

# Package ‘pre’

May 7, 2018

**Title** Prediction Rule Ensembles

**Version** 0.5.0

**Author** Marjolein Fokkema [aut, cre],  
Benjamin Christoffersen [aut]

**Maintainer** Marjolein Fokkema <m.fokkema@fsw.leidenuniv.nl>

**Description** Derives prediction rule ensembles (PREs). Largely follows the procedure for deriving PREs as described in Friedman & Popescu (2008; <DOI:10.1214/07-AOAS148>), with adjustments and improvements. The main function `pre()` derives prediction rule ensembles consisting of rules and/or linear terms for continuous, binary, count, multinomial, and multivariate continuous responses. Function `gpe()` derives generalized prediction ensembles, consisting of rules, hinge and linear functions of the predictor variables.

**URL** <https://github.com/marjoleinF/pre>

**BugReports** <https://github.com/marjoleinF/pre/issues>

**Depends** R (>= 3.1.0)

**Imports** earth, Formula, glmnet, graphics, methods, partykit, rpart,  
stringr

**Suggests** akima, datasets, doParallel, foreach, glmertree, grid,  
mlbench, testthat

**License** GPL-2 | GPL-3

**Encoding** UTF-8

**LazyData** true

**RoxygenNote** 6.0.1

**NeedsCompilation** no

**Repository** CRAN

**Date/Publication** 2018-05-07 19:51:48 UTC

## R topics documented:

bsnullinteract	2
carrillo	3
coef.gpe	4
coef.pre	5
corplot	6
cvpre	7
gpe	8
gpe_cv.glmnet	10
gpe_rules_pre	10
gpe_sample	12
gpe_trees	12
importance	14
interact	16
maxdepth_sampler	18
pairplot	19
plot.pre	21
pre	22
predict.gpe	26
predict.pre	26
print.gpe	27
print.pre	28
rTerm	29
singleplot	30
<b>Index</b>	<b>32</b>

---

bsnullinteract	<i>Compute bootstrapped null interaction prediction rule ensembles</i>
----------------	--

---

### Description

bsnullinteract generates bootstrapped null interaction models, which can be used to derive a reference distribution of the test statistic calculated with [interact](#).

### Usage

```
bsnullinteract(object, nsamp = 10, parallel = FALSE,
  penalty.par.val = "lambda.1se", verbose = FALSE)
```

### Arguments

object	object of class <a href="#">pre</a> .
nsamp	numeric. Number of bootstrapped null interaction models to be derived.
parallel	logical. Should parallel foreach be used to generate initial ensemble? Must register parallel beforehand, such as doMC or others.

`penalty.par.val` character or numeric. Which value of the penalty parameter criterion should be used? The value yielding minimum cv error ("`lambda.min`") or penalty parameter yielding error within 1 standard error of minimum cv error ("`lambda.1se`")? Alternatively, a numeric value may be specified, corresponding to one of the values of `lambda` in the sequence used by `glmnet`, for which estimated cv error can be inspected by inspecting `object$glmnet.fit` and running `plot(object$glmnet.fit)`.

`verbose` logical. should progress be printed to the command line?

**Details**

Computationally intensive.

**Value**

A list of length `nsamp` with null interaction models, to be used as input for `interact`.

**See Also**

[pre](#), [interact](#)

**Examples**

```
set.seed(42)
airq.ens <- pre(Ozone ~ ., data=airquality[complete.cases(airquality),])
nullmods <- bsnullinteract(airq.ens)
interact(airq.ens, nullmods = nullmods, col = c("#7FBFF5", "#8CC876"))
```

---

carrillo

*Data on personality characteristics and depressive symptom severity*

---

**Description**

Dataset from a study by Carrillo et al. (2001), who assessed the extent to which the subscales of the NEO Personality Inventory (NEO-PI; Costa and McCrae 1985) could predict depressive symptomatology, as measured by the Beck Depression Inventory (BDI; Beck, Steer, and Carbin 1988). The NEO-PI assesses five major personality dimensions (Neuroticism, Extraversion, Openness to Experience, Agreeableness and Conscientiousness). Each of these dimensions consist of six specific subtraits (facets). The NEO-PI and BDI were administered to 112 Spanish respondents. Respondents' age in years and sex were also recorded and included in the dataset.

**Usage**

```
data(carrillo)
```

**Format**

A data frame with 112 observations and 26 variables

## Details

- neuroticism facet and total scores: n1, n2, n3, n4, n5, n6, ntot
- extraversion facet and total scores: e1, e2, e3, e4, e5, e6, etot
- openness to experience facet and total scores: open1, open2, open3, open4, open5, open6, opentot
- altruism total score: altot
- conscientiousness total score: contot
- depression symptom severity: bdi
- sex: sexo
- age in years: edad

## References

Beck, A.T., Steer, R.A. & Carbin, M.G. (1988). Psychometric properties of the Beck Depression Inventory: Twenty-five years of evaluation. *Clinical Psychology Review*, 8(1), 77-100.

Carrillo, J. M., Rojo, N., Sanchez-Bernardos, M. L., & Avia, M. D. (2001). Openness to experience and depression. *European Journal of Psychological Assessment*, 17(2), 130.

Costa, P.T. & McCrae, R.R. (1985). *The NEO Personality Inventory*. Psychological Assessment Resources, Odessa, FL.

## Examples

```
data("carrillo")
summary(carrillo)
```

---

coef.gpe

*Coefficients for a General Prediction Ensemble (gpe)*

---

## Description

coef function for [gpe](#)

## Usage

```
## S3 method for class 'gpe'
coef(object, penalty.par.val = "lambda.1se", ...)
```

**Arguments**

object            object of class [pre](#)  
 penalty.par.val    character. Penalty parameter criterion to be used for selecting final model: lambda giving minimum cv error ("lambda.min") or lambda giving cv error that is within 1 standard error of minimum cv error ("lambda.1se"). Alternatively, a numeric value may be specified, corresponding to one of the values of lambda in the sequence used by glmnet, for which estimated cv error can be inspected by running `object$glmnet.fit` and `plot(object$glmnet.fit)`.  
 ...                additional arguments to be passed to [coef.glmnet](#).

**See Also**

[coef.pre](#)

---

coef.pre	<i>Coefficients for the final prediction rule ensemble</i>
----------	--

---

**Description**

`coef.pre` returns coefficients for prediction rules and linear terms in the final ensemble

**Usage**

```
## S3 method for class 'pre'
coef(object, penalty.par.val = "lambda.1se", ...)
```

**Arguments**

object            object of class [pre](#)  
 penalty.par.val    character. Penalty parameter criterion to be used for selecting final model: lambda giving minimum cv error ("lambda.min") or lambda giving cv error that is within 1 standard error of minimum cv error ("lambda.1se"). Alternatively, a numeric value may be specified, corresponding to one of the values of lambda in the sequence used by glmnet, for which estimated cv error can be inspected by running `object$glmnet.fit` and `plot(object$glmnet.fit)`.  
 ...                additional arguments to be passed to [coef.glmnet](#).

**Details**

In some cases, duplicated variable names may appear in the model. For example, the first variable is a factor named 'V1' and there are also variables named 'V10' and/or 'V11' and/or 'V12' (etc). Then for for the binary factor V1, dummy contrast variables will be created, named 'V10', 'V11', 'V12' (etc). As should be clear from this example, this yields duplicated variable names, which may yield problems, for example in the calculation of predictions and importances, later on. This can be prevented by renaming factor variables with numbers in their name, prior to analysis.

**Value**

returns a dataframe with 3 columns: coefficient, rule (rule or variable name) and description (NA for linear terms, conditions for rules).

**See Also**

[pre](#), [plot.pre](#), [cvpre](#), [importance](#), [predict.pre](#), [interact](#), [print.pre](#)

**Examples**

```
set.seed(42)
airq.ens <- pre(Ozone ~ ., data = airquality[complete.cases(airquality),])
coefs <- coef(airq.ens)
```

---

corplot	<i>Plot correlations between baselearners in a prediction rule ensemble (ore)</i>
---------	---

---

**Description**

corplot plots correlations between baselearners in a prediction rule ensemble

**Usage**

```
corplot(object, penalty.par.val = "lambda.1se", colors = NULL,
  fig.plot = c(0, 0.85, 0, 1), fig.legend = c(0.8, 0.95, 0, 1),
  legend.breaks = seq(-1, 1, by = 0.1))
```

**Arguments**

object	object of class pre
penalty.par.val	character or numeric. Value of the penalty parameter $\lambda$ to be used for selecting the final ensemble. The ensemble with penalty parameter criterion yielding minimum cv error ("lambda.min") is taken, by default. Alternatively, the penalty parameter yielding error within 1 standard error of minimum cv error ("lambda.1se"), or a numeric value may be specified, corresponding to one of the values of lambda in the sequence used by glmnet, for which estimated cv error can be inspected by running <code>x\$glmnet.fit</code> and <code>plot(x\$glmnet.fit)</code> .
colors	vector of contiguous colors to be used for plotting. If <code>colors = NULL</code> (default), <code>colorRampPalette</code> is used to generate a sequence of 200 colors going from red to white to blue. A different set of plotting colors can be specified here, for example: <code>cm.colors(100)</code> , <code>colorspace::rainbow_hcl(100)</code> or <code>colorRampPalette(c("red", "yellow", "green"))(100)</code> .
fig.plot	plotting region to be used for correlation plot. See <code>fig</code> under <a href="#">par</a> .

`fig.legend` plotting region to be used for legend. See `fig` under `par`.

`legend.breaks` numeric vector of breakpoints to be depicted in the plot's legend. Should be a sequence from -1 to 1.

### See Also

See [rainbow\\_hcl](#) and [colorRampPalette](#).

### Examples

```
set.seed(42)
airq.ens <- pre(Ozone ~ ., data = airquality[complete.cases(airquality),])
corplot(airq.ens)
```

---

cvpre

*Full k-fold cross validation of a prediction rule ensemble (pre)*

---

### Description

cvpre performs k-fold cross validation on the dataset used to create the prediction rule ensemble, providing an estimate of predictive accuracy on future observations.

### Usage

```
cvpre(object, k = 10, verbose = FALSE, pclass = 0.5,
       penalty.par.val = "lambda.1se", parallel = FALSE)
```

### Arguments

`object` An object of class [pre](#).

`k` integer. The number of cross validation folds to be used.

`verbose` logical. Should progress of the cross validation be printed to the command line?

`pclass` numeric. Only used for classification. Cut-off value for the predicted probabilities that should be used to classify observations to the second class.

`penalty.par.val` numeric or character. Calculate cross-validated error for ensembles with penalty parameter criterion giving minimum cv error ("lambda.min") or giving cv error that is within 1 standard error of minimum cv error ("lambda.1se")? Alternatively, a numeric value may be specified, corresponding to one of the values of lambda in the sequence used by `glmnet`, for which estimated cv error can be inspected by running `object$glmnet.fit` and `plot(object$glmnet.fit)`.

`parallel` logical. Should parallel foreach be used? Must register parallel beforehand, such as `doMC` or others.

**Value**

A list with three objects: `$cvpreds` (a vector with cross-validated predicted  $y$  values), `$ss` (a vector indicating the cross-validation subsample each training observation was assigned to) and `$accuracy`. For continuous outputs, accuracy is a list with elements `$MSE` (mean squared error on test observations), `$MAE` (mean absolute error on test observations). For classification, accuracy is a list with elements `$SEL` (mean squared error on predicted probabilities), `$AEL` (mean absolute error on predicted probabilities), `$MCR` (average misclassification error rate) and `$table` (table with proportions of (in)correctly classified observations per class).

**See Also**

[pre](#), [plot.pre](#), [coef.pre](#), [importance](#), [predict.pre](#), [interact](#), [print.pre](#)

**Examples**

```
set.seed(42)
airq.ens <- pre(Ozone ~ ., data = airquality[complete.cases(airquality),])
airq.cv <- cvpre(airq.ens)
```

---

gpe

*Derive a General Prediction Ensemble (gpe)*


---

**Description**

Provides an interface for deriving sparse prediction ensembles where basis functions are selected through L1 penalization.

**Usage**

```
gpe(formula, data, base_learners = list(gpe_trees(), gpe_linear()),
    weights = rep(1, times = nrow(data)), sample_func = gpe_sample(),
    verbose = FALSE, penalized_trainer = gpe_cv.glmnet(), model = TRUE)
```

**Arguments**

formula	Symbolic description of the model to be fit of the form $y \sim x_1 + x_2 + \dots + x_n$ . If the output variable (left-hand side of the formula) is a factor, an ensemble for binary classification is created. Otherwise, an ensemble for prediction of a continuous variable is created.
data	data.frame containing the variables in the model.
base_learners	List of functions which has formal arguments <code>formula</code> , <code>data</code> , <code>weights</code> , <code>sample_func</code> , <code>verbose</code> , and <code>family</code> and returns a vector of characters with terms for the final formula passed to <code>cv.glmnet</code> . See <a href="#">gpe_linear</a> , <a href="#">gpe_trees</a> , and <a href="#">gpe_earth</a> .
weights	Case weights with length equal to number of rows in data.



sample_func	Function used to sample when learning with base learners. The function should have formal argument <code>n</code> and <code>weights</code> and return a vector of indices. See <a href="#">gpe_sample</a> .
verbose	TRUE if comments should be posted throughout the computations.
penalized_trainer	Function with formal arguments <code>x</code> , <code>y</code> , <code>weights</code> , <code>family</code> which returns a fit object. This can be changed to test other "penalized trainers" (like other function that perform an L1 penalty or L2 penalty and elastic net penalty). Not using <a href="#">cv.glmnet</a> may cause other function for gpe objects to fail. See <a href="#">gpe_cv.glmnet</a> .
model	TRUE if the data should added to the returned object.

## Details

Provides a more general framework for making a sparse prediction ensemble than [pre](#).

By default, a similar fit to [pre](#) is obtained. In addition, multivariate adaptive regression splines (Friedman, 1991) can be included with `gpe_earth`. See examples.

Other customs base learners can be implemented. See [gpe\\_trees](#), [gpe\\_linear](#) or [gpe\\_earth](#) for details of the setup. The sampling function given by `sample_func` can also be replaced by a custom sampling function. See [gpe\\_sample](#) for details of the setup.

## Value

An object of class `gpe`.

## References

Friedman, J. H., & Popescu, B. E. (2008). Predictive learning via rule ensembles. *The Annals of Applied Statistics*, 2(3), 916-954. Friedman, J. H. (1991). Multivariate adaptive regression splines. *The Annals of Statistics*, 19(1), 1-67.

## See Also

[pre](#), [gpe\\_trees](#), [gpe\\_linear](#), [gpe\\_earth](#), [gpe\\_sample](#), [gpe\\_cv.glmnet](#)

## Examples

```
## Not run:
## Obtain similar fit to \code{\link{pre}}:
gpe.rules <- gpe(Ozone ~ ., data = airquality[complete.cases(airquality)],
  base_learners = list(gpe_linear(), gpe_trees()))
gpe.rules

## Also include products of hinge functions using MARS:
gpe.hinge <- gpe(Ozone ~ ., data = airquality[complete.cases(airquality)],
  base_learners = list(gpe_linear(), gpe_trees(), gpe_earth()))

## End(Not run)
```

---

gpe_cv.glmnet	<i>Default penalized trainer for gpe</i>
---------------	--

---

**Description**

Default "penalizer function" generator [gpe](#) which uses [cv.glmnet](#).

**Usage**

```
gpe_cv.glmnet(...)
```

**Arguments**

... arguments to [cv.glmnet](#). `x`, `y`, `weights` and `family` will not be used.

**Value**

Returns a function with formal arguments `x`, `y`, `weights`, `family` and returns a fit object.

**See Also**

[gpe](#)

---

gpe_rules_pre	<i>Get rule learner for gpe which mimics behavior of pre</i>
---------------	--

---

**Description**

`gpe_rules_pre` generates a learner function which generates rules like `pre`, which can be supplied to the `gpe` `base_learner` argument

**Usage**

```
gpe_rules_pre(learnrate = 0.01, par.init = FALSE, mtry = Inf,  
maxdepth = 3L, ntrees = 500, tree.control = ctree_control(),  
use.grad = TRUE, removeduplicates = TRUE, removecomplements = TRUE,  
tree.unbiased = TRUE)
```

**Arguments**

<code>learnrate</code>	numeric value $> 0$ . Learning rate or boosting parameter.
<code>par.init</code>	logical. Should parallel foreach be used to generate initial ensemble? Only used when <code>learnrate == 0</code> . Note: Must register parallel beforehand, such as <code>doMC</code> or others. Furthermore, setting <code>par.init = TRUE</code> will likely increase computation time for smaller datasets.
<code>mtry</code>	positive integer. Number of randomly selected predictor variables for creating each split in each tree. Ignored when <code>tree.unbiased=FALSE</code> .
<code>maxdepth</code>	positive integer. Maximum number of conditions in a rule. If <code>length(maxdepth) == 1</code> , it specifies the maximum depth of of each tree grown. If <code>length(maxdepth) == ntrees</code> , it specifies the maximum depth of every consecutive tree grown.
<code>ntrees</code>	positive integer value. Number of trees to generate for the initial ensemble.
<code>tree.control</code>	list with control parameters to be passed to the tree fitting function, generated using <code>ctree_control</code> , <code>mob_control</code> (if <code>use.grad = FALSE</code> ), or <code>rpart.control</code> (if <code>tree.unbiased = FALSE</code> ).
<code>use.grad</code>	logical. Should gradient boosting with regression trees be employed when <code>learnrate &gt; 0</code> ? That is, use <code>ctree</code> as in Friedman (2001), but without the line search. If <code>FALSE</code> . By default set to <code>TRUE</code> , as this yields shorter computation times. If set to <code>FALSE</code> , <code>glmtree</code> with intercept only models in the nodes will be employed. This will yield longer computation times, but may increase accuracy. See details below for possible combinations with <code>family</code> , <code>use.grad</code> and <code>learnrate</code> .
<code>removeduplicates</code>	logical. Remove rules from the ensemble which are identical to an earlier rule?
<code>removecomplements</code>	logical. Remove rules from the ensemble which are identical to (1 - an earlier rule)?
<code>tree.unbiased</code>	logical. Should an unbiased tree generation algorithm be employed for rule generation? Defaults to <code>TRUE</code> , if set to <code>FALSE</code> , rules will be generated employing the CART algorithm (which suffers from biased variable selection) as implemented in <code>rpart</code> . See details below for possible combinations with <code>family</code> , <code>use.grad</code> and <code>learnrate</code> .

**Examples**

```
## Not run:
## Obtain same fits with pre and gpe
set.seed(42)
gpe.mod <- gpe(Ozone ~ ., data = airquality[complete.cases(airquality),],
              base_learners = list(gpe_rules_pre(), gpe_linear()))
set.seed(42)
pre.mod <- pre(Ozone ~ ., data = airquality[complete.cases(airquality),],)

## End(Not run)
```

---

gpe\_sample                      *Sampling Function Generator for gpe*

---

### Description

Provides a sample function for [gpe](#).

### Usage

```
gpe_sample(sampfrac = 0.5)
```

### Arguments

sampfrac                      Fraction of n to use for sampling. It is the  $\eta/N$  in Friedman & Popescu (2008).

### Value

Returns a function that takes an n argument for the number of observations and a weights argument for the case weights. The function returns a vector of indices.

### References

Friedman, J. H., & Popescu, B. E. (2008). Predictive learning via rule ensembles. *The Annals of Applied Statistics*, 2(3), 916-954.

### See Also

[gpe](#)

---

gpe\_trees                      *Learner Functions Generators for gpe*

---

### Description

Functions to get "learner" functions for [gpe](#).

### Usage

```
gpe_trees(..., remove_duplicates_complements = TRUE, mtry = Inf,
  ntrees = 500, maxdepth = 3L, learnrate = 0.01, parallel = FALSE,
  use_grad = TRUE, tree.control = ctree_control(mtry = mtry, maxdepth =
  maxdepth))
```

```
gpe_linear(..., winsfrac = 0.025, normalize = TRUE)
```

```
gpe_earth(..., degree = 3, nk = 8, normalize = TRUE, ntrain = 100,
  learnrate = 0.1, cor_thresh = 0.99)
```

**Arguments**

...	Currently not used.
remove_duplicates_complements	TRUE. Should rules with complementary or duplicate support be removed?
mtry	Number of input variables randomly sampled as candidates at each node for random forest like algorithms. The argument is passed to the tree methods in the partykit package.
ntrees	Number of trees to fit. Will not have an effect if tree.control is used.
maxdepth	Maximum depth of trees. Will not have an effect if tree.control is used.
learnrate	Learning rate for methods. Corresponds to the $\nu$ parameter in Friedman & Popescu (2008).
parallel	TRUE. Should basis functions be found in parallel?
use_grad	TRUE. Should binary outcomes use gradient boosting with regression trees when learnrate > 0? That is, use <code>ctree</code> instead of <code>glmtree</code> as in Friedman (2001) with a second order Taylor expansion instead of first order as in Chen and Guestrin (2016).
tree.control	<code>ctree_control</code> with options for the <code>ctree</code> function.
winsfrac	Quantile to winsorize linear terms. The value should be in [0, 0.5)
normalize	TRUE. Should value be scaled by .4 times the inverse standard deviation? If TRUE, gives linear terms the same influence as a typical rule.
degree	Maximum degree of interactions in <code>earth</code> model.
nk	Maximum number of basis functions in <code>earth</code> model.
ntrain	Number of models to fit.
cor_thresh	A threshold on the pairwise correlation for removal of basis functions. This is similar to <code>remove_duplicates_complements</code> . One of the basis functions in pairs where the correlation exceeds the threshold is excluded. NULL implies no exclusion. Setting a value closer to zero will decrease the time needed to fit the final model.

**Details**

`gpe_trees` provides learners for tree method. Either `ctree` or `glmtree` from the partykit package will be used.

`gpe_linear` provides linear terms for the gpe.

`gpe_earth` provides basis functions where each factor is a hinge function. The model is estimated with `earth`.

**Value**

A function that has formal arguments `formula`, `data`, `weights`, `sample_func`, `verbose`, `family`, ... The function returns a vector with character where each element is a term for the final formula in the call to `cv.glmnet`

## References

- Hothorn, T., & Zeileis, A. (2015). partykit: A modular toolkit for recursive partytioning in R. *Journal of Machine Learning Research*, 16, 3905-3909.
- Friedman, J. H. (1991). Multivariate adaptive regression splines. *The Annals Statistics*, 19(1), 1-67.
- Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *The Annals of Applied Statistics*, 29(5), 1189-1232.
- Friedman, J. H. (1993). Fast MARS. Dept. of Statistics Technical Report No. 110, Stanford University.
- Friedman, J. H., & Popescu, B. E. (2008). Predictive learning via rule ensembles. *The Annals of Applied Statistics*, 2(3), 916-954.
- Chen T., & Guestrin C. (2016). Xgboost: A scalable tree boosting system. *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2016.

## See Also

[gpe](#), [rTerm](#), [lTerm](#), [eTerm](#)

---

importance	<i>Calculate importances of baselearners (rules and linear terms) and input variables in a prediction rule ensemble (pre)</i>
------------	---

---

## Description

importance calculates importances for rules, linear terms and input variables in the prediction rule ensemble (pre), and creates a bar plot of variable importances.

## Usage

```
importance(object, standardize = FALSE, global = TRUE,
  quantprobs = c(0.75, 1), penalty.par.val = "lambda.1se", round = NA,
  plot = TRUE, ylab = "Importance", main = "Variable importances",
  diag.xlab = TRUE, diag.xlab.hor = 0, diag.xlab.vert = 2, cex.axis = 1,
  ...)
```

## Arguments

- |             |   |
|-------------|---|
| object      | an object of class <a href="#">pre</a>  |
| standardize | logical. Should baselearner importances be standardized with respect to the outcome variable? If TRUE, baselearner importances have a minimum of 0 and a maximum of 1. Only used for ensembles with numeric (non-count) response variables. |
| global      | logical. Should global importances be calculated? If FALSE, local importances will be calculated, given the quantiles of the predictions $F(x)$ in quantprobs.  |

quantprobs	optional numeric vector of length two. Only used when <code>global = FALSE</code> . Probabilities for calculating sample quantiles of the range of $F(X)$ , over which local importances are calculated. The default provides variable importances calculated over the 25% highest values of $F(X)$ .
penalty.par.val	character or numeric. Should model be selected with lambda yielding minimum cv error ("lambda.min"), or lambda giving cv error that is within 1 standard error of minimum cv error ("lambda.1se")? Alternatively, a numeric value may be specified, corresponding to one of the values of lambda in the sequence used by <code>glmnet</code> .
round	integer. Number of decimal places to round numeric results to. If NA (default), no rounding is performed.
plot	logical. Should variable importances be plotted?
ylab	character string. Plotting label for y-axis. Only used when <code>plot = TRUE</code> .
main	character string. Main title of the plot. Only used when <code>plot = TRUE</code> .
diag.xlab	logical. Should variable names be printed diagonally (that is, in a 45 degree angle)? Alternatively, variable names may be printed vertically by specifying <code>diag.xlab = FALSE, las = 2</code> .
diag.xlab.hor	numeric. Horizontal adjustment for lining up variable names with bars in the plot if variable names are printed diagonally.
diag.xlab.vert	positive integer. Vertical adjustment for position of variable names, if printed diagonally. Corresponds to the number of character spaces added after variable names.
cex.axis	numeric. The magnification to be used for axis annotation relative to the current setting of <code>cex</code> .
...	further arguments to be passed to <code>barplot</code> (only used when <code>plot = TRUE</code> ).

## Value

A list with two dataframes: `$baseimps`, giving the importances for baselearners in the ensemble, and `$varimps`, giving the importances for all predictor variables.

## Examples

```
set.seed(42)
airq.ens <- pre(Ozone ~ ., data = airquality[complete.cases(airquality),])
# calculate global importances:
importance(airq.ens)
# calculate local importances (default: over 25% highest predicted values):
importance(airq.ens, global = FALSE)
# calculate local importances (custom: over 25% lowest predicted values):
importance(airq.ens, global = FALSE, quantprobs = c(0, .25))
```

---

interact	<i>Calculate interaction statistics for variables in a prediction rule ensemble (pre)</i>
----------	---

---

## Description

interact calculates test statistics for assessing the strength of interactions between a set of user-specified input variable(s), and all other input variables.

## Usage

```
interact(object, varnames = NULL, nullmods = NULL,
  penalty.par.val = "lambda.1se", quantprobs = c(0.05, 0.95), plot = TRUE,
  col = c("#8CC876", "#7FBFF5"), ylab = "Interaction strength",
  main = "Interaction test statistics", se.linewidth = 0.05,
  parallel = FALSE, k = 10, verbose = FALSE, ...)
```

## Arguments

object	an object of class <code>pre</code> .
varnames	character vector. Names of variables for which interaction statistics should be calculated. If <code>NULL</code> , interaction statistics for all predictor variables with non-zero coefficients will be calculated (which may take a long time).
nullmods	object with bootstrapped null interaction models, resulting from application of <code>bsnullinteract</code> .
penalty.par.val	character. Which value of the penalty parameter criterion should be used? The value yielding minimum cv error ("lambda.min") or penalty parameter yielding error within 1 standard error of minimum cv error ("lambda.1se")? Alternatively, a numeric value may be specified, corresponding to one of the values of lambda in the sequence used by <code>glmnet</code> , for which estimated cv error can be inspected by running <code>object\$glmnet.fit</code> and <code>plot(object\$glmnet.fit)</code> .
quantprobs	numeric vector of length two. Probabilities that should be used for plotting the range of bootstrapped null interaction model statistics. Only used when <code>nullmods</code> argument is specified and <code>plot = TRUE</code> . The default yields sample quantiles corresponding to .05 and .95 probabilities.
plot	logical. Should interaction statistics be plotted?
col	character vector of length one or two. The first value specifies the color to be used for plotting the interaction statistic from the training data, the second color is used for plotting the interaction statistic from the bootstrapped null interaction models. Only used when <code>plot = TRUE</code> and Only the first element is used if <code>nullmods = NULL</code> .
ylab	character string. Label to be used for plotting y-axis.
main	character. Main title for the bar plot.



<code>se.linewidth</code>	numeric. Width of the whiskers of the plotted standard error bars (in inches).
<code>parallel</code>	logical. Should parallel foreach be used? Must register parallel beforehand, such as <code>doMC</code> or others.
<code>k</code>	integer. Calculating interaction test statistics is a computationally intensive, so calculations are split up in several parts to prevent memory allocation errors. If a memory allocation error still occurs, increase <code>k</code> .
<code>verbose</code>	logical. Should progress information be printed to the command line?
<code>...</code>	Additional arguments to be passed to <code>barplot</code> .

### Details

Can be computationally intensive, especially when `nullmods` is specified, in which case setting `parallel = TRUE` may improve speed.

### Value

Function `interact()` returns and plots interaction statistics for the specified predictor variables. If `nullmods` is not specified, it returns and plots only the interaction test statistics for the specified fitted prediction rule ensemble. If `nullmods` is specified, the function returns a list, with elements `$fittedH2`, containing the interaction statistics of the fitted ensemble, and `$nullH2`, which contains the interaction test statistics for each of the bootstrapped null interaction models.

If `plot = TRUE` (the default), a barplot is created with the interaction test statistic from the fitted prediction rule ensemble. If `nullmods` is specified, bars representing the median of the distribution of interaction test statistics of the bootstrapped null interaction models are plotted. In addition, error bars representing the quantiles of the distribution (their value specified by the `quantprobs` argument) are plotted. These allow for testing the null hypothesis of no interaction effect for each of the input variables.

Note that the error rates of null hypothesis tests of interaction effects have not yet been studied in detail, but likely depend on the number of generated bootstrapped null interaction models as well as the complexity of the fitted ensembles. Users are therefore advised to test for the presence of interaction effects by setting the `nsamp` argument of the function `bsnullinteract`  $\geq 100$ .

### See Also

[pre](#), [bsnullinteract](#)

### Examples

```
set.seed(42)
airq.ens <- pre(Ozone ~ ., data=airquality[complete.cases(airquality),])
interact(airq.ens, c("Temp", "Wind", "Solar.R"))
```

---

maxdepth_sampler	<i>Sampling function generator for specifying varying maximum tree depth in a prediction rule ensemble (pre)</i>
------------------	--

---

### Description

maxdepth\_sampler generates a random sampling function, governed by a pre-specified average tree depth.

### Usage

```
maxdepth_sampler(av.no.term.nodes = 4L, av.tree.depth = NULL)
```

### Arguments

av.no.term.nodes  
integer of length one. Specifies the average number of terminal nodes in trees used for rule induction.

av.tree.depth  
integer of length one. Specifies the average maximum tree depth in trees used for rule induction.

### Details

The original RuleFit implementation varying tree sizes for rule induction. Furthermore, it defined tree size in terms of the number of terminal nodes. In contrast, function [pre](#) defines the maximum tree size in terms of a (constant) tree depth. Function maxdepth\_sampler allows for mimicing the behavior of the original RuleFit implementation. In effect, the maximum tree depth is sampled from an exponential distribution with learning rate  $\frac{1}{\bar{L}-2}$ , where  $(\bar{L}) \geq 2$  represents the average number of terminal nodes for trees in the ensemble. See Friedman & Popescu (2008, section 3.3).

### Value

Returns a random sampling function with single argument 'ntrees', which can be supplied to the maxdepth argument of function [pre](#) to specify varying tree depths.

### References

Friedman, J. H., & Popescu, B. E. (2008). Predictive learning via rule ensembles. *The Annals of Applied Statistics*, 2(3), 916-954.

### See Also

[pre](#)

**Examples**

```

## RuleFit default is max. 4 terminal nodes, on average:
func1 <- maxdepth_sampler()
set.seed(42)
func1(10)
mean(func1(1000))

## Max. 16 terminal nodes, on average (equals average maxdepth of 4):
func2 <- maxdepth_sampler(av.no.term.nodes = 16L)
set.seed(42)
func2(10)
mean(func2(1000))

## Max. tree depth of 3, on average:
func3 <- maxdepth_sampler(av.tree.depth = 3)
set.seed(42)
func3(10)
mean(func3(1000))

## Max. 2 of terminal nodes, on average (always yields maxdepth of 1):
func4 <- maxdepth_sampler(av.no.term.nodes = 2L)
set.seed(42)
func4(10)
mean(func4(1000))

## Not run:
## Create rule ensemble with varying maxdepth:
set.seed(42)
airq.ens <- pre(Ozone ~ ., data = airquality[complete.cases(airquality)],
               maxdepth = func1)

airq.ens

## End(Not run)

```

---

pairplot

*Create partial dependence plot for a pair of predictor variables in a prediction rule ensemble (pre)*

---

**Description**

pairplot creates a partial dependence plot to assess the effects of a pair of predictor variables on the predictions of the ensemble

**Usage**

```

pairplot(object, varnames, type = "both", penalty.par.val = "lambda.1se",
         nvals = c(20, 20), pred.type = "response", ...)

```

**Arguments**

<code>object</code>	an object of class <code>pre</code>
<code>varnames</code>	character vector of length two. Currently, pairplots can only be requested for non-nominal variables. If <code>varnames</code> specifies the name(s) of variables of class "factor", an error will be printed.
<code>type</code>	character string. Type of plot to be generated. <code>type = "heatmap"</code> yields a heatmap plot, <code>type = "contour"</code> yields a contour plot, <code>type = "both"</code> yields a heatmap plot with added contours, <code>type = "perspective"</code> yields a three dimensional plot.
<code>penalty.par.val</code>	character. Should model be selected with lambda giving minimum cv error ("lambda.min"), or lambda giving cv error that is within 1 standard error of minimum cv error ("lambda.1se")? Alternatively, a numeric value may be specified, corresponding to one of the values of lambda in the sequence used by <code>glmnet</code> , for which estimated cv error can be inspected by running <code>object\$glmnet.fit</code> and <code>plot(object\$glmnet.fit)</code> .
<code>nvals</code>	optional numeric vector of length 2. For how many values of <code>x1</code> and <code>x2</code> should partial dependence be plotted? If <code>NULL</code> , all observed values for the two predictor variables specified will be used (see details).
<code>pred.type</code>	character string. Type of prediction to be plotted on z-axis. <code>pred.type = "response"</code> gives fitted values for continuous outputs and fitted probabilities for nominal outputs. <code>pred.type = "link"</code> gives fitted values for continuous outputs and linear predictor values for nominal outputs.
<code>...</code>	Additional arguments to be passed to <code>image</code> , <code>contour</code> or <code>persp</code> (depending on whether <code>type</code> is specified to be "heatmap", "contour", "both" or "perspective").

**Details**

By default, partial dependence will be plotted for each combination of 20 values of the specified predictor variables. When `nvals = NULL` is specified a dependence plot will be created for every combination of the unique observed values of the two predictor variables specified. Therefore, using `nvals = NULL` will often result in long computation times, and / or memory allocation errors. Also, `pre` ensembles derived from training datasets that are very wide or long may result in long computation times and / or memory allocation errors. In such cases, reducing the values supplied to `nvals` will reduce computation time and / or memory allocation errors. When the `nvals` argument is supplied, values for the minimum, maximum, and `nvals - 2` intermediate values of the predictor variable will be plotted. Furthermore, if none of the variables specified appears in the final prediction rule ensemble, an error will occur.

**Note**

Function `pairplot` uses package `akima` to construct interpolated surfaces and has an ACM license that restricts applications to non-commercial usage, see <https://www.acm.org/publications/policies/software-copyright-notice> Function `pairplot` prints a note referring to this ACM licence.

**See Also**[pre](#), [singleplot](#)**Examples**

```
set.seed(42)
airq.ens <- pre(Ozone ~ ., data = airquality[complete.cases(airquality),])
pairplot(airq.ens, c("Temp", "Wind"))
```

plot.pre

*Plot method for class pre***Description**

plot.pre creates one or more plots depicting the rules in the final ensemble as simple decision trees.

**Usage**

```
## S3 method for class 'pre'
plot(x, penalty.par.val = "lambda.1se", linear.terms = TRUE,
     nterms = NULL, ask = FALSE, exit.label = "0", standardize = FALSE,
     plot.dim = c(3, 3), ...)
```

**Arguments**

x	an object of class <a href="#">pre</a> .
penalty.par.val	character. Which value of the penalty parameter criterion should be used? The value yielding minimum cv error ("lambda.min") or penalty parameter yielding error within 1 standard error of minimum cv error ("lambda.1se")? Alternatively, a numeric value may be specified, corresponding to one of the values of lambda in the sequence used by glmnet, for which estimated cv error can be inspected by running <code>x\$glmnet.fit</code> and <code>plot(x\$glmnet.fit)</code> .
linear.terms	logical. Should linear terms be included in the plot?
nterms	numeric. The total number of terms (or rules, if <code>linear.terms = FALSE</code> ) being plotted. Default is NULL, resulting in all terms of the final ensemble to be plotted.
ask	logical. Should user be prompted before starting a new page of plots?
exit.label	character string. Label to be printed in nodes to which the rule does not apply ("exit nodes")?
standardize	logical. Should printed importances be standardized? See <a href="#">importance</a> .
plot.dim	integer vector of length two. Specifies the number of rows and columns in the plot. The default yields a plot with three rows and three columns, depicting nine baselearners per plotting page.
...	Arguments to be passed to <a href="#">gpar</a> .

**See Also**

[pre](#), [print.pre](#)

**Examples**

```
set.seed(42)
airq.ens <- pre(Ozone ~ ., data = airquality[complete.cases(airquality),])
plot(airq.ens)
```

---

```
pre
```

---

*Derive a prediction rule ensemble*

---

**Description**

pre derives a sparse ensemble of rules and/or linear functions for prediction of a continuous or binary outcome.

**Usage**

```
pre(formula, data, family = gaussian, use.grad = TRUE, weights,
    type = "both", sampfrac = 0.5, maxdepth = 3L, learnrate = 0.01,
    mtry = Inf, ntrees = 500, removecomplements = TRUE,
    removeduplicates = TRUE, winsfrac = 0.025, normalize = TRUE,
    standardize = FALSE, nfolds = 10L, tree.control, tree.unbiased = TRUE,
    verbose = FALSE, par.init = FALSE, par.final = FALSE, ...)
```

**Arguments**

formula	a symbolic description of the model to be fit of the form $y \sim x_1 + x_2 + \dots + x_n$ . Response (left-hand side of the formula) should be of class numeric (for continuous outcomes), integer (for count outcomes) or a factor. In addition, a multivariate continuous response may be specified as follows: $y_1 + y_2 + y_3 \sim x_1 + x_2 + x_3$ . If the response is a factor, an ensemble for classification will be derived. Otherwise, an ensemble for prediction of a numeric response is created. If the outcome is a non-negative count, this should be specified by setting <code>family = "poisson"</code> . Note that input variables may not have 'rule' as (part of) their name, and the formula may not exclude the intercept (that is, $+ 0$ or $- 1$ may not be used in the right-hand side of the formula).
data	data.frame containing the variables in the model. Response must be a factor for binary classification, numeric for (count) regression. Input variables must be of class numeric, factor or ordered factor.
family	specification of a glm family. Can be a character string (i.e., "gaussian", "binomial", "poisson", "multinomial", or "mgaussian") or a corresponding family object (e.g., gaussian, binomial or poisson, see <a href="#">family</a> ). Specification is required only for non-negative count responses, e.g., <code>family = "poisson"</code> .

Otherwise, the program will try to make an informed guess: `family = "gaussian"` will be employed a numeric, `family = "binomial"` will be employed if a binary factor. `family = "multinomial"` will be employed if a factor with  $> 2$  levels, and `family = "mgaussian"` will be employed if multiple continuous response variables were specified.

<code>use.grad</code>	logical. Should gradient boosting with regression trees be employed when <code>learnrate &gt; 0</code> ? That is, use <code>ctree</code> as in Friedman (2001), but without the line search. If FALSE. By default set to TRUE, as this yields shorter computation times. If set to FALSE, <code>glmtree</code> with intercept only models in the nodes will be employed. This will yield longer computation times, but may increase accuracy. See details below for possible combinations with <code>family</code> , <code>use.grad</code> and <code>learnrate</code> .
<code>weights</code>	an optional vector of observation weights to be used for deriving the ensemble.
<code>type</code>	character. Specifies type of base learners to be included in the ensemble. Defaults to "both" (initial ensemble will include both rules and linear functions). Other options are "rules" (prediction rules only) or "linear" (linear functions only).
<code>sampfrac</code>	numeric value $> 0$ and $\leq 1$ . Specifies the fraction of randomly selected training observations used to produce each tree. Values $< 1$ will result in sampling without replacement (i.e., subsampling), a value of 1 will result in sampling with replacement (i.e., bootstrap sampling). Alternatively, a sampling function may be supplied, which should take arguments <code>n</code> (sample size) and <code>weights</code> .
<code>maxdepth</code>	positive integer. Maximum number of conditions in a rule. If <code>length(maxdepth) == 1</code> , it specifies the maximum depth of each tree grown. If <code>length(maxdepth) == ntrees</code> , it specifies the maximum depth of every consecutive tree grown. Alternatively, a random sampling function may be supplied, which takes argument <code>ntrees</code> and returns integer values. See also <code>maxdepth_sampler</code> .
<code>learnrate</code>	numeric value $> 0$ . Learning rate or boosting parameter.
<code>mtry</code>	positive integer. Number of randomly selected predictor variables for creating each split in each tree. Ignored when <code>tree.unbiased=FALSE</code> .
<code>ntrees</code>	positive integer value. Number of trees to generate for the initial ensemble.
<code>removecomplements</code>	logical. Remove rules from the ensemble which are identical to (1 - an earlier rule)?
<code>removeduplicates</code>	logical. Remove rules from the ensemble which are identical to an earlier rule?
<code>winsfrac</code>	numeric value $> 0$ and $\leq 0.5$ . Quantiles of data distribution to be used for winsorizing linear terms. If set to 0, no winsorizing is performed. Note that ordinal variables are included as linear terms in estimating the regression model and will also be winsorized.
<code>normalize</code>	logical. Normalize linear variables before estimating the regression model? Normalizing gives linear terms the same a priori influence as a typical rule, by dividing the (winsorized) linear term by 2.5 times its SD.
<code>standardize</code>	logical. Should rules and linear terms be standardized to have SD equal to 1 before estimating the regression model? This will also standardize the dummified factors, users are advised to use the default <code>standardize = FALSE</code> .

<code>nfolds</code>	positive integer. Number of cross-validation folds to be used for selecting the optimal value of the penalty parameter $\lambda$ in selecting the final ensemble.
<code>tree.control</code>	list with control parameters to be passed to the tree fitting function, generated using <code>ctree_control</code> , <code>mob_control</code> (if <code>use.grad = FALSE</code> ), or <code>rpart.control</code> (if <code>tree.unbiased = FALSE</code> ).
<code>tree.unbiased</code>	logical. Should an unbiased tree generation algorithm be employed for rule generation? Defaults to TRUE, if set to FALSE, rules will be generated employing the CART algorithm (which suffers from biased variable selection) as implemented in <code>rpart</code> . See details below for possible combinations with <code>family</code> , <code>use.grad</code> and <code>learnrate</code> .
<code>verbose</code>	logical. Should information on the initial and final ensemble be printed to the command line?
<code>par.init</code>	logical. Should parallel foreach be used to generate initial ensemble? Only used when <code>learnrate == 0</code> . Note: Must register parallel beforehand, such as <code>doMC</code> or others. Furthermore, setting <code>par.init = TRUE</code> will likely increase computation time for smaller datasets.
<code>par.final</code>	logical. Should parallel foreach be used to perform cross validation for selecting the final ensemble? Must register parallel beforehand, such as <code>doMC</code> or others.
<code>...</code>	Additional arguments to be passed to <code>cv.glmnet</code> .

## Details

Observations with missing values will be removed prior to analysis.

In some cases, duplicated variable names may appear in the model. For example, the first variable is a factor named 'V1' and there are also variables named 'V10' and/or 'V11' and/or 'V12' (etc). Then for the binary factor V1, dummy contrast variables will be created, named 'V10', 'V11', 'V12' (etc). As should be clear from this example, this yields duplicated variable names, which may yield problems, for example in the calculation of predictions and importances, later on. This can be prevented by renaming factor variables with numbers in their name, prior to analysis.

The table below provides an overview of combinations of response variable types, `use.grad`, `tree.unbiased` and `learnrate` settings that are supported, and the tree induction algorithm that will be employed as a result:

<code>use.grad</code>	<code>tree.unbiased</code>	<code>learnrate</code>	<code>family</code>	<code>tree alg.</code>	Response variable format
TRUE	TRUE	0	gaussian	ctree	Single, numeric (non-integer)
TRUE	TRUE	0	mgaussian	ctree	Multiple, numeric (non-integer)
TRUE	TRUE	0	binomial	ctree	Single, factor with 2 levels
TRUE	TRUE	0	multinomial	ctree	Single, factor with >2 levels
TRUE	TRUE	0	poisson	ctree	Single, integer
TRUE	TRUE	>0	gaussian	ctree	Single, numeric (non-integer)
TRUE	TRUE	>0	mgaussian	ctree	Multiple, numeric (non-integer)
TRUE	TRUE	>0	binomial	ctree	Single, factor with 2 levels
TRUE	TRUE	>0	multinomial	ctree	Single, factor with >2 levels
TRUE	TRUE	>0	poisson	ctree	Single, integer



FALSE	TRUE	0	gaussian	glmtree	Single, numeric (non-integer)
FALSE	TRUE	0	binomial	glmtree	Single, factor with 2 levels
FALSE	TRUE	0	poisson	glmtree	Single, integer
FALSE	TRUE	>0	gaussian	glmtree	Single, numeric (non-integer)
FALSE	TRUE	>0	binomial	glmtree	Single, factor with 2 levels
FALSE	TRUE	>0	poisson	glmtree	Single, integer
TRUE	FALSE	0	gaussian	rpart	Single, numeric (non-integer)
TRUE	FALSE	0	binomial	rpart	Single, factor with 2 levels
TRUE	FALSE	0	multinomial	rpart	Single, factor with >2 levels
TRUE	FALSE	0	poisson	rpart	Single, integer
FALSE	FALSE	>0	gaussian	rpart	Single, numeric (non-integer)
FALSE	FALSE	>0	binomial	rpart	Single, factor with 2 levels
FALSE	FALSE	>0	poisson	rpart	Single, integer

## Value

An object of class `pre`, which contains the initial ensemble of rules and/or linear terms and the final ensembles for a wide range of penalty parameter values. By default, the final ensemble employed by all of the other methods and functions in package `pre` is selected using the 'minimum cross validated error plus 1 standard error' criterion. All functions and methods take a `penalty.parameter.value` argument, which can be used to select a more or less sparse final ensembles. Users can assess the trade-off between sparsity and accuracy provided by every possible value of the penalty parameter ( $\lambda$ ) by running `object$glmnet.fit` and `plot(object$glmnet.fit)`.

## Note

The code for deriving rules from the nodes of trees was taken from an internal function of the `partykit` package of Achim Zeileis and Torsten Hothorn.

## References

Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *The Annals of Applied Statistics*, 29(5), 1189-1232. Friedman, J. H., & Popescu, B. E. (2008). Predictive learning via rule ensembles. *The Annals of Applied Statistics*, 2(3), 916-954. Hothorn, T., & Zeileis, A. (2015). `partykit`: A modular toolkit for recursive partytioning in R. *Journal of Machine Learning Research*, 16, 3905-3909.

## See Also

[print.pre](#), [plot.pre](#), [coef.pre](#), [importance](#), [predict.pre](#), [interact](#), [cvpre](#)

## Examples

```
set.seed(42)
```

```
airq.ens <- pre(Ozone ~ ., data = airquality[complete.cases(airquality),], verbose = TRUE)
```

---

predict.gpe	<i>Predicted values based on gpe ensemble</i>
-------------	---

---

### Description

Predict function for [gpe](#)

### Usage

```
## S3 method for class 'gpe'
predict(object, newdata = NULL, type = "link",
        penalty.par.val = "lambda.1se", ...)
```

### Arguments

object	of class <a href="#">gpe</a>
newdata	optional new data to compute predictions for
type	argument passed to <a href="#">predict.cv.glmnet</a>
penalty.par.val	argument passed to s argument of <a href="#">predict.cv.glmnet</a>
...	Unused

### Details

The initial training data is used if newdata = NULL.

### See Also

[gpe](#)

---

predict.pre	<i>Predicted values based on final unbiased prediction rule ensemble</i>
-------------	--

---

### Description

predict.pre generates predictions based on the final prediction rule ensemble, for training or new (test) observations

### Usage

```
## S3 method for class 'pre'
predict(object, newdata = NULL, type = "link",
        penalty.par.val = "lambda.1se", ...)
```

**Arguments**

object	object of class <a href="#">pre</a> .
newdata	optional dataframe of new (test) observations, including all predictor variables used for deriving the prediction rule ensemble.
type	character string. The type of prediction required; the default type = "link" is on the scale of the linear predictors. Alternatively, for count and factor outputs, type = "response" may be specified to obtain the fitted mean and fitted probabilities, respectively; type = "class" returns the predicted class membership.
penalty.par.val	character or numeric. Penalty parameter criterion to be used for selecting final model: lambda giving minimum cv error ("lambda.min") or lambda giving cv error that is within 1 standard error of minimum cv error ("lambda.1se"). Alternatively, a numeric value may be specified, corresponding to one of the values of lambda in the sequence used by <code>glmnet</code> , for which estimated cv error can be inspected by running <code>object\$glmnet.fit</code> and <code>plot(object\$glmnet.fit)</code> .
...	further arguments to be passed to <a href="#">predict.cv.glmnet</a> .

**Details**

If `newdata` is not provided, predictions for training data will be returned.

**See Also**

[pre](#), [plot.pre](#), [coef.pre](#), [importance](#), [cvpre](#), [interact](#), [print.pre](#), [predict.cv.glmnet](#)

**Examples**

```
set.seed(1)
train <- sample(1:sum(complete.cases(airquality)), size = 100)
set.seed(42)
airq.ens <- pre(Ozone ~ ., data = airquality[complete.cases(airquality),][train,])
predict(airq.ens)
predict(airq.ens, newdata = airquality[complete.cases(airquality),][-train,])
```

---

print.gpe

---

*Print a General Prediction Ensemble (gpe)*


---

**Description**

Print a General Prediction Ensemble (gpe)

**Usage**

```
## S3 method for class 'gpe'
print(x, penalty.par.val = "lambda.1se",
      digits = getOption("digits"), ...)
```

**Arguments**

x	An object of class <a href="#">pre</a> .
penalty.par.val	character or numeric. Information for which final prediction rule ensemble should be printed? The ensemble with penalty parameter criterion yielding minimum cv error ("lambda.min") or penalty parameter yielding error within 1 standard error of minimum cv error ("lambda.1se")? Alternatively, a numeric value may be specified, corresponding to one of the values of lambda in the sequence used by glmnet, for which estimated cv error can be inspected by inspecting <code>x\$glmnet.fit</code> and <code>plot(x\$glmnet.fit)</code> .
digits	Number of decimal places to print
...	Additional arguments, currently not used.

**See Also**

[print.pre](#)

---

print.pre

*Print method for objects of class pre*

---

**Description**

`print.pre` prints information about the generated prediction rule ensemble to the command line

**Usage**

```
## S3 method for class 'pre'
print(x, penalty.par.val = "lambda.1se",
      digits = getOption("digits"), ...)
```

**Arguments**

x	An object of class <a href="#">pre</a> .
penalty.par.val	character or numeric. Information for which final prediction rule ensemble should be printed? The ensemble with penalty parameter criterion yielding minimum cv error ("lambda.min") or penalty parameter yielding error within 1 standard error of minimum cv error ("lambda.1se")? Alternatively, a numeric value may be specified, corresponding to one of the values of lambda in the sequence used by glmnet, for which estimated cv error can be inspected by inspecting <code>x\$glmnet.fit</code> and <code>plot(x\$glmnet.fit)</code> .
digits	Number of decimal places to print
...	Additional arguments, currently not used.

**Details**

Note that the cv error is estimated with data that was also used for learning rules and may be too optimistic. Use `cvpre()` to obtain a more realistic estimate of future prediction error.

**Value**

Prints information about the fitted prediction rule ensemble.

**See Also**

[pre](#), [plot.pre](#), [coef.pre](#), [importance](#), [predict.pre](#), [interact](#), [cvpre](#)

**Examples**

```
set.seed(42)
airq.ens <- pre(Ozone ~ ., data = airquality[complete.cases(airquality),])
print(airq.ens)
```

---

rTerm

*Wrapper Functions for terms in gpe*


---

**Description**

Wrapper functions for terms in `gpe`.

**Usage**

```
rTerm(x)
```

```
lTerm(x, lb = -Inf, ub = Inf, scale = 1/0.4)
```

```
eTerm(x, scale = 1/0.4)
```

**Arguments**

<code>x</code>	Input symbol.
<code>lb</code>	Lower quantile when winsorizing. <code>-Inf</code> yields no winsorizing in the lower tail.
<code>ub</code>	Lower quantile when winsorizing. <code>Inf</code> yields no winsorizing in the upper tail.
<code>scale</code>	Inverse value to time <code>x</code> by. Usually the standard deviation is used. $0.4/scale$ is used as the multiplier as suggested in Friedman & Popescu (2008) and gives each linear term the same a-priori influence as a typical rule.

**Details**

The motivation to use wrappers is to ease getting the different terms as shown in the examples and to simplify the formula passed to `cv.glmnet` in `gpe`. `lTerm` potentially rescales and/or winsorizes `x` depending on the input. `eTerm` potentially rescale `x` depending on the input.

**Value**

x potentially transformed with additional information provided in the attributes.

**References**

Friedman, J. H., & Popescu, B. E. (2008). Predictive learning via rule ensembles. *The Annals of Applied Statistics*, 2(3), 916-954.

**See Also**

[gpe](#), [gpe\\_trees](#) [gpe\\_linear](#) [gpe\\_earth](#)

**Examples**

```
mt <- terms(
  ~ rTerm(x1 < 0) + rTerm(x2 > 0) + lTerm(x3) + eTerm(x4),
  specials = c("rTerm", "lTerm", "eTerm"))
attr(mt, "specials")
# $rTerm
# [1] 1 2
#
# $lTerm
# [1] 3
#
# $eTerm
# [1] 4
```

---

singleplot

*Create partial dependence plot for a single variable in a prediction rule ensemble (pre)*

---

**Description**

singleplot creates a partial dependence plot, which shows the effect of a predictor variable on the ensemble's predictions

**Usage**

```
singleplot(object, varname, penalty.par.val = "lambda.1se", nvals = NULL,
  type = "response", ...)
```

**Arguments**

object            an object of class [pre](#)

<code>varname</code>	character vector of length one, specifying the variable for which the partial dependence plot should be created. <code>penalty.par.val</code> character. Penalty parameter criterion to be used for selecting final model: <code>lambda</code> giving minimum cv error ("lambda.min") or <code>lambda</code> giving cv error that is within 1 standard error of minimum cv error ("lambda.1se"). Alternatively, a numeric value may be specified, corresponding to one of the values of <code>lambda</code> in the sequence used by <code>glmnet</code> , for which estimated cv error can be inspected by running <code>object\$glmnet.fit</code> and <code>plot(object\$glmnet.fit)</code> .
<code>penalty.par.val</code>	character. Penalty parameter criterion to be used for selecting final model: <code>lambda</code> giving minimum cv error ("lambda.min") or <code>lambda</code> giving cv error that is within 1 standard error of minimum cv error ("lambda.1se"). Alternatively, a numeric value may be specified, corresponding to one of the values of <code>lambda</code> in the sequence used by <code>glmnet</code> , for which estimated cv error can be inspected by running <code>object\$glmnet.fit</code> and <code>plot(object\$glmnet.fit)</code> .
<code>nvals</code>	optional numeric vector of length one. For how many values of <code>x</code> should the partial dependence plot be created?
<code>type</code>	character string. Type of prediction to be plotted on y-axis. <code>type = "response"</code> gives fitted values for continuous outputs and fitted probabilities for nominal outputs. <code>type = "link"</code> gives fitted values for continuous outputs and linear predictor values for nominal outputs.
<code>...</code>	Further arguments to be passed to <code>plot.default</code> .

### Details

By default, a partial dependence plot will be created for each unique observed value of the specified predictor variable. When the number of unique observed values is large, this may take a long time to compute. In that case, specifying the `nvals` argument can substantially reduce computing time. When the `nvals` argument is supplied, values for the minimum, maximum, and (`nvals - 2`) intermediate values of the predictor variable will be plotted. Note that `nvals` can be specified only for numeric and ordered input variables. If the plot is requested for a nominal input variable, the `nvals` argument will be ignored and a warning is printed.

### See Also

[pre](#), [pairplot](#)

### Examples

```
set.seed(42)
airq.ens <- pre(Ozone ~ ., data = airquality[complete.cases(airquality),])
singleplot(airq.ens, "Temp")
```

# Index

## \*Topic **datasets**

carrillo, 3

bsnullinteract, 2, 17

carrillo, 3

coef.glmnet, 5

coef.gpe, 4

coef.pre, 5, 5, 8, 25, 27, 29

colorRampPalette, 7

contour, 20

corplot, 6

ctree, 11, 13, 23

ctree\_control, 11, 13, 24

cv.glmnet, 9, 10, 13, 24, 29

cvpre, 6, 7, 25, 27, 29

earth, 13

eTerm, 14

eTerm(rTerm), 29

family, 22

glmtree, 11, 13, 23

gpar, 21

gpe, 4, 8, 10, 12, 14, 26, 29, 30

gpe\_cv.glmnet, 9, 10

gpe\_earth, 8, 9, 30

gpe\_earth(gpe\_trees), 12

gpe\_linear, 8, 9, 30

gpe\_linear(gpe\_trees), 12

gpe\_rules\_pre, 10

gpe\_sample, 9, 12

gpe\_trees, 8, 9, 12, 30

image, 20

importance, 6, 8, 14, 21, 25, 27, 29

interact, 2, 3, 6, 8, 16, 25, 27, 29

lTerm, 14

lTerm(rTerm), 29

maxdepth\_sampler, 18, 23

mob\_control, 11, 24

pairplot, 19, 31

par, 6, 7

persp, 20

plot.default, 31

plot.pre, 6, 8, 21, 25, 27, 29

pre, 2, 3, 5–9, 14, 16–18, 20–22, 22, 27–31

predict.cv.glmnet, 26, 27

predict.gpe, 26

predict.pre, 6, 8, 25, 26, 29

print.gpe, 27

print.pre, 6, 8, 22, 25, 27, 28, 28

rainbow\_hcl, 7

rpart, 11, 24

rpart.control, 11, 24

rTerm, 14, 29

singleplot, 21, 30