Not run)
```

SNPdata

Simulated data of wheat yield and genomic markers

Description

Data set containing simulated data of wheat yield in g/m^2 of 250 wheat lines and 5,000 SNP markers being coded as 0 for homozygous form of the major allele and 2 for homozygous form of the minor allele.

Usage

```
data(SNPdata)
```

Format

An object of class "data.frame"

yield Simulated wheat yield in g/m^2

SNP_0001 to SNP_5000 Simulated values for 5,000 single nucleotide polymorphism (SNP) markers

References

This artificial data set was created for the optRF package.

Examples

```
data(SNPdata)
SNPdata[1:5,1:5]
```

Index

* datasets

SNPdata, [10](#)

estimate_numtrees, [2](#)

estimate_stability, [3](#)

measure_stability, [4](#)

opt_importance, [2](#), [3](#), [5](#), [9](#)

opt_prediction, [2](#), [3](#), [7](#), [9](#)

plot_stability, [8](#)

SNPdata, [10](#)