

# Package: varPro (via r-universe)

June 5, 2026

**Version** 3.1.0

**Date** 2026-04-19

**Title** Model-Independent Variable Selection via the Rule-Based Variable Priority

**Author** Min Lu [aut], Aster K. Shear [aut], Udaya B. Kogalur [aut, cre], Hemant Ishwaran [aut]

**Maintainer** Udaya B. Kogalur <ubk@kogalur.com>

**BugReports** <https://github.com/kogalur/varPro/issues/>

**Depends** R (>= 4.3.0),

**Imports** randomForestSRC (>= 3.4.5), glmnet, parallel, foreach, gbm, BART, survival

**Suggests** mlbench, doMC, caret, MASS, igraph

**SystemRequirements** OpenMP

**Description** A new framework of variable selection, which instead of generating artificial covariates such as permutation importance and knockoffs, creates release rules to examine the affect on the response for each covariate where the conditional distribution of the response variable can be arbitrary and unknown.

**License** GPL (>=3)

**URL** <https://www.varprotools.org/> <https://www.luminwin.net/>  
<https://ishwaran.org/>

**Config/pak/sysreqs** cmake libglpk-dev make libicu-dev libuv1-dev libxml2-dev libx11-dev

**Repository** <https://kogalur.r-universe.dev>

**Date/Publication** 2026-04-19 14:30:46 UTC

**RemoteUrl** <https://github.com/kogalur/varpro>

**RemoteRef** HEAD

**RemoteSha** dbbb4d6351c6eeb93ac44659b5b8441773d42ce8

## Contents

alzheimers	2
cv.varpro	4
gliomas	8
hrrecov	9
importance	12
importance.varpro	15
isopro	18
ivarpro	22
outpro	28
partial.ivarpro	31
partialpro	35
plot.ivarpro	39
plot.partialpro	43
predict.isopro	46
predict.ivarpro	48
predict.uvarpro	51
predict.varpro	53
uvarpro	55
varpro	61
varpro.news	69
varpro.strength	69
<b>Index</b>	<b>71</b>

---

alzheimers	<i>Alzheimer's Disease Dataset</i>
------------	------------------------------------

---

### Description

Health, lifestyle, and clinical data for 2,149 individuals used for studying Alzheimer's Disease. Variables include demographics, cognitive assessments, medical conditions, and symptoms.

### Usage

```
data(alzheimers)
```

### Format

A data frame with 2,149 observations on the following variables:

**Age:** Age in years (60 to 90).

**Gender:** Gender (0 = Male, 1 = Female).

**Ethnicity:** Ethnicity (0 = Caucasian, 1 = African American, 2 = Asian, 3 = Other).

**EducationLevel:** Education level (0 = None, 1 = High School, 2 = Bachelor's, 3 = Higher).

**BMI:** Body Mass Index (15 to 40).

**Smoking:** Smoking status (0 = No, 1 = Yes).

**AlcoholConsumption:** Weekly alcohol consumption in units (0 to 20).

**PhysicalActivity:** Weekly physical activity in hours (0 to 10).

**DietQuality:** Diet quality score (0 to 10).

**SleepQuality:** Sleep quality score (4 to 10).

**FamilyHistoryAlzheimers:** Family history of Alzheimer's (0 = No, 1 = Yes).

**CardiovascularDisease:** Cardiovascular disease (0 = No, 1 = Yes).

**Diabetes:** Diabetes (0 = No, 1 = Yes).

**Depression:** Depression (0 = No, 1 = Yes).

**HeadInjury:** History of head injury (0 = No, 1 = Yes).

**Hypertension:** Hypertension (0 = No, 1 = Yes).

**SystolicBP:** Systolic blood pressure (90 to 180 mmHg).

**DiastolicBP:** Diastolic blood pressure (60 to 120 mmHg).

**CholesterolTotal:** Total cholesterol (150 to 300 mg/dL).

**CholesterolLDL:** LDL cholesterol (50 to 200 mg/dL).

**CholesterolHDL:** HDL cholesterol (20 to 100 mg/dL).

**CholesterolTriglycerides:** Triglycerides (50 to 400 mg/dL).

**MMSE:** Mini-Mental State Examination score (0 to 30). Lower is worse.

**FunctionalAssessment:** Functional score (0 to 10). Lower is worse.

**MemoryComplaints:** Memory complaints (0 = No, 1 = Yes).

**BehavioralProblems:** Behavioral problems (0 = No, 1 = Yes).

**ADL:** Activities of Daily Living score (0 to 10). Lower is worse.

**Confusion:** Presence of confusion (0 = No, 1 = Yes).

**Disorientation:** Presence of disorientation (0 = No, 1 = Yes).

**PersonalityChanges:** Presence of personality changes (0 = No, 1 = Yes).

**DifficultyCompletingTasks:** Difficulty completing tasks (0 = No, 1 = Yes).

**Forgetfulness:** Forgetfulness (0 = No, 1 = Yes).

**Diagnosis:** Alzheimer's diagnosis (No, Yes).

## Details

This dataset is suitable for modeling Alzheimer's risk, performing exploratory analysis, and evaluating statistical and machine learning algorithms. All individuals are uniquely identified and evaluated on a standardized set of clinical and behavioral measures.

## Source

Rabie El Kharoua (2024). Alzheimer's Disease Dataset. Available from Kaggle at <https://www.kaggle.com/datasets/rabieelkharoua/alzheimers-disease-dataset>

**Examples**

```
## load the data
data(alzheimers, package = "varPro")
o <- varpro(Diagnosis~., alzheimers)
imp <- importance(o)
print(imp)
```

cv.varpro

*Cross-Validated Cutoff Value for Variable Priority (VarPro)***Description**

Selects Cutoff Value for Variable Priority (VarPro).

**Usage**

```
cv.varpro(formula, data, nvar = 30, ntree = 150,
          local.std = TRUE, zcut = seq(0.1, 2, length = 50), nblocks = 10,
          split.weight = TRUE, split.weight.method = NULL, sparse = TRUE,
          nodesize = NULL, max.rules.tree = 150, max.tree = min(150, ntree),
          verbose = FALSE, seed = NULL, fast = FALSE, crps = FALSE, ...)
```

**Arguments**

formula	Model formula specifying the outcome and predictors.
data	Training data set (data frame).
nvar	Maximum number of variables to return.
ntree	Number of trees to grow.
local.std	Use locally standardized importance values?
zcut	Grid of positive cutoff values used for selecting top variables.
nblocks	Number of blocks (folds) for cross-validation.
split.weight	Use guided tree-splitting? Variables are selected for splitting with probability proportional to split-weights, obtained by default from a preliminary lasso+tree step.
split.weight.method	Character string or vector specifying how split-weights are generated. Defaults to lasso+tree.
sparse	Use sparse split-weights?
nodesize	Minimum terminal node size. If not specified, an internal function sets the value based on sample size and data dimension.
max.rules.tree	Maximum number of rules per tree.
max.tree	Maximum number of trees used for rule extraction.
verbose	Print verbose output?

seed	Seed for reproducibility.
fast	Use <code>rfsrc.fast</code> in place of <code>rfsrc</code> ? May improve speed at the cost of accuracy.
crps	Use CRPS (continuous ranked probability score) instead of Harrell's C-index for evaluating survival performance? Applies only to survival families.
...	Additional arguments passed to <code>varpro</code> .

### Details

Applies `VarPro` and then selects from a grid of cutoff values the cutoff value for identifying variables that minimizes out-of-sample performance (error rate) of a random forest where the forest is fit to the top variables identified by the given cutoff value.

Additionally, a "conservative" and "liberal" list of variables are returned using a one standard deviation rule. The conservative list comprises variables using the largest cutoff with error rate within one standard deviation from the optimal cutoff error rate, whereas the liberal list uses the smallest cutoff value with error rate within one standard deviation of the optimal cutoff error rate.

For class imbalanced settings (two class problems where relative frequency of labels is skewed towards one class) the code automatically switches to random forest quantile classification (RFQ; see O'Brien and Ishwaran, 2019) under the `gmean` (geometric mean) performance metric.

### Value

Output containing importance values for the optimized cutoff value. A conservative and liberal list of variables is also returned.

Note that importance values are returned in terms of the original features and not their hot-encodings. For importance in terms of hot-encodings, use the built-in wrapper `get.vimp` (see example below).

### Author(s)

Min Lu and Hemant Ishwaran

### References

Lu, M. and Ishwaran, H. (2024). Model-independent variable selection via the rule-based variable priority. *arXiv e-prints*, pp.arXiv-2409.

O'Brien R. and Ishwaran H. (2019). A random forests quantile classifier for class imbalanced data. *Pattern Recognition*, 90, 232-249.

### See Also

[importance.varpro](#) [uvarpro](#) [varpro](#)

### Examples

```
## -----
## van de Vijver microarray breast cancer survival data
## high dimensional example
## -----
```

```

data(vdv, package = "randomForestSRC")
o <- cv.varpro(Surv(Time, Censoring) ~ ., vdv)
print(o)

## -----
## boston housing
## -----

data(BostonHousing, package = "mlbench")
print(cv.varpro(medv~., BostonHousing))

## -----
## boston housing - original/hot-encoded vimp
## -----

## load the data
data(BostonHousing, package = "mlbench")

## convert some of the features to factors
Boston <- BostonHousing
Boston$zn <- factor(Boston$zn)
Boston$chas <- factor(Boston$chas)
Boston$lstat <- factor(round(0.2 * Boston$lstat))
Boston$nox <- factor(round(20 * Boston$nox))
Boston$rm <- factor(round(Boston$rm))

## make cv call
o <-cv.varpro(medv~., Boston)
print(o)

## importance original variables (default)
print(get.orgvimp(o, pretty = FALSE))

## importance for hot-encoded variables
print(get.vimp(o, pretty = FALSE))

## -----
## multivariate regression example: boston housing
## vimp is collapsed across the outcomes
## -----

data(BostonHousing, package = "mlbench")
print(cv.varpro(cbind(lstat, nox) ~., BostonHousing))

## -----
## iris
## -----

print(cv.varpro(Species~., iris))

## -----
## friedman 1
## -----

```

```

print(cv.varpro(y~., data.frame(mlbench::mlbench.friedman1(1000))))

##-----
## class imbalanced problem
##
## - simulation example using the caret R-package
## - creates imbalanced data by randomly sampling the class 1 values
##
##-----

if (library("caret", logical.return = TRUE)) {

  ## experimental settings
  n <- 5000
  q <- 20
  ir <- 6
  f <- as.formula(Class ~ .)

  ## simulate the data, create minority class data
  d <- twoClassSim(n, linearVars = 15, noiseVars = q)
  d$Class <- factor(as.numeric(d$Class) - 1)
  idx.0 <- which(d$Class == 0)
  idx.1 <- sample(which(d$Class == 1), sum(d$Class == 1) / ir , replace = FALSE)
  d <- d[c(idx.0,idx.1),, drop = FALSE]
  d <- d[sample(1:nrow(d)), ]

  ## cv.varpro call
  print(cv.varpro(f, d))

}

## -----
## pbc survival with rmst vector
## note that vimp is collapsed across the rmst values
## similar to mv-regression
## -----

data(pbc, package = "randomForestSRC")
print(cv.varpro(Surv(days, status)~., pbc, rmst = c(500, 1000)))

## -----
## peak V02 with cutoff selected using fast option
## (a) C-index (default) (b) CRPS performance metric
## -----

data(peakV02, package = "randomForestSRC")
f <- as.formula(Surv(ttodead, died)~.)

## Harrel's C-index (default)
print(cv.varpro(f, peakV02, ntree = 100, fast = TRUE))

```

```

## Harrel's C-index with smaller bootstrap
print(cv.varpro(f, peakV02, ntree = 100, fast = TRUE, sampsize = 100))

## CRPS with smaller bootstrap
print(cv.varpro(f, peakV02, crps = TRUE, ntree = 100, fast = TRUE, sampsize = 100))

## -----
## largish data set: illustrates various options to speed up calculations
## -----

## roughly impute the data
data(housing, package = "randomForestSRC")
housing2 <- roughfix(housing)

## use bigger nodesize
print(cv.varpro(SalePrice~., housing2, fast = TRUE, ntree = 50, nodesize = 150))

## use smaller bootstrap
print(cv.varpro(SalePrice~., housing2, fast = TRUE, ntree = 50, nodesize = 150, sampsize = 250))

```

---

gliomas

*Diffuse Adult Glioma*


---

## Description

Subset of the data used in Ceccarelli et al. (2016) for molecular profiling of adult diffuse gliomas. As part of the analysis, the authors developed a supervised analysis using DNA methylation data. Their original dataset was collected from a core set of 25,978 CpG probes which was reduced to eliminate sites that were methylated. This reduced set of 1206 probes from 880 tissues makes up part of the features of this data. Also included are clinical data and other molecular data collected for the samples. The outcome is a supervised class label developed in the study with labels: Classic-like, Codel, G-CIMP-high, G-CIMP-low, LGM6-GBM, Mesenchymal-like and PA-like.

## References

Ceccarelli, M., Barthel, F.P., Malta, T.M., Sabedot, T.S., Salama, S.R., Murray, B.A., Morozova, O., Newton, Y., Radenbaugh, A., Pagnotta, S.M. et al. (2016). Molecular profiling reveals biologically discrete subsets and pathways of progression in diffuse glioma. *Cell*, 164, 550-563.

## Examples

```

data(glioma, package = "varPro")
o <- varpro(y~., glioma, nodesize=2, max.tree=250)
imp <- importance(o)
print(head(imp$unconditional))
print(imp$conditional.z)

```

hrrecov

*Exercise Heart Rate Recovery and Mortality***Description**

A survival dataset from a cohort of patients referred for symptom-limited exercise testing, originally used to study exercise heart rate recovery as a predictor of all-cause mortality.

**Usage**

```
data(hrrecov)
```

**Format**

A data frame with 23701 observations on 83 variables.

`t_todead` Follow-up time to death or censoring, in years (rounded to 2 decimals).

`died` Event indicator (1 = death, 0 = right-censored).

`hrrecov` Heart rate recovery (beats/min): peak heart rate minus heart rate 1 minute into recovery.

`lowrec` Indicator for low/abnormal heart rate recovery.

`lowcri` Indicator for low chronotropic response index/chronotropic incompetence.

`peak_hr` Peak heart rate during exercise (beats/min).

`peak_met` Peak workload achieved (metabolic equivalents, METs).

`fitness` Fitness category (ordinal 1-5; coding as in source data).

`heart_ra` Resting heart rate prior to testing (beats/min).

`sbprest` Resting systolic blood pressure (mm Hg).

`dbprest` Resting diastolic blood pressure (mm Hg).

`age` Age at exercise test (years).

`gender` Sex indicator (1=men).

`race` Race category (integer code; coding as in source data).

`black` Indicator for Black race (0/1).

`height` Height (m).

`weight` Weight (kg).

`bmi` Body mass index ( $\text{kg}/\text{m}^2$ ).

`bsa` Body surface area.

`wght` Weight-to-height ratio (weight/height;  $\text{kg}/\text{m}$ ).

`obese` Indicator for obesity (0/1).

`priorcad` Known or suspected coronary artery disease prior to test (0/1).

`mihist` History of myocardial infarction (0/1).

`pcabg` Prior coronary artery bypass grafting (CABG) (0/1).

ppci Prior percutaneous coronary intervention (PCI) (0/1).  
cva History of cerebrovascular accident/stroke (0/1).  
tia History of transient ischemic attack (0/1).  
pvd Peripheral vascular disease (0/1).  
diabetes History of diabetes mellitus (0/1).  
insulin Insulin therapy (0/1).  
htn History of hypertension (0/1).  
htnrx Antihypertensive treatment (0/1).  
hichol History of high cholesterol/hyperlipidemia (0/1).  
smknow Smoking history/status indicator.  
asthma History of asthma (0/1).  
copd History of chronic obstructive pulmonary disease (0/1).  
esrd End-stage renal disease (0/1).  
lv\_dysf Left ventricular dysfunction (0/1).  
ecgmi ECG evidence of myocardial infarction (0/1).  
ecglvh ECG evidence of left ventricular hypertrophy (0/1).  
lbbb Left bundle branch block (0/1).  
rbbb Right bundle branch block (0/1).  
restst Resting ST-segment abnormality (0/1).  
stnond Non-diagnostic ST-segment response (0/1).  
stabn ST-segment abnormality during testing (0/1).  
stabnb ST-segment abnormality subtype/flag.  
stabnv ST-segment abnormality subtype/flag.  
rtachy Tachycardia indicator.  
typcp Typical chest pain indicator.  
ntangina Angina history/symptom indicator.  
ttangina Angina during treadmill test (0/1).  
ttclaud Claudication during treadmill test (0/1).  
image Exercise test performed with imaging (0/1).  
acei Angiotensin-converting enzyme (ACE) inhibitor use (0/1).  
aspirin Aspirin use (0/1).  
betablok Beta-blocker use (0/1).  
dilver Diltiazem/verapamil use (0/1).  
nifed Nifedipine use (0/1).  
diuretic Diuretic use (0/1).  
lipidrx Lipid-lowering therapy use (0/1).  
nitrates Nitrate therapy use (0/1).

bdilat Bronchodilator use (0/1).  
rs\_ami Reason for referral/testing: acute myocardial infarction (0/1).  
rs\_mi Reason for referral/testing: myocardial infarction (0/1).  
rs\_ptca Reason for referral/testing: PTCA/angioplasty (0/1).  
rs\_cabg Reason for referral/testing: CABG (0/1).  
rs\_arrth Reason for referral/testing: arrhythmia (0/1).  
rs\_htran Reason for referral/testing: heart transplant (0/1).  
bpvc\_rst Bigeminal premature ventricular contractions at rest (0/1).  
bpvc\_ex Bigeminal premature ventricular contractions during exercise (0/1).  
bpvc\_rec Bigeminal premature ventricular contractions during recovery (0/1).  
fpvc\_rst Frequent premature ventricular contractions at rest (0/1).  
fpvc\_ex Frequent premature ventricular contractions during exercise (0/1).  
fpvc\_rec Frequent premature ventricular contractions during recovery (0/1).  
nsvt\_rst Non-sustained ventricular tachycardia at rest (0/1).  
nsvt\_ex Non-sustained ventricular tachycardia during exercise (0/1).  
nsvt\_rec Non-sustained ventricular tachycardia during recovery (0/1).  
vtrp\_rst Ventricular triplets at rest (0/1).  
vtrp\_ex Ventricular triplets during exercise (0/1).  
vtrp\_rec Ventricular triplets during recovery (0/1).  
hbm2\_rec Heart rate-related recovery measure at 2 minutes.  
exec Indicator variable (0/1); coding as in source data.  
aso Indicator variable (0/1); coding as in source data.

## Details

Heart rate recovery (hrrecov) is defined as peak heart rate minus the heart rate measured 1 minute into recovery. The survival outcome is all-cause mortality with right-censoring. Unless otherwise noted, indicator variables are coded 0/1 (0 = no/absent, 1 = yes/present).

## References

Ishwaran, H., Blackstone, E. H., Pothier, C. E. and Lauer, M. S. (2004). Relative risk forests for exercise heart rate recovery as a predictor of mortality. *Journal of the American Statistical Association*, 99, 591-600.

## Examples

```
data(hrrecov)
vp <- varpro(Surv(ttodead, died)~., hrrecov, ntree=50, split.weight=FALSE)
ivp <- ivarpro(vp)
plot(ivp, var = "peak_met", col.var = "fitness", size.var = "peak_hr",
      col.legend.n = 7, smooth.n = 7, x.dist = "auto")
```

importance

*Calculate Importance for VarPro and UVarPro Objects***Description**

Calculates variable importance from compatible objects such as varpro and uvarpro.

**Usage**

```
importance(x, ...)

## S3 method for class 'varpro'
importance(x, local.std = TRUE, y.external = NULL,
  cutoff = 0.79, trim = 0.1, plot.it = FALSE, conf = TRUE, sort = TRUE,
  ylab = if (conf) "Importance" else "Standardized Importance",
  max.rules.tree, max.tree,
  ...)

## S3 method for class 'uvarpro'
importance(x, local.std = FALSE, y.external = NULL,
  cutoff = 0.79, trim = 0.1, plot.it = FALSE, conf = TRUE, sort = TRUE,
  ylab = if (conf) "Importance" else "Standardized Importance",
  max.rules.tree, max.tree,
  ...)

## S3 method for class 'rhf'
importance(x, local.std = TRUE, y.external = NULL,
  cutoff = 0.79, trim = 0.1, plot.it = FALSE, conf = TRUE, sort = TRUE,
  ylab = if (conf) "Importance" else "Standardized Importance",
  max.rules.tree, max.tree,
  ...)
```

**Arguments**

x	A varpro, uvarpro or rhf object.
local.std	Logical. If TRUE, uses locally standardized importance values. Ignored for uvarpro objects.
y.external	Optional user-supplied response vector. Must match the expected dimension and outcome family. Ignored for uvarpro objects.
cutoff	Threshold used to highlight significant variables in the importance plot. Applies only when plot.it = TRUE.
trim	Windsorization trim value used to robustify the mean and standard deviation calculations.

<code>plot.it</code>	Logical. If TRUE, generates a plot of importance values.
<code>conf</code>	Logical. If TRUE, displays importance values with standard errors as a boxplot (providing an informal confidence region). If FALSE, plots standardized importance values.
<code>sort</code>	Logical. If TRUE, sorts results in decreasing order of importance.
<code>ylab</code>	Character string specifying the y-axis label.
<code>max.rules.tree</code>	Optional. Maximum number of rules per tree. Defaults to the value stored in the fitted object if unspecified.
<code>max.tree</code>	Optional. Maximum number of trees used for rule extraction. Defaults to the value from the fitted object if unspecified.
<code>...</code>	Additional arguments passed to internal methods.

### Details

This page documents the public `importance()` generic together with the methods for `varpro` and `uvarpro` objects.

The supervised `varpro` method calculates standardized importance values for identifying and ranking variables. Optionally, graphical output is provided, including confidence-style boxplots.

### Value

Invisibly, a table summarizing the results. Contains mean importance mean, the standard deviation `std`, and standardized importance `z`.

For classification, conditional `z` tables are additionally provided, where the  $z$  standardized importance values are conditional on the class label.

See `cv.varpro` for a data-driven cross-validation method for selecting the cutoff value, `cutoff`, in supervised `varpro` analyses.

### Author(s)

Min Lu and Hemant Ishwaran

### References

Lu, M. and Ishwaran, H., (2024). Model-independent variable selection via the rule-based variable priority. arXiv e-prints, pp.arXiv-2409.

### See Also

[cv.varpro](#) [varpro](#) [uvarpro](#)

### Examples

```
## -----
## toy example - needed to pass CRAN test
## -----
```

```

## mtcars regression
o <- varpro(mpg ~ ., mtcars, ntree = 1)
imp <- importance(o, local.std = FALSE)
print(imp)

## -----
## iris example
## -----

## apply varpro to the iris data
o <- varpro(Species ~ ., iris, max.tree = 5)

## print/plot the results
imp <- importance(o, plot.it = TRUE)
print(imp)

## -----
## boston housing: regression
## -----

data(BostonHousing, package = "mlbench")

## call varpro
o <- varpro(medv~., BostonHousing)

## extract importance values
imp <- importance(o)
print(imp)

## plot the results
imp <- importance(o, plot.it = TRUE)
print(imp)

## -----
## illustrates y-external: regression example
## -----

## friedman1 - standard application of varpro
d <- data.frame(mlbench::mlbench.friedman1(250),noise=matrix(runif(250*10,-1,1),250))
o <- varpro(y~.,d)
print(importance(o))

## importance using external rf predictor
print(importance(o,y.external=randomForestSRC::rfsrc(y~.,d)$predicted.oob))

## importance using external lm predictor
print(importance(o,y.external=lm(y~.,d)$fitted))

```

```

## importance using external randomized predictor
print(importance(o,y.external=sample(o$y)))

## -----
## illustrates y-external: classification example
## -----

## iris - standard application of varpro
o <- varpro(Species~.,iris)
print(importance(o))

## importance using external rf predictor
print(importance(o,y.external=randomForestSRC::rfsrc(Species~.,iris)$class.oob))

## importance using external randomized predictor
print(importance(o,y.external=sample(o$y)))

## -----
## illustrates y-external: survival
## -----
data(pbc, package = "randomForestSRC")
o <- varpro(Surv(days, status)~., pbc)
print(importance(o))

## importance using external rsf predictor
print(importance(o,y.external=randomForestSRC::rfsrc(Surv(days, status)~., pbc)$predicted.oob))

## importance using external randomized predictor
print(importance(o,y.external=sample(o$y)))

```

---

importance.varpro      *Calculate VarPro Importance*

---

## Description

Calculates variable importance using results from previous varpro call.

## Usage

```

## S3 method for class 'varpro'
importance(o, local.std = TRUE, y.external = NULL,
  cutoff = 0.79, trim = 0.1, plot.it = FALSE, conf = TRUE, sort = TRUE,
  ylab = if (conf) "Importance" else "Standardized Importance",
  max.rules.tree, max.tree,
  ...)

```

**Arguments**

<code>o</code>	varpro object returned from a previous call to <code>varpro</code> .
<code>local.std</code>	Logical. If TRUE, uses locally standardized importance values.
<code>y.external</code>	Optional user-supplied response vector. Must match the expected dimension and outcome family.
<code>cutoff</code>	Threshold used to highlight significant variables in the importance plot. Applies only when <code>plot.it = TRUE</code> .
<code>trim</code>	Windsorization trim value used to robustify the mean and standard deviation calculations.
<code>plot.it</code>	Logical. If TRUE, generates a plot of importance values.
<code>conf</code>	Logical. If TRUE, displays importance values with standard errors as a boxplot (providing an informal confidence region). If FALSE, plots standardized importance values.
<code>sort</code>	Logical. If TRUE, sorts results in decreasing order of importance.
<code>ylab</code>	Character string specifying the y-axis label.
<code>max.rules.tree</code>	Optional. Maximum number of rules per tree. Defaults to the value stored in the varpro object if unspecified.
<code>max.tree</code>	Optional. Maximum number of trees used for rule extraction. Defaults to the value from the varpro object if unspecified.
<code>...</code>	Additional arguments passed to internal methods.

**Details**

Calculates standardized importance values for identifying and ranking variables. Optionally, graphical output is provided, including confidence-style boxplots.

**Value**

Invisibly, table summarizing the results. Contains mean importance 'mean', the standard deviation 'std', and standardized importance 'z'.

For classification, conditional 'z' tables are additionally provided, where the  $z$  standardized importance values are conditional on the class label.

See `cv.varpro` for a data driven cross-validation method for selecting the cutoff value, `cutoff`.

**Author(s)**

Min Lu and Hemant Ishwaran

**References**

Lu, M. and Ishwaran, H., (2024). Model-independent variable selection via the rule-based variable priority. arXiv e-prints, pp.arXiv-2409.

**See Also**

[cv.varpro varpro](#)

**Examples**

```
## -----  
## toy example - needed to pass CRAN test  
## -----  
  
## mtcars regression  
o <- varpro(mpg ~ ., mtcars, ntree = 1)  
imp <- importance(o, local.std = FALSE)  
print(imp)  
  
## -----  
## iris example  
## -----  
  
## apply varpro to the iris data  
o <- varpro(Species ~ ., iris, max.tree = 5)  
  
## print/plot the results  
imp <- importance(o, plot.it = TRUE)  
print(imp)  
  
## -----  
## boston housing: regression  
## -----  
  
data(BostonHousing, package = "mlbench")  
  
## call varpro  
o <- varpro(medv ~ ., BostonHousing)  
  
## extract importance values  
imp <- importance(o)  
print(imp)  
  
## plot the results  
imp <- importance(o, plot.it = TRUE)  
print(imp)  
  
## -----  
## illustrates y-external: regression example  
## -----  
  
## friedman1 - standard application of varpro  
d <- data.frame(mlbench::mlbench.friedman1(250), noise=matrix(runif(250*10,-1,1),250))  
o <- varpro(y~.,d)  
print(importance(o))
```

```

## importance using external rf predictor
print(importance(o,y.external=randomForestSRC::rfsrc(y~.,d)$predicted.oob))

## importance using external lm predictor
print(importance(o,y.external=lm(y~.,d)$fitted))

## importance using external randomized predictor
print(importance(o,y.external=sample(o$y)))

## -----
## illustrates y-external: classification example
## -----

## iris - standard application of varpro
o <- varpro(Species~.,iris)
print(importance(o))

## importance using external rf predictor
print(importance(o,y.external=randomForestSRC::rfsrc(Species~.,iris)$class.oob))

## importance using external randomized predictor
print(importance(o,y.external=sample(o$y)))

## -----
## illustrates y-external: survival
## -----
data(pbc, package = "randomForestSRC")
o <- varpro(Surv(days, status)~., pbc)
print(importance(o))

## importance using external rsf predictor
print(importance(o,y.external=randomForestSRC::rfsrc(Surv(days, status)~., pbc)$predicted.oob))

## importance using external randomized predictor
print(importance(o,y.external=sample(o$y)))

```

---

isopro

*Identify Anomalous Data*


---

## Description

Use isolation forests to identify rare/anomalous data.

## Usage

```
isopro(object,
        method = c("unsupv", "rnd", "auto"),
```

```
sampsize = function(x){min(2^6, .632 * x)},
ntree = 500, nodesize = 1,
formula = NULL, data = NULL, ...)
```

### Arguments

object	varpro object returned from a previous call.
method	Isolation forest method. Options are "unsupv" (unsupervised analysis, default), "rnd" (pure random splitting), and "auto" (auto-encoder, a type of multivariate forest).
sampsize	Function or numeric value specifying the sample size used for constructing each tree. Sampling is without replacement.
ntree	Number of trees to grow.
nodesize	Minimum terminal node size.
formula	Formula used for supervised isolation forest. Ignored if object is provided.
data	Data frame used to fit the isolation forest. Ignored if object is provided.
...	Additional arguments passed to rfsrc.

### Details

Isolation Forest (Liu et al., 2008) is a random forest-based method for detecting anomalous observations. In its original form, trees are constructed using pure random splits, with each tree built from a small subsample of the data, typically much smaller than the standard 0.632 fraction used in random forests. The idea is that anomalous or rare observations are more likely to be isolated early, requiring fewer splits to reach terminal nodes. Thus, observations with relatively small depth values (i.e., shallow nodes) are considered anomalies.

There are several ways to apply the method:

- The default approach is to supply a formula and data to build a supervised isolation forest. If only data is provided (i.e., no response), an unsupervised analysis is performed. In this case, the method option is used to specify the type of isolation forest (e.g., "unsupv", "rnd", or "auto").
- If both a formula and data are provided, a supervised model is fit. In this case, method is ignored. While less conventional, this approach may be useful in certain applications.
- Alternatively, a varpro object may be supplied, but other configurations are also supported. In this setting, isolation forest is applied to the reduced feature matrix extracted from theobject. This is similar to using the data option alone but with the advantage of prior dimension reduction.

Users are encouraged to experiment with the choice of method, as the original isolation forest ("rnd") performs well in many scenarios but can be improved upon in others. For example, in some cases, "unsupv" or "auto" may yield better detection performance.

In terms of computational cost, "rnd" is the fastest, followed by "unsupv". The slowest is "auto", which is best suited for low-dimensional settings.

**Value**

Trained isolation forest and anomaly scores.

**Author(s)**

Min Lu and Hemant Ishwaran

**References**

Liu, Fei Tony, Kai Ming Ting, and Zhi-Hua Zhou. (2008). Isolation forest. 2008 Eighth IEEE International Conference on Data Mining. IEEE.

Ishwaran H. (2025). Multivariate Statistics: Classical Foundations and Modern Machine Learning, CRC (Chapman and Hall), in press.

**See Also**

[predict.isopro](#) [uvarpro](#) [varpro](#)

**Examples**

```
## -----
##
## satellite data: convert some of the classes to "outliers"
## unsupervised isopro analysis
##
## -----

## load data, make three of the classes into outliers
data(Satellite, package = "mlbench")
is.outlier <- is.element(Satellite$classes,
                        c("damp grey soil", "cotton crop", "vegetation stubble"))

## remove class labels, make unsupervised data
x <- Satellite[, names(Satellite)[names(Satellite) != "classes"]]

## isopro calls
i.rnd <- isopro(data=x, method = "rnd", sampsize=32)
i.uns <- isopro(data=x, method = "unsupv", sampsize=32)
i.aut <- isopro(data=x, method = "auto", sampsize=32)

## AUC and precision recall (computed using true class label information)
perf <- cbind(get.iso.performance(is.outlier,i.rnd$showbad),
             get.iso.performance(is.outlier,i.uns$showbad),
             get.iso.performance(is.outlier,i.aut$showbad))
colnames(perf) <- c("rnd", "unsupv", "auto")
print(perf)

## -----
```

```
##
## boston housing analysis
## isopro analysis using a previous VarPro (supervised) object
##
## -----

data(BostonHousing, package = "mlbench")

## call varpro first and then isopro
o <- varpro(medv~., BostonHousing)
o.iso <- isopro(o)

## identify data with extreme percentiles
print(BostonHousing[o.iso$howbad <= quantile(o.iso$howbad, .01),])

## -----
##
## boston housing analysis
## supervised isopro analysis - direct call using formula/data
##
## -----

data(BostonHousing, package = "mlbench")

## direct approach uses formula and data options
o.iso <- isopro(formula=medv~., data=BostonHousing)

## identify data with extreme percentiles
print(BostonHousing[o.iso$howbad <= quantile(o.iso$howbad, .01),])

## -----
##
## monte carlo experiment to study different methods
## unsupervised isopro analysis
##
## -----

## monte carlo parameters
nrep <- 25
n <- 1000

## simulation function
twodimsim <- function(n=1000) {
  cluster1 <- data.frame(
    x = rnorm(n, -1, .4),
    y = rnorm(n, -1, .2)
  )
  cluster2 <- data.frame(
    x = rnorm(n, +1, .2),
    y = rnorm(n, +1, .4)
  )
  outlier <- data.frame(
```

```

    x = -1,
    y = 1
  )
x <- data.frame(rbind(cluster1, cluster2, outlier))
is.outlier <- c(rep(FALSE, 2 * n), TRUE)
list(x=x, is.outlier=is.outlier)
}

## monte carlo loop
hbad <- do.call(rbind, lapply(1:nrep, function(b) {
  cat("iteration:", b, "\n")
  ## draw the data
  sim0 <- twodimsim(n)
  x <- sim0$x
  is.outlier <- sim0$is.outlier
  ## iso pro calls
  i.rnd <- isopro(data=x, method = "rnd")
  i.uns <- isopro(data=x, method = "unsupv")
  i.aut <- isopro(data=x, method = "auto")
  ## save results
  c(tail(i.rnd$howbad,1),
    tail(i.uns$howbad,1),
    tail(i.aut$howbad,1))
}))

## compare performance
colnames(hbad) <- c("rnd", "unsupv", "auto")
print(summary(hbad))
boxplot(hbad,col="blue",ylab="outlier percentile value")

```

---

ivarpro

*Individual Variable Priority (iVarPro): Case-Specific Variable Importance*


---

## Description

Individual Variable Priority (iVarPro) computes case-specific (individual-level) variable importance scores. For each observation in the data and for each predictor identified by the VarPro analysis, iVarPro returns a local gradient-based priority measure that quantifies how sensitive that case's prediction is to changes in that variable.

## Usage

```

ivarpro(object,
        adaptive = TRUE,
        cut = NULL,
        cut.max = 1,

```

```

ncut = 51,
nmin = 20, nmax = 150,
y.external = NULL,
noise.na = TRUE,
max.rules.tree = NULL,
max.tree = NULL,
use.loo = TRUE,
use.abs = FALSE,
path.store.membership = TRUE,
save.data = TRUE,
save.model = TRUE,
scale = c("local", "global", "none")

```

### Arguments

object	varpro object from a previous call to varpro, or a rfsrc object.
adaptive	Logical. If FALSE and cut is not supplied, the cut grid is constructed as <code>seq(0, cut.max, length.out = ncut)</code> . If TRUE (default) and cut is not supplied, a data-adaptive upper bound for the neighborhood scale is computed from the sample size using a simple bandwidth-style rule-of-thumb, and cut is constructed as a sequence from 0 to this data-adaptive maximum (subject to <code>cut.max</code> ). This provides a convenient way to automatically sharpen the local neighborhood for case-specific gradients when the sample size is moderate to large.
cut	Optional user-supplied sequence of $\lambda$ values used to relax the constraint region in the local linear regression model. For continuous release variables, each value in cut is calibrated so that <code>cut = 1</code> corresponds to one standard deviation of the release coordinate. If cut is supplied, it is used as-is and the arguments <code>cut.max</code> , <code>ncut</code> , and <code>adaptive</code> are ignored. For binary or one-hot encoded release variables, the full released region is used and cut does not control neighborhood size.
cut.max	Maximum value of the $\lambda$ grid used to define the local neighborhood for continuous release variables when cut is not supplied. By default, cut is constructed as <code>seq(0, cut.max, length.out = ncut)</code> (or up to a data-adaptive value if <code>adaptive = TRUE</code> ). Smaller values of <code>cut.max</code> yield more local, sharper case-specific gradients, while larger values yield smoother, more global behavior.
ncut	Length of the cut grid when cut is not supplied. The grid is constructed as <code>seq(0, cut.max, length.out = ncut)</code> (or up to an adaptively chosen maximum if <code>adaptive = TRUE</code> ).
nmin	Minimum number of observations required for fitting a local linear model.
nmax	Maximum number of observations allowed for fitting a local linear model. Internally, <code>nmax</code> is capped at 10% of the sample size.
y.external	Optional user-supplied response vector or matrix to use as the dependent variable in the local linear regression. Must have the same number of rows as the feature matrix and match the dimension and type expected for the outcome family.
noise.na	Logical. If TRUE (default), gradients for noisy or non-signal variables are set to NA; if FALSE, they are set to zero.

<code>max.rules.tree</code>	Maximum number of rules per tree. If unspecified, the value from the <code>varpro</code> object is used, while for <code>rfsrc</code> objects, a default value is used.
<code>max.tree</code>	Maximum number of trees used to extract rules. If unspecified, the value from the <code>varpro</code> object is used, while for <code>rfsrc</code> objects, a default value is used.
<code>use.loo</code>	Logical. If <code>TRUE</code> (default), leave-one-out cross-validation is used to select the best neighborhood size (i.e., the best value in <code>cut</code> ) for each rule and release variable. If <code>FALSE</code> , the neighborhood is chosen to use the largest available sample that satisfies <code>nmin</code> and <code>nmax</code> .
<code>use.abs</code>	Use the absolute gradient for individual importance? Default is <code>FALSE</code> which uses the actual gradient.
<code>path.store.membership</code>	Store the rule membership indices (OOB case IDs) in the returned object for later ladder/band calculations? Setting <code>FALSE</code> can substantially reduce memory usage when the number of rules is large, but disables ladder-based bands in <code>plot.ivarpro()</code> and prevents <code>ivarpro_band()</code> from being used. Default is <code>TRUE</code> .
<code>save.data</code>	Save the x and y data (default is <code>TRUE</code> )? Used for downstream plots.
<code>save.model</code>	Save the original object (default is <code>TRUE</code> )? Used for <code>predict.ivarpro</code> to obtain gradient values of test data.
<code>scale</code>	Character. Controls how the local slope is scaled when forming the rule-level gradient. <code>"local"</code> (default) uses the neighborhood-based standardization described in Lu and Ishwaran (2025). <code>"global"</code> uses a fixed global scale (the overall standard deviation of the predictor in the training data) rather than a neighborhood-specific scale. <code>"none"</code> returns the unscaled slope in the original predictor units (or the unscaled finite difference for binary 0/1 release variables).

## Details

Understanding individual-level (case-specific) variable importance is important in applications where decisions are made at the level of a single person, unit, or record. A predictor may have only a modest average effect, yet be highly influential for certain cases, or the direction of its effect may differ across individuals.

The VarPro framework summarizes population-level importance by defining feature-space regions using rule-based splitting and computing importance using only observed data. `iVarPro` (Lu and Ishwaran, 2025) extends this idea to the individual level by quantifying how sensitive each case's prediction is to small changes in a predictor identified by the VarPro rule set.

For each VarPro rule, `iVarPro` considers the corresponding rule-defined region and then *releases* the rule along the rule's release coordinate. Intuitively, releasing a region means keeping the other rule constraints in place while allowing additional variation in the released variable, which provides the information needed to estimate a local directional effect. A simple local linear regression is then fit on this released region, and the resulting slope is used as a local, gradient-based priority score. Case-specific scores are obtained by aggregating the relevant rule-level gradients over the rules that apply to each case.

**Scaling / standardization.** By default, `iVarPro` reports a *standardized* local gradient: the local linear regression is fit using a standardized release coordinate so that `cut = 1` corresponds to one standard deviation of the release variable within the released region. This yields a dimensionless

effect-size interpretation and helps make iVarPro values comparable across predictors with different units and variability (see Section~2.6 of Lu and Ishwaran, 2025). The argument `scale` controls how this scaling is applied: "local" uses the neighborhood-specific scale (default), "global" uses a fixed global scale for the predictor (overall SD in the training data), and "none" returns the unscaled slope in the original predictor units (or the unscaled finite-difference for binary 0/1 release variables). When using iVarPro-derived gradients for interaction screening, comparing results across scale settings can help distinguish structural interactions in the prediction surface from scale-induced modulation due to neighborhood-dependent standardization.

**Neighborhood size and `cut.max`.** For continuous release variables, the size of the local neighborhood used for slope estimation is controlled by `cut` (constructed from `cut.max` and `ncut` when not supplied). Smaller neighborhoods produce more local behavior and can better reflect sharp changes, while larger neighborhoods produce smoother, more global behavior. When `use.loo = TRUE`, the neighborhood size is chosen in a data-driven way using a leave-one-out criterion; when `use.loo = FALSE`, the choice is based on meeting the requested sample-size bounds `nmin` and `nmax`. When `adaptive = TRUE` and `cut` is not supplied, an additional sample-size based rule is used to limit the maximum neighborhood scale (subject to `cut.max`).

For binary or one-hot encoded release variables, iVarPro interprets the local effect as a scaled finite difference between the two levels (0 and 1), conditional on the other rule constraints; in this case `cut` does not control the neighborhood along the binary coordinate.

**Cut.max ladder (neighborhood sensitivity).** Because the choice of neighborhood scale can affect the estimated local gradients, iVarPro also records a *ladder* of rule-level gradient estimates across the candidate neighborhood sizes defined by the `cut` grid. These ladder values can be summarized (e.g., ranges or quantiles) and used to visualize how sensitive case-specific gradients are to the neighborhood choice, without repeatedly refitting iVarPro for many different `cut.max` values. Ladder-based case summaries require rule membership information; set `path.store.membership = TRUE` to enable ladder bands and related summaries, or leave it `FALSE` to reduce memory usage when the number of rules is very large. See examples below.

**Settings that are currently handled.** The flexibility of this framework makes it suitable for quantifying case-specific variable importance in regression, classification, and survival settings. Currently, multivariate forests are not handled.

## Value

For univariate outcomes (and two-class classification treated as a single score), a numeric `data.frame` of dimension  $n \times p$  containing case-specific (individual-level) variable priority values, where  $n$  is the number of observations and  $p$  is the number of predictors in `object$xvar.names`.

- Each row corresponds to a case (observation) in the original data.
- Each column corresponds to a predictor variable in `object$xvar.names`.

The entry in row  $i$  and column  $j$  is the iVarPro importance score for variable  $j$  for case  $i$ , measuring the local sensitivity of that case's prediction to changes in that variable. Predictors that are never used as release variables in the VarPro rule set may appear with constant NA values (when `noise.na = TRUE`) or constant zero values (when `noise.na = FALSE`).

**Ladder/path information.** The returned object carries an attribute "ivarpro.path" containing additional information used for ladder-based summaries and plotting. In particular:

`cut` The full `cut` grid used to evaluate candidate local neighborhoods.

`cut.ladder` The interior values of `cut` (excluding the endpoints) used for the `cut.max` ladder path.

`scale` The scaling option used to compute rule-level gradients ("local", "global", or "none").

`rule.imp.ladder` A numeric matrix of dimension  $R \times L$  storing rule-level gradients selected under each ladder truncation, where  $R$  is the number of retained rules and  $L = \text{length}(\text{cut.ladder})$ .

`rule.variable` Integer vector of length  $R$  giving the release-variable index for each retained rule.

`oobMembership` Optional list (length  $R$ ) giving the OOB comparison/released membership indices for each retained rule; included only when `path.store.membership = TRUE`.

Additional tuning flags and rule metadata (e.g., `use.loo`, `adaptive`, and `tree/branch` identifiers) may also be included for diagnostics.

### Author(s)

Min Lu and Hemant Ishwaran

### References

Lu, M. and Ishwaran, H. (2025). Individual variable priority: a model-independent local gradient method for variable importance. *Artificial Intelligence Review*, 58:407.

### See Also

[varpro](#) [plot.ivarpro](#)

### Examples

```
## -----
##
## survival example with shap-like plot
##
## -----

data(peakV02, package = "randomForestSRC")
o <- varpro(Surv(ttodead, died)~., peakV02, ntree = 50)

## canonical standard analysis
imp1 <- ivarpro(o)
shap.ivarpro(imp1)

## non-adaptive analysis
imp2 <- ivarpro(o, adaptive = FALSE)
shap.ivarpro(imp2)

## non-adaptive using a small cut.max
imp3 <- ivarpro(o, cut.max = 0.5, adaptive = FALSE)
shap.ivarpro(imp3)

## -----
```

```

##
## synthetic regression example with plot
##
## -----

## true regression function
true.function <- function(which.simulation) {
  if (which.simulation == 1) {
    function(x1, x2) { 1 * (x2 <= .25) +
      15 * x2 * (x1 <= .5 & x2 > .25) +
      (7 * x1 + 7 * x2) * (x1 > .5 & x2 > .25) }
  }
  else if (which.simulation == 2) {
    function(x1, x2) { r <- x1^2 + x2^2; 5 * r * (r <= .5) }
  }
  else {
    function(x1, x2) { 6 * x1 * x2 }
  }
}

## simulation function
simfunction <- function(n = 1000, true.function, d = 20, sd = 1) {
  d <- max(2, d)
  X <- matrix(runif(n * d, 0, 1), ncol = d)
  dta <- data.frame(list(
    x = X,
    y = true.function(X[, 1], X[, 2]) + rnorm(n, sd = sd)
  ))
  colnames(dta)[1:d] <- paste("x", 1:d, sep = "")
  dta
}

## simulate the data
which.simulation <- 1
df <- simfunction(n = 500, true.function(which.simulation))

## varpro analysis
vp <- varpro(y ~ ., df)

## ivarpro analysis
imp <- ivarpro(vp)
## ivarpro analysis using a fixed global scale (optional)
imp.global <- ivarpro(vp, scale = "global")

## ivarpro analysis returning unscaled slopes (optional)
imp.none <- ivarpro(vp, scale = "none")

## plot of x2
plot(imp, var="x2")

## plot of x2 without ladder band
plot(imp, var="x2", ladder=FALSE)

```

```

## optional: use only a subset of ladder cuts
plot(imp, var="x2", ladder=TRUE, ladder.cuts=1:8)

## plot with color/size using x1 (color) and y (size)
plot(imp, var="x2", col.var="x1", size.var="y")

## -----
##
## survival example with plot
##
## -----

data(peakV02, package = "randomForestSRC")

## varpro/importance call
vp <- varpro(Surv(ttodead, died)~., peakV02)
ivp <- ivarpro(vp, adaptive = FALSE, cut.max = 2)

## plot of peak vo2
## color displays interval (a measure of exercise time)
## size displays "y" which is predicted mortality in survival
plot(ivp, var="peak.vo2", col.var="interval", size.var="y")

## same but using beta blockers for color
plot(ivp, var="peak.vo2", col.var="betablok", size.var="y")

```

---

outpro

---

*Model and subsapce aware out-of-distribution (OOD) scoring with outPro*


---

## Description

outpro computes an out-of-distribution (OOD) score for new inputs using a fitted model, integrating variable prioritization and local neighborhoods derived from the model. The procedure is model aware and subspace aware: it scores departures in the coordinates that the model has learned to rely on, rather than relying on a global distance in the full feature space. Applicable across all outcome types.

## Usage

```

outpro(object,
        newdata,
        neighbor = NULL,
        distancef = "prod",
        reduce = TRUE,
        cutoff = NULL,
        max.rules.tree = 150,

```

```

max.tree = 150)

outpro.null(object,
             nulldata = NULL,
             neighbor = NULL,
             distancef = "prod",
             reduce = TRUE,
             cutoff = .79,
             max.rules.tree = 150,
             max.tree = 150)

```

### Arguments

object	A fitted varpro object or an rfsrc object with classes <code>c("rfsrc", "grow")</code> .
newdata	New data to score. If omitted, the training design matrix is used. For varpro objects, encodings are aligned to training with <code>get.hotencode.test</code> .
neighbor	Number of training neighbors per case, as determined by the model structure. If NULL, a default of $\min(n/10, 5000)$ is used where $n$ is the number of training rows.
distancef	Distance function for aggregation. One of "prod", "euclidean", "mahalanobis", "manhattan", "minkowski", "kernel". The default is "prod".
reduce	Controls variable selection. If TRUE with a varpro object, uses model based prioritization with threshold cutoff. A character vector selects variables by name. A named numeric vector supplies variable weights. Otherwise all predictors are used with unit weights.
cutoff	Threshold used with varpro variable importance $z$ . If NULL, a default based on the number of predictors is used: .79 when the number of predictors is not large, else 0.
max.rules.tree	Maximum number of rules per tree for neighbor extraction.
max.tree	Maximum number of trees to use for neighbor extraction.
nulldata	For <code>outpro.null</code> , optional data representing an in distribution reference. If omitted, the training design matrix is used.

### Details

Out-of-distribution (OOD) detection is essential for determining when a supervised model encounters inputs that differ in ways that matter for prediction. The approach here embeds variable prioritization directly in the detection step, constructing localized, task relevant neighborhoods from the fitted model and aggregating coordinate wise deviations within the selected subspace to obtain a distance value for an input.

For a varpro object, variable prioritization is obtained from the model and controlled by cutoff. For an rfsrc object, all predictors are used unless a reduction is supplied. Distances are computed after standardizing the selected variables with training means and scales. Variables with zero standard deviation in the training data are removed automatically before scoring.

The multiplicative "prod" metric uses a small  $\epsilon$  to avoid zero multiplicands. Since differences are measured on a standardized scale,  $\epsilon$  is set automatically by default as a small fraction of the median

absolute coordinate difference across variables and neighbors; users can keep the default or pass a custom value via `out.distance` if calling it directly.

The Mahalanobis option uses absolute differences by design and the covariance of standardized training features. A small ridge is added to the covariance for numerical stability.

## Value

`outpro` returns a list with components:

- `distance`: numeric vector of length `nrow(newdata)` with one score per case.
- `distance.object`: ingredients used for distance computation, including
  - `score`: neighbor frames returned by `varpro.strength`.
  - `neighbor`: neighbor count per case.
  - `xvar.names`: selected variable names after zero sd removal.
  - `xvar.wt`: variable weights used after normalization.
  - `dist.xvar`: list of absolute coordinate difference matrices (neighbors by cases) in standardized units.
  - `xorg.scale`, `xnew.scale`: standardized training and test matrices for the selected variables.
  - `means`, `sds`: training means and scales for the selected variables.
  - `dropped.zero.sd.variables`: variables removed due to zero standard deviation in training.
- `distance.args`: list of metric arguments actually used, including `distancef`, `weights.used`, `normalize.weights`, `p`, and `epsilon.used`.
- `score`: the neighbor information returned by `varpro.strength`.
- `neighbor`: neighbor setting used.
- `cutoff`: cutoff used for variable prioritization.
- `oob.bits`: indicator of whether scoring was done on training rows or new data.
- `selected.variables`: the variables used in scoring after all filters.
- `selected.weights`: the normalized squared weights for the selected variables.
- `means`, `sds`: duplicates for convenience.
- `call`: the matched call.

`outpro.null` returns the same list with two additional components:

- `cdf`: the empirical distribution function of `distance`.
- `quantile`: the empirical cumulative probability for each scored case.

## Background

The method follows a model centered view of out-of-distribution (OOD) detection that is both model aware and subspace aware. Variable prioritization is embedded directly in the detection process to focus on coordinates that matter for prediction and to discount nuisance directions. Scoring does not rely on global feature density estimation. The implementation uses a random forest engine whose rule based structure provides localized neighborhoods reflecting the learned predictive mapping.

**See Also**

[varpro, rfsrc.](#)

**Examples**

```
## -----

## fit a varPro model
data(BostonHousing, package = "mlbench")
smp <- sample(1:nrow(BostonHousing), size = nrow(BostonHousing) * .75)
train.data <- BostonHousing[smp,]
test.data <- BostonHousing[-smp,]
vp <- varpro(medv ~ ., data = train.data)

## Score new data with default multiplicative metric
op <- outpro(vp, newdata = test.data)
head(op$distance)

## Calibrate a null distribution using training data
op.null <- outpro.null(vp)
head(op.null$quantile)
```

---

partial.ivarpro

*Partial iVarPro Plot for a Single Variable*

---

**Description**

partial.ivarpro() draws a base-graphics scatterplot of case-specific iVarPro gradients for one predictor: x is the predictor value and y is the iVarPro local gradient. Points can be colored (col.var), sized (size.var), optionally jittered, and optionally overlaid with loess smooths stratified by col.var.

**Usage**

```
partial.ivarpro(ivar,
                var,
                col.var = NULL,
                size.var = NULL,
                x = NULL,
                ladder = FALSE,
                ladder.cuts = NULL,
                ladder.max.segments = 3000,
                pch = 16,
                cex = 0.8,
```

```

cex.range = c(0.5, 2),
main = NULL,
xlab = NULL,
ylab = "iVarPro gradient",
legend = TRUE,
... )

```

## Arguments

<code>ivar</code>	An object returned by <code>ivarpro</code> . If <code>x</code> is not supplied, <code>partial.ivarpro()</code> uses <code>attr(ivar, "data")</code> when available.
<code>var</code>	Variable to plot (name or column index). Must exist in both <code>ivar</code> and <code>x</code> .
<code>col.var</code>	Optional variable (column name in <code>x</code> ) used for coloring. If treated as categorical, colors are assigned per level; if continuous, a color ramp is used and the legend shows selected quantiles.
<code>size.var</code>	Optional variable (column name in <code>x</code> ) used to scale point sizes to <code>cex.range</code> .
<code>x</code>	Optional data.frame or matrix of original feature values.
<code>ladder</code>	Logical. If TRUE, attempt to draw a ladder-based band (vertical segments spanning the range across <code>ladder.cuts</code> ) when path membership information is present.
<code>ladder.cuts</code>	Optional subset of ladder indices/values used for the band.
<code>ladder.max.segments</code>	Maximum number of ladder segments to draw.
<code>pch</code>	Point character for the default point style.
<code>cex</code>	Base point size (used when <code>size.var</code> is not supplied).
<code>cex.range</code>	Min/max point sizes used when <code>size.var</code> is supplied.
<code>main</code>	Main title (defaults to <code>paste0(var, " vs iVarPro gradient")</code> ).
<code>xlab</code>	X-axis label (defaults to <code>var</code> ).
<code>ylab</code>	Y-axis label (defaults to <code>"iVarPro gradient"</code> ).
<code>legend</code>	Logical. If TRUE and <code>col.var</code> is supplied, draw a legend describing the color mapping.
<code>...</code>	Additional arguments passed to <code>graphics::plot()</code> , plus these optional controls: <ul style="list-style-type: none"> <li><code>smooth</code> Logical. Draw loess smooth curves. Default TRUE. For categorical <code>col.var</code>, one curve per level. For continuous <code>col.var</code>, curves are drawn for strata defined by quantiles.</li> <li><code>smooth.span</code> Loess span (default 0.75). Other loess controls: <code>smooth.degree</code>, <code>smooth.family</code>, <code>smooth.lwd</code>, <code>smooth.lty</code>, <code>smooth.alpha</code>, <code>smooth.min.n</code>, <code>smooth.n.grid</code>.</li> <li><code>jitter</code> Logical or numeric. Default TRUE. Adds horizontal jitter to <code>x</code> to reduce overplotting. Related options: <code>jitter.amount</code>, <code>jitter.fraction</code>, <code>jitter.seed</code>.</li> </ul>

`x.dist` Character vector controlling an x-axis distribution strip. Default "none". Supported values include "rug", "hist", "density", and "auto" (defaults to `c("hist", "rug")`). Use `x.dist = c("hist", "rug")` (or `c("density", "rug")`) to combine.

Related options: `x.dist.side`, `x.dist.height`, `x.dist.pad`, `x.dist.col`, `x.dist.border`, `x.dist.lwd`, `x.dist.lty`, `x.dist.bins`, `x.dist.adjust`, `x.dist.n`, `x.dist.rug.col`, `x.dist.rug.lwd`, `x.dist.rug.ticks`, `x.dist.rug.max`. By default, an outline is drawn to keep the strip visible; set `x.dist.lwd = 0` to suppress the outline.

`col.legend.probs`, `col.legend.n` Quantiles shown in the legend when `col.var` is continuous. Default is 5 quantiles: `c(0.05, 0.25, 0.5, 0.75, 0.95)`.

`smooth.probs`, `smooth.n` Quantiles defining strata for smooth curves when `col.var` is continuous. Default matches the legend quantiles.

`zero.line` Logical. Default TRUE. Adds a dashed reference line at gradient 0 ( $y = 0$ ). Related options: `zero.line.col`, `zero.line.lty`, `zero.line.lwd`.

`col.var.discrete.max` If `col.var` is numeric with at most this many distinct values, treat it as categorical (default 10).

`col.style` How categorical colors are rendered: "auto" (default), "solid", "outline", "binary". Related options include `col.outline`, `col.binary.pch`, `col.dodge`.

## Details

Jitter is applied for visualization only; loess smooths and the optional x-axis distribution strip (`x.dist`) are computed using theunjittered x-values.

## Value

Invisibly returns TRUE.

## See Also

[ivarpro](#), [varpro](#), [shap.ivarpro](#)

## Examples

```
## -----
##
## Survival example: peakV02 partial plot
##
## -----

data(peakV02, package = "randomForestSRC")

ipv <- ivarpro(varpro(Surv(ttodead, died) ~ ., peakV02))

## Continuous col.var: legend + smooth strata default to 5 quantiles
partial.ivarpro(ipv, var = "peak.vo2", col.var = "interval", size.var = "y")
```

```

## Add an x-axis distribution strip (histogram + rug)
partial.ivarpro(ipv, var = "peak.vo2", col.var = "interval", size.var = "y",
  x.dist = c("hist", "rug"))

## Increase legend/smooth strata (e.g., 7 quantiles)
partial.ivarpro(ipv, var = "peak.vo2", col.var = "interval", size.var = "y",
  col.legend.n = 7, smooth.n = 7)

## Classic 3-quantile view (5%, 50%, 95%)
partial.ivarpro(ipv, var = "peak.vo2", col.var = "interval", size.var = "y",
  col.legend.probs = c(0.05, 0.5, 0.95),
  smooth.probs = c(0.05, 0.5, 0.95))

## Factor col.var example: one smooth per level
partial.ivarpro(ipv, var = "peak.vo2", col.var = "betablok", size.var = "y")

## -----
##
## multiclass example: iris (use target= to choose a class)
##
## -----

data(iris)

vp.ir <- varpro(Species ~ ., iris, ntree = 50)
ivp.ir <- ivarpro(vp.ir)

## Plot gradients for the "setosa" class (target selects the list element)
partial.ivarpro(ivp.ir, var = "Petal.Length", target = "setosa",
  col.var = "Species", x = iris)

## Alternatively, color by the predicted class probability stored in ivp.ir
partial.ivarpro(ivp.ir, var = "Petal.Length", target = "setosa",
  col.var = "y.setosa")

## -----
##
## multiclass example: wine (advanced example)
##
## -----

data(wine, package = "randomForestSRC")

## Give the class labels nicer names than "3", "4", ..., "9"
wine$quality <- factor(wine$quality)
levels(wine$quality) <- paste0("Q", levels(wine$quality)) # "Q3" "Q4" ... "Q9"

vp <- varpro(quality ~ ., wine, ntree = 50)
ivp <- ivarpro(vp)

## Available targets correspond to class probability columns
names(ivp)

```

```

## Build a plotting data.frame that contains:
## - predictors (from attr(ivp,"data"))
## - predicted class probabilities (y.Q3, y.Q4, ..., y.Q9)
## - the observed class label (quality)
xdat <- attr(ivp, "data")
xdat$quality <- wine$quality

## Plot gradients for the "Q7" class:
## y-axis: d P(Y="Q7" | x) / d alcohol (iVarPro gradient)
## x-axis: alcohol
## color: observed class label
partial.ivarpro(ivp, var = "alcohol", target = "Q7",
               x = xdat, col.var = "quality")

## Alternatively, color by the model's predicted probability for the same class.
## (These columns come from attr(ivp,"data") as y.<class>)
partial.ivarpro(ivp, var = "alcohol", target = "Q7", col.var = "y.Q7")

```

---

partialpro

*Partial Effects for Variable(s)*


---

## Description

Obtain the partial effect of x-variables from a VarPro analysis.

## Usage

```

partialpro(object, xvar.names, nvar,
           target, learner, newdata, method = c("unsupv", "rnd", "auto"),
           verbose = FALSE, ...)

```

## Arguments

object	varpro object returned from a previous call to varpro.
xvar.names	Names of the x-variables to use.
nvar	Number of variables to include. Defaults to all.
target	For classification, specifies the class for which the partial effect is computed. Can be an integer or character label. Defaults to the last class.
learner	Optional function specifying a user-defined prediction model. See Details.
newdata	Optional data frame containing test features. If not provided, the training data is used.
method	Isolation forest method used for Unlimited Virtual Twins (UVT). Options are "unsupv" (default), "rnd" (pure random splitting), and "auto" (autoencoder). See isopro for details.

verbose      Print verbose output?  
 ...          Additional hidden options: "cut", "nsmp", "nvirtual", "nmin", "alpha",  
               "df", "sampsiz", "ntree", "nodesize", "mse.tolerance".

## Details

Computes partial effects for selected variables based on a VarPro analysis. If a variable was filtered out during VarPro (e.g., due to noise), its partial effect cannot be computed.

Partial effects are derived using predictions from the forest built during VarPro. These predictions are restricted using Unlimited Virtual Twins (UVT), which apply an isolation forest criterion to filter unlikely combinations of partial data. The filtering threshold is governed by the internal cut parameter. Isolation forests are constructed via `isopro`.

Interpretation of partial effects depends on the outcome type:

- For regression: effects are on the response scale.
- For survival: effects are either on mortality (default) or RMST (if specified in the original `varpro` call).
- For classification: effects are log-odds for the specified target class.

Partial effects are estimated locally using polynomial linear models fit to the predicted values. The degrees of freedom for the local model are controlled by the `df` option (default = 2, i.e., quadratic).

By default, predictions use the forest from the VarPro object. Alternatively, users may supply a custom prediction function via `learner`. This function should accept a data frame of features and return:

- A numeric vector for regression or survival outcomes.
- A matrix of class probabilities (one column per class, in original class order) for classification.
- If `newdata` is missing, the function should return predictions on the original training data.

See the examples for use cases with external learners, including:

1. Random forest (external to VarPro),
2. Gradient tree boosting,
3. Bayesian Additive Regression Trees (BART).

## Value

Named list, with entries containing the partial plot information for a variable.

## Author(s)

Min Lu and Hemant Ishwaran

## References

Ishwaran H. (2025). *Multivariate Statistics: Classical Foundations and Modern Machine Learning*, CRC (Chapman and Hall), in press.

**See Also**[varpro isopro](#)**Examples**

```

##-----
##
## Boston housing
##
##-----

library(mlbench)
data(BostonHousing)
oldpar <- par(mfrow=c(2,3))
plot((oo.boston<-partialpro(varpro(medv~.,BostonHousing),nvar=6)))
par(oldpar)

##-----
##
## Boston housing using newdata option
##
##
##-----

library(mlbench)
data(BostonHousing)
o <- varpro(medv~.,BostonHousing)
oldpar <- par(mfrow=c(2,3))
plot(partialpro(o,nvar=3))
## same but using newdata (set to first 6 cases of the training data)
plot(partialpro(o,newdata=o$x[1:6,],nvar=3))
par(oldpar)

##-----
##
## Boston housing with externally constructed rf learner
##
##-----

## varpro analysis
library(mlbench)
data(BostonHousing)
o <- varpro(medv~.,BostonHousing)

## default partial pro call
pro <- partialpro(o, nvar=3)

## partial pro call using built in rf learner
mypro <- partialpro(o, nvar=3, learner=rf.learner(o))

```

```

## compare the two
oldpar <- par(mfrow=c(2,3))
plot(pro)
plot(mypro, ylab="external rf learner")
par(oldpar)

##-----
##
## Boston housing: tree gradient boosting learner, bart learner
##
##-----

if (library("gbm", logical.return=TRUE) &&
    library("BART", logical.return=TRUE)) {

## varpro analysis
library(parallel)
library(mlbench)
data(BostonHousing)
o <- varpro(medv~.,BostonHousing)

## default partial pro call
pro <- partialpro(o, nvar=3)

## partial pro call using built in gradient boosting learner
mypro <- partialpro(o, nvar=3, learner=gbm.learner(o, n.trees=1000, n.cores=get.mc.cores()))

## partial pro call using built in bart learner
mypro2 <- partialpro(o, nvar=3, learner=bart.learner(o, mc.cores=get.mc.cores()))

## compare the learners
oldpar <- par(mfrow=c(3,3))
plot(pro)
plot(mypro, ylab="external boosting learner")
plot(mypro2, ylab="external bart learner")
par(oldpar)
}

##-----
##
## peak vo2 with 5 year rmst
##
##-----

data(peakV02, package = "randomForestSRC")
oldpar <- par(mfrow=c(2,3))
plot((oo.peak<-partialpro(varpro(Surv(ttodead,died)~.,peakV02,rmst=5),nvar=6)))
par(oldpar)

##-----
##
## veteran data set with celltype as a factor
##

```

```

##-----
data(veteran, package = "randomForestSRC")
dta <- veteran
dta$celltype <- factor(dta$celltype)
oldpar <- par(mfrow=c(2,3))
plot((oo.veteran<-partialpro(varpro(Surv(time, status)~., dta), nvar=6)))
par(oldpar)

##-----
##
## iris: classification analysis showing partial effects for all classes
##
##-----

o.iris <- varpro(Species~.,iris)
yl <- paste("log-odds", levels(iris$Species))
oldpar <- par(mfrow=c(3,2))
plot((oo.iris.1 <- partialpro(o.iris, target=1, nvar=2)),ylab=yl[1])
plot((oo.iris.2 <- partialpro(o.iris, target=2, nvar=2)),ylab=yl[2])
plot((oo.iris.3 <- partialpro(o.iris, target=3, nvar=2)),ylab=yl[3])
par(oldpar)

##-----
##
## iowa housing data
##
##-----

## quickly impute the data; log transform the outcome
data(housing, package = "randomForestSRC")
housing <- randomForestSRC::impute(SalePrice~., housing, splitrule="random", nimpute=1)
dta <- data.frame(data.matrix(housing))
dta$y <- log(housing$SalePrice)
dta$SalePrice <- NULL

## partial effects analysis
o.housing <- varpro(y~., dta, nvar=Inf)
oo.housing <- partialpro(o.housing,nvar=15)
oldpar <- par(mfrow=c(3,5))
plot(oo.housing)
par(oldpar)

```

## Description

plot.ivarpro() draws a base-graphics scatterplot of case-specific iVarPro gradients for one predictor: the x-axis is the predictor value and the y-axis is the iVarPro local gradient. Points can be colored (col.var), sized (size.var), optionally jittered, and optionally overlaid with loess smooths stratified by col.var.

## Usage

```
## S3 method for class 'ivarpro'
plot(x,
     var,
     col.var = NULL,
     size.var = NULL,
     data = NULL,
     target = NULL,
     ladder = FALSE,
     ladder.cuts = NULL,
     ladder.max.segments = 3000,
     pch = 16,
     cex = 0.8,
     cex.range = c(0.5, 2),
     main = NULL,
     xlab = NULL,
     ylab = "iVarPro gradient",
     legend = TRUE,
     ...)
```

## Arguments

x	An object returned by <code>ivarpro</code> or <code>predict.ivarpro</code> . If data is not supplied, <code>plot.ivarpro()</code> uses <code>attr(x, "data")</code> when available.
var	Variable to plot (name or column index). Must exist in both the iVarPro object and data.
col.var	Optional variable (column name in data) used for coloring. If treated as categorical, colors are assigned per level; if continuous, a color ramp is used and the legend shows selected quantiles.
size.var	Optional variable (column name in data) used to scale point sizes to <code>cex.range</code> .
data	Optional data.frame or matrix of feature values used for the x-axis variable and any auxiliary variables referenced by <code>col.var</code> or <code>size.var</code> .
target	Optional target/outcome coordinate when x is a list-valued iVarPro result (for example multiclass or multivariate output). May be specified as a single name or index. If omitted, the first element is used.
ladder	Logical. If TRUE, attempt to draw a ladder-based band (vertical segments spanning the range across <code>ladder.cuts</code> ) when path membership information is present.
ladder.cuts	Optional subset of ladder indices/values used for the band.

ladder.max.segments	Maximum number of ladder segments to draw.
pch	Point character for the default point style.
cex	Base point size (used when size.var is not supplied).
cex.range	Min/max point sizes used when size.var is supplied.
main	Main title (defaults to paste0(var, " vs iVarPro gradient")).
xlab	X-axis label (defaults to var).
ylab	Y-axis label (defaults to "iVarPro gradient").
legend	Logical. If TRUE and col.var is supplied, draw a legend describing the color mapping.
...	Additional arguments passed to graphics::plot(), plus these optional controls:
smooth	Logical. Draw loess smooth curves. Default TRUE. For categorical col.var, one curve per level. For continuous col.var, curves are drawn for strata defined by quantiles.
smooth.span	Loess span (default 0.75). Other loess controls: smooth.degree, smooth.family, smooth.lwd, smooth.lty, smooth.alpha, smooth.min.n, smooth.n.grid.
jitter	Logical or numeric. Default TRUE. Adds horizontal jitter to the x-values to reduce overplotting. Related options: jitter.amount, jitter.fraction, jitter.seed.
x.dist	Character vector controlling an x-axis distribution strip. Default "none". Supported values include "rug", "hist", "density", and "auto" (defaults to c("hist", "rug")). Use x.dist = c("hist", "rug") (or c("density", "rug")) to combine. Related options: x.dist.side, x.dist.height, x.dist.pad, x.dist.col, x.dist.border, x.dist.lwd, x.dist.lty, x.dist.bins, x.dist.adjust, x.dist.n, x.dist.rug.col, x.dist.rug.lwd, x.dist.rug.ticks, x.dist.rug.max. By default, an outline is drawn to keep the strip visible; set x.dist.lwd = 0 to suppress the outline.
col.legend.probs, col.legend.n	Quantiles shown in the legend when col.var is continuous. Default is 5 quantiles: c(0.05, 0.25, 0.5, 0.75, 0.95).
smooth.probs, smooth.n	Quantiles defining strata for smooth curves when col.var is continuous. Default matches the legend quantiles.
zero.line	Logical. Default TRUE. Adds a dashed reference line at gradient 0 (y = 0). Related options: zero.line.col, zero.line.lty, zero.line.lwd.
col.var.discrete.max	If col.var is numeric with at most this many distinct values, treat it as categorical (default 10).
col.style	How categorical colors are rendered: "auto" (default), "solid", "outline", "binary". Related options include col.outline, col.binary.pch, col.dodge.

## Details

Jitter is applied for visualization only; loess smooths and the optional x-axis distribution strip (x.dist) are computed using the unjittered x-values.

**Value**

Invisibly returns TRUE.

**See Also**

[ivarpro](#), [varpro](#), [shap.ivarpro](#)

**Examples**

```
## -----
##
## Survival example: peakV02 plot
##
## -----

data(peakV02, package = "randomForestSRC")

ipv <- ivarpro(varpro(Surv(ttodead, died) ~ ., peakV02))

## Continuous col.var: legend + smooth strata default to 5 quantiles
plot(ipv, var = "peak.vo2", col.var = "interval", size.var = "y")

## Add an x-axis distribution strip (histogram + rug)
plot(ipv, var = "peak.vo2", col.var = "interval", size.var = "y",
      x.dist = c("hist", "rug"))

## Increase legend/smooth strata (e.g., 7 quantiles)
plot(ipv, var = "peak.vo2", col.var = "interval", size.var = "y",
      col.legend.n = 7, smooth.n = 7)

## Classic 3-quantile view (5%, 50%, 95%)
plot(ipv, var = "peak.vo2", col.var = "interval", size.var = "y",
      col.legend.probs = c(0.05, 0.5, 0.95),
      smooth.probs = c(0.05, 0.5, 0.95))

## Factor col.var example: one smooth per level
plot(ipv, var = "peak.vo2", col.var = "betablok", size.var = "y")

## -----
##
## multiclass example: iris (use target= to choose a class)
##
## -----

data(iris)

vp.ir <- varpro(Species ~ ., iris, ntree = 50)
ivp.ir <- ivarpro(vp.ir)

## Plot gradients for the "setosa" class (target selects the list element)
plot(ivp.ir, var = "Petal.Length", target = "setosa",
```

```

    col.var = "Species", data = iris)

## Alternatively, color by the predicted class probability stored in ivp.iv
plot(ivp.iv, var = "Petal.Length", target = "setosa",
     col.var = "y.setosa")

## -----
##
## multiclass example: wine (advanced example)
##
## -----

data(wine, package = "randomForestSRC")

## Give the class labels nicer names than "3", "4", ..., "9"
wine$quality <- factor(wine$quality)
levels(wine$quality) <- paste0("Q", levels(wine$quality)) # "Q3" "Q4" ... "Q9"

vp <- varpro(quality ~ ., wine, ntree = 50)
ivp <- ivarpro(vp)

## Available targets correspond to class probability columns
names(ivp)

## Build a plotting data.frame that contains:
## - predictors (from attr(ivp,"data"))
## - predicted class probabilities (y.Q3, y.Q4, ..., y.Q9)
## - the observed class label (quality)
xdat <- attr(ivp, "data")
xdat$quality <- wine$quality

## Plot gradients for the "Q7" class:
## y-axis:  $d P(Y="Q7" | x) / d \text{alcohol}$  (iVarPro gradient)
## x-axis: alcohol
## color: observed class label
plot(ivp, var = "alcohol", target = "Q7",
     data = xdat, col.var = "quality")

## Alternatively, color by the model's predicted probability for the same class.
## (These columns come from attr(ivp,"data") as y.<class>)
plot(ivp, var = "alcohol", target = "Q7", col.var = "y.Q7")

```

---

plot.partialpro

*Plot method for partialpro objects*


---

## Description

Plot partial effects from a previous partialpro analysis.

**Usage**

```
## S3 method for class 'partialpro'
plot(x, xvar.names, nvar,
     parametric = FALSE, se = TRUE,
     causal = FALSE, subset = NULL, plot.it = TRUE, ...)
```

**Arguments**

<code>x</code>	A <code>partialpro</code> object returned from a previous call to <code>partialpro</code> .
<code>xvar.names</code>	Names (or integer indices) of the x-variables to plot. Defaults to all variables.
<code>nvar</code>	Number of variables to plot. Defaults to all variables.
<code>parametric</code>	Logical. Set to <code>TRUE</code> only if the partial effect is believed to follow a polynomial form.
<code>se</code>	Display standard errors?
<code>causal</code>	Display causal estimator?
<code>subset</code>	Optional conditioning factor. Not applicable if <code>parametric = TRUE</code> . May also be a logical or integer vector to subset the analysis.
<code>plot.it</code>	If <code>FALSE</code> , no plot is produced; instead, the internal plotting objects are returned.
<code>...</code>	Additional arguments passed to <code>plot</code> .

**Details**

Generates smoothed partial-effect plots for continuous variables. The solid black line represents the estimated partial effect; dashed red lines show an approximate plus-minus standard error band. These standard errors are intended as heuristic guides and should be interpreted cautiously.

Partial effects are estimated nonparametrically using locally fitted polynomial models. This is the default behavior and is recommended when effects are expected to be nonlinear. Use `parametric = TRUE` if the underlying effect is believed to follow a global polynomial form.

For binary variables, partial effects are shown as boxplots, with whiskers reflecting variability analogous to standard error.

The causal estimator, when requested, displays the baseline-subtracted local effect.

Conditioning is supported via the `subset` option. When supplied as a factor (with length equal to the original data), the plot is stratified by its levels. Alternatively, `subset` can be a logical or integer vector indicating the cases to include in the analysis.

**Value**

If `plot.it = TRUE`, the method is called for its side effect of producing plots and returns `NULL`.

If `plot.it = FALSE`, the method returns a named list of internal plot objects, one per requested variable, containing the partial-effect curves and associated summaries used for plotting.

**Author(s)**

Min Lu and Hemant Ishwaran

## References

Ishwaran H. (2025). *Multivariate Statistics: Classical Foundations and Modern Machine Learning*, CRC (Chapman and Hall), in press.

## See Also

[partialpro](#)

## Examples

```
##-----
##
## Boston housing
##
##-----

library(mlbench)
data(BostonHousing)
o.boston <- varpro(medv~., BostonHousing)
oo.boston <- partialpro(o.boston, nvar = 4, learner = rf.learner(o.boston))

oldpar <- par(mfrow = c(2, 4))

## parametric local estimation
plot(oo.boston, parametric = TRUE, ylab = "parametric est.")

## non-parametric local estimation (default)
plot(oo.boston, parametric = FALSE, ylab = "non-parametric est.")

par(oldpar)

##-----
##
## Boston housing with subsetting
##
##-----

library(mlbench)
data(BostonHousing)
o.boston <- varpro(medv~., BostonHousing)
oo.boston <- partialpro(o.boston, nvar = 3, learner = rf.learner(o.boston))

## subset analysis
price <- BostonHousing$medv
pricef <- factor(price > median(price), labels = c("low priced", "high priced"))
oldpar <- par(mfrow = c(1, 1))
plot(oo.boston, subset = pricef, nvar = 1)
par(oldpar)

##-----
##
## veteran data with subsetting using celltype as a factor
```

```
##
##-----

data(veteran, package = "randomForestSRC")
dta <- veteran
dta$celltype <- factor(dta$celltype)
o.vet <- varpro(Surv(time, status) ~ ., dta)
oo.vet <- partialpro(o.vet, nvar = 6, nsmp = Inf, learner = rf.learner(o.vet))

## partial effects, with subsetting
oldpar <- par(mfrow = c(2, 3))
plot(oo.vet, subset = dta$celltype)
par(oldpar)

## causal effects, with subsetting
oldpar <- par(mfrow = c(2, 3))
plot(oo.vet, subset = dta$celltype, causal = TRUE)
par(oldpar)

## retrieve plotting objects without drawing
obj <- plot(oo.vet, subset = dta$celltype, plot.it = FALSE)
str(obj, max.level = 1)
```

---

predict.isopro

*Prediction for Isopro for Identifying Anomalous Data*

---

## Description

Use isolation forests to identify rare/anomalous values using test data.

## Usage

```
## S3 method for class 'isopro'
predict(object, newdata, quantiles = TRUE, ...)
```

## Arguments

object	isopro object returned from a previous call.
newdata	Optional test data. If not provided, the training data is used.
quantiles	Logical. If TRUE (default), returns quantile values; if FALSE, returns case depth values.
...	Additional arguments passed to internal methods.

**Details**

Uses a previously constructed `isopro` object to assess anomalous observations in the test data. By default, returns quantile values representing the depth of each test observation relative to the original training data. Smaller values indicate greater outlyingness.

To return raw depth values instead of quantiles, set `quantiles = FALSE`.

**Value**

Anomaly scores for the test data (or training data).

**Author(s)**

Min Lu and Hemant Ishwaran

**References**

Liu, Fei Tony, Kai Ming Ting, and Zhi-Hua Zhou. (2008). Isolation forest. 2008 Eighth IEEE International Conference on Data Mining. IEEE.

Ishwaran H. (2025). Multivariate Statistics: Classical Foundations and Modern Machine Learning, CRC (Chapman and Hall), in press.

**See Also**

[isopro](#) [uvarpro](#) [varpro](#)

**Examples**

```
## -----
##
## boston housing
## unsupervised isopro analysis
##
## -----

## training
data(BostonHousing, package = "mlbench")
o <- isopro(data=BostonHousing)

## make fake data
fake <- do.call(rbind, lapply(1:nrow(BostonHousing), function(i) {
  fakei <- BostonHousing[i,]
  fakei$lstat <- quantile(BostonHousing$lstat, .99)
  fakei$nox <- quantile(BostonHousing$nox, .99)
  fakei
}))

## compare depth values for fake data to training data
depth.fake <- predict(o, fake)
depth.train <- predict(o)
depth.data <- rbind(data.frame(whichdata="fake", depth=depth.fake),
```

```

                                data.frame(whichdata="train", depth=depth.train))
boxplot(depth~whichdata, depth.data, xlab="data", ylab="depth quantiles")

## -----
##
## boston housing
## isopro supervised analysis with different split rules
##
## -----

data(BostonHousing, package="mlbench")

## supervised isopro analysis using different splitrules
o <- isopro(formula=medv~.,data=BostonHousing)
o.hvwt <- isopro(formula=medv~.,data=BostonHousing,splitrule="mse.hvwt")
o.unwt <- isopro(formula=medv~.,data=BostonHousing,splitrule="mse.unwt")

## make fake data
fake <- do.call(rbind, lapply(1:nrow(BostonHousing), function(i) {
  fakei <- BostonHousing[i,]
  fakei$lstat <- quantile(BostonHousing$lstat, .99)
  fakei$nox <- quantile(BostonHousing$nox, .99)
  fakei
}))

## compare depth values for fake data to training data
depth.train <- predict(o)
depth.hvwt.train <- predict(o.hvwt)
depth.unwt.train <- predict(o.unwt)
depth.fake <- predict(o, fake)
depth.hvwt.fake <- predict(o.hvwt, fake)
depth.unwt.fake <- predict(o.unwt, fake)
depth.data <- rbind(data.frame(whichdata="fake", depth=depth.fake),
                    data.frame(whichdata="fake.hvwt", depth=depth.hvwt.fake),
                    data.frame(whichdata="fake.unwt", depth=depth.unwt.fake),
                    data.frame(whichdata="train", depth=depth.train),
                    data.frame(whichdata="train.hvwt", depth=depth.hvwt.train),
                    data.frame(whichdata="train.unwt", depth=depth.unwt.train))
boxplot(depth~whichdata, depth.data, xlab="data", ylab="depth quantiles")

```

**Description**

Calculates case-specific iVarPro gradients on a new feature matrix by reusing the rule-level gradients stored in a previously computed iVarPro object and recomputing rule membership for the new cases.

**Usage**

```
## S3 method for class 'ivarpro'
predict(object,
  newdata = NULL,
  model = NULL,
  noise.na = NULL,
  path.store.membership = FALSE,
  save.data = TRUE,
  ...)
```

**Arguments**

object	An object returned by <code>ivarpro</code> (a numeric <code>data.frame</code> for univariate outcomes, or a list of such <code>data.frames</code> for multivariate outcomes). The object must contain the <code>"ivarpro.path"</code> attribute with stored rule metadata and rule-level gradients.
newdata	Optional <code>data.frame</code> of predictors for which to compute case-specific gradients. If supplied, gradients are computed for <code>newdata</code> using full membership (all eligible trees). If <code>NULL</code> (default), the function attempts a restore prediction for the training data, using OOB membership so that the returned gradients match the original <code>ivarpro()</code> output when possible.
model	Optional model override. If <code>NULL</code> (default), the model is taken from <code>attr(object, "model")</code> . The stored model can be either a <code>varpro</code> object or a <code>randomForestSRC</code> <code>rfsrc</code> grow object.
noise.na	Logical controlling how cells with no usable contributing rules are handled. If <code>NULL</code> (default), inherits the setting stored in <code>attr(object, "ivarpro.path")\$noise.na</code> . If <code>TRUE</code> , such cells are set to <code>NA</code> ; if <code>FALSE</code> , they are set to zero.
path.store.membership	Logical. If <code>TRUE</code> , store the rule membership indices (case IDs) for the prediction in <code>attr(out, "ivarpro.path")\$oobMembership</code> . This enables ladder-based bands in <code>plot.ivarpro</code> and summaries via <code>ivarpro_band()</code> for the predicted gradients, at the cost of additional memory. Default is <code>FALSE</code> .
save.data	Logical. If <code>TRUE</code> (default) and <code>newdata</code> is supplied, save the predictor data used for prediction as <code>attr(out, "data")</code> , enabling downstream plotting functions (e.g., <code>plot.ivarpro</code> ) to retrieve data automatically.
...	Additional arguments passed to <code>randomForestSRC::predict.rfsrc()</code> .

**Details**

A previously computed `iVarPro` object contains (i) a set of retained `VarPro` rules identified by tree and node/branch IDs, (ii) a release-variable index for each rule, and (iii) an estimated rule-level gradient. For new cases, the function obtains terminal node membership for each tree using `randomForestSRC::predict.rfsrc(membership = TRUE)` and determines which stored rules apply to each case by matching the case's terminal node to the stored rule node/branch ID for that tree. Case-specific gradients are then computed by aggregating the stored rule-level gradients over all rules that apply to the case and release the corresponding variable.

When `newdata` is not supplied, `predict()` attempts to return OOB case-specific gradients for the original training data. When the `iVarPro` object was created with membership information stored in its path (`path.store.membership = TRUE`, the default for `ivarpro()`), the restore prediction will reproduce the original `iVarPro` matrix.

If `save.data = TRUE` and `newdata` is supplied, `attr(out, "data")` contains the predictor matrix used for prediction, so downstream wrappers such as `plot.ivarpro` can be used directly. If you want ladder bands on predicted gradients, set `path.store.membership = TRUE` when predicting.

## Value

Returns an object with the same shape as object:

- For univariate outcomes, a numeric `data.frame` of dimension  $n_{test} \times p$  containing predicted case-specific gradients.
- For multivariate outcomes, a named list of such `data.frames`, one per outcome coordinate.

The returned object includes an `"ivarpro.path"` attribute containing rule metadata and (optionally) prediction-time membership indices.

## See Also

[ivarpro](#), [plot.ivarpro](#), [varpro](#)

## Examples

```
## -----
##
## Restore mode example (training OOB gradients are reproduced)
##
## -----

data(peakVO2, package = "randomForestSRC")

## Fit VarPro and iVarPro
vp <- varpro(Surv(ttodead, died) ~ ., peakVO2, ntree = 50)
imp <- ivarpro(vp)

## Restore prediction: should match the original iVarPro matrix
p.restore <- predict(imp)
all.equal(p.restore, imp, check.attributes = FALSE)

## Downstream plotting using the restored gradients
plot(p.restore, var = "peak.vo2",
      col.var = "interval", size.var = "y",
      data = attr(imp, "data"))

## -----
##
## Synthetic example (Friedman #1) with prediction on new data
```

```

##
## -----
if (requireNamespace("mlbench", quietly = TRUE)) {

  set.seed(123)

  ## training data
  tr <- mlbench::mlbench.friedman1(500, sd = 1)
  train <- data.frame(tr$x, y = tr$y)
  colnames(train)[1:ncol(tr$x)] <- paste0("x", 1:ncol(tr$x))

  ## ivarpro fit on training data
  vp <- varpro(y ~ ., train, ntree = 100)
  imp <- ivarpro(vp)

  ## test data
  te <- mlbench::mlbench.friedman1(1e4, sd = 1)
  test <- data.frame(te$x)
  colnames(test) <- paste0("x", 1:ncol(te$x))

  ## predicted gradients on test data
  p.test <- predict(imp, newdata = test)

  ## partial plot directly on predicted object
  plot(p.test, var = "x1", col.var = "x2")

  ## push x1 outside the original Friedman support [0, 1]
  ## a heuristic way to test out-of-distribution (OOD)
  test.ood <- test
  test.ood$x1 <- runif(nrow(test), min = -0.5, max = 1.5)

  ## predicted gradients on OOD test data
  p.test.ood <- predict(imp, newdata = test.ood)

  ## partial plot showing support for x1
  plot(p.test.ood, var = "x1", col.var = "x2",
       x.dist = c("density", "rug"))
  ## reference lines for the original training support of x1
  abline(v = c(0, 1), lty = 2)

}

```

---

predict.uvarpro

*Prediction on Test Data using Unsupervised VarPro*


---

### Description

Obtain predicted values on test data for unsupervised forests.

**Usage**

```
## S3 method for class 'uvarpro'
predict(object, newdata, ...)
```

**Arguments**

object	Unsupervised VarPro object from a previous call to <code>uvarpro</code> . Only applies if <code>method = "auto"</code> was used.
newdata	Optional test data. If not provided, the training data is used.
...	Additional arguments passed to internal methods.

**Details**

Applies to unsupervised VarPro objects built using the autoencoder (`method = "auto"`). The object contains a multivariate random forest used to generate predictions for the test data.

**Value**

Returns a matrix of predicted values, where each column corresponds to a feature (with one-hot encoding applied). The result includes the following attributes:

1. `mse`: Standardized mean squared error averaged across features.
2. `mse.all`: Standardized mean squared error for each individual feature.

**Author(s)**

Min Lu and Hemant Ishwaran

**See Also**

[uvarpro](#)

**Examples**

```
## -----
##
## boston housing
## obtain predicted values for the training data
##
## -----

## unsupervised varpro on boston housing
data(BostonHousing, package = "mlbench")
o <- uvarpro(data=BostonHousing)

## predicted values for the training features
print(head(predict(o)))

## -----
##
```

```

## mtcars
## obtain predicted values for test data
## also illustrates hot-encoding working on test data
##
## -----

## mtcars with some factors
d <- data.frame(mpg=mtcars$mpg,lapply(mtcars[, c("cyl", "vs", "carb")], as.factor))

## training
o <- uvarpro(d[1:20,])

## predicted values on test data
print(predict(o, d[-(1:20),]))

## predicted values on bad test data with strange factor values
dbad <- d[-(1:20),]
dbad$carb <- as.character(dbad$carb)
dbad$carb <- sample(LETTERS, size = nrow(dbad))
print(predict(o, dbad))

```

---

predict.varpro

*Prediction on Test Data using VarPro*


---

## Description

Obtain predicted values on test data for VarPro object.

## Usage

```

## S3 method for class 'varpro'
predict(object, newdata, ...)

```

## Arguments

object	VarPro object returned from a previous call to varpro.
newdata	Optional test data. If not provided, predictions are computed using the training data (out-of-bag).
...	Additional arguments passed to internal methods.

## Details

VarPro uses rules extracted from a random forest built using guided tree-splitting, where variables are selected based on split-weights computed in a preprocessing step.

**Value**

Returns predicted values for the input data. If `newdata` is provided, predictions are made on that data; otherwise, out-of-bag predictions for the training data are returned.

**Author(s)**

Min Lu and Hemant Ishwaran

**References**

Lu, M. and Ishwaran, H. (2024). Model-independent variable selection via the rule-based variable priority. arXiv e-prints, pp.arXiv-2409.

**See Also**

[varpro](#)

**Examples**

```
## -----
## toy example - needed to pass CRAN test
## -----

## train call
o <- varpro(mpg~., mtcars[1:20,], ntree = 1)

## predict call
print(predict(o, mtcars[-(1:20),]))

## -----
##
## boston housing regression
## obtain predicted values for the training data
##
## -----

## varpro applied to boston housing data
data(BostonHousing, package = "mlbench")
o <- varpro(medv~., BostonHousing)

## predicted values for the training features
print(head(predict(o)))

## -----
##
## iris classification
## obtain predicted values for test data
##
## -----
```

```

## varpro applied to iris data
trn <- sample(1:nrow(iris), size = 100, replace = FALSE)
o <- varpro(Species~., iris[trn,])

## predicted values on test data
print(data.frame(Species=iris[-trn, "Species"], predict(o, iris[-trn,])))

## -----
##
## mtcars regression: illustration of hot-encoding on test data
##
## -----

## mtcars with some factors
d <- data.frame(mpg=mtcars$mpg,lapply(mtcars[, c("cyl", "vs", "carb")], as.factor))

## varpro on training data
o <- varpro(mpg~., d[1:20,])

## predicted values on test data
print(predict(o, d[-(1:20),]))

## predicted values on bad test data with strange factor values
dbad <- d[-(1:20),]
dbad$carb <- as.character(dbad$carb)
dbad$carb <- sample(LETTERS, size = nrow(dbad))
print(predict(o, dbad))

```

---

uvarpro

*Unsupervised Variable Selection using Variable Priority (UVarPro)*


---

## Description

Performs unsupervised variable selection by extending the VarPro framework to forests grown without labels. UVarPro identifies features that explain structure in the data through region-release contrasts, with importance assessed using entropy or lasso-based methods.

## Usage

```

uvarpro(data,
  method = c("auto", "unsupv", "rnd"),
  ntree = 200, nodesize = NULL,
  max.rules.tree = 20, max.tree = 200,
  verbose = FALSE, seed = NULL,
  ...)

```

**Arguments**

<code>data</code>	Data frame containing the unsupervised data.
<code>method</code>	Type of forest used. Options are "auto" (auto-encoder), "unsupv" (unsupervised analysis), and "rnd" (pure random forest).
<code>ntree</code>	Number of trees to grow.
<code>nodesize</code>	Minimum terminal node size. If not specified, an internal function selects an appropriate value based on sample size and dimension.
<code>max.rules.tree</code>	Maximum number of rules per tree.
<code>max.tree</code>	Maximum number of trees used to extract rules.
<code>verbose</code>	Print verbose output?
<code>seed</code>	Seed for reproducibility.
<code>...</code>	Additional arguments passed to <code>rfsrc</code> .

**Details**

UVarPro performs unsupervised variable selection by applying the VarPro framework to random forests trained on unlabeled data (Zhou et al., 2025). The procedure has two components: (i) the construction of an unsupervised forest, and (ii) the evaluation of variable importance based on region-release contrasts, in direct analogy to the supervised setting in VarPro.

The forest construction is controlled by the `method` argument. By default, `method = "auto"` fits a random forest autoencoder, which regresses each selected variable on itself, a specialized form of multivariate forest modeling. Alternatives include "unsupv", which uses pseudo-responses and multivariate splits to build an unsupervised forest (Tang and Ishwaran, 2017), and "rnd", which uses completely random splits. For large datasets, the autoencoder may be slower, while "unsupv" and "rnd" are often much faster.

Variable importance is assessed using region-release contrasts formed by the forest. By default, the importance function returns an entropy-based criterion. This measure compares the variability within each region to the variability across the combined sample, effectively acting like a ratio of between to within sums of squares. Importance values are averaged across many region-release rules, providing a rough but fast estimate of how strongly a feature contributes to distinguishing regions. See examples below.

In addition to this default entropy measure, UVarPro supports custom user-defined entropy functions to create alternative importance metrics.

A more sophisticated procedure, described in Zhou et al. (2026), reframes each region-release contrast as a supervised classification task, with membership in the region serving as the class label. Variable effects are estimated using lasso-based logistic regression, and coefficients are aggregated over the collection of region-release tasks. This produces a sparser and often more interpretable assessment of importance compared to the entropy method. Although more computationally intensive, the lasso-driven approach can provide sharper separation of relevant and irrelevant features. See examples below for details.

**Value**

A `uvarpro` object.

**Author(s)**

Min Lu and Hemant Ishwaran

**References**

Tang F. and Ishwaran H. (2017). Random forest missing data algorithms. *Statistical Analysis and Data Mining*, 10:363-377.

Zhou L., Lu M. and Ishwaran H. (2026). Variable priority for unsupervised variable selection. *Pattern Recognition*, 172:112727.

**See Also**

[varpro](#)

**Examples**

```
## -----  
## toy example - needed to pass CRAN test  
## -----  
  
## mtcars unsupervised regression  
o <- uvarpro(mtcars, ntree = 1)  
  
## -----  
## boston housing: default call  
## -----  
  
data(BostonHousing, package = "mlbench")  
  
## default call  
o <- uvarpro(BostonHousing)  
print(importance(o))  
  
## -----  
## boston housing: using method="unsupv"  
## -----  
  
data(BostonHousing, package = "mlbench")  
  
## unsupervised splitting  
o <- uvarpro(BostonHousing, method = "unsupv")  
print(importance(o))  
  
## -----  
## boston housing: illustrates hot-encoding  
## -----  
  
## load the data
```

```

data(BostonHousing, package = "mlbench")

## convert some of the features to factors
Boston <- BostonHousing
Boston$zn <- factor(Boston$zn)
Boston$chas <- factor(Boston$chas)
Boston$lstat <- factor(round(0.2 * Boston$lstat))
Boston$nox <- factor(round(20 * Boston$nox))
Boston$rm <- factor(round(Boston$rm))

## call unsupervised varpro and print importance
print(importance(o <- uvarpro(Boston)))

## get top variables
get.topvars(o)

## map importance values back to original features
print(get.orgvimp(o))

## same as above ... but for all variables
print(get.orgvimp(o, pretty = FALSE))

## -----
## latent variable simulation
## -----

n <- 1000
w <- rnorm(n)
x <- rnorm(n)
y <- rnorm(n)
z <- rnorm(n)
ei <- matrix(rnorm(n * 20, sd = sqrt(.1)), ncol = 20)
e21 <- rnorm(n, sd = sqrt(.4))
e22 <- rnorm(n, sd = sqrt(.4))
wi <- w + ei[, 1:5]
xi <- x + ei[, 6:10]
yi <- y + ei[, 11:15]
zi <- z + ei[, 16:20]
h1 <- w + x + e21
h2 <- y + z + e22
dta <- data.frame(w=w,wi=wi,x=x,xi=xi,y=y,yi=yi,z=z,zi=zi,h1=h1,h2=h2)

## default call
print(importance(uvarpro(dta)))

## -----
## glass (remove outcome)
## -----

data(Glass, package = "mlbench")

```

```

## remove the outcome
Glass$Type <- NULL

## get importance
o <- uvarpro(Glass)
print(importance(o))

## compare to PCA
(biplot(prcomp(o$x, scale = TRUE)))

## -----
## iowa housing - illustrates lasso importance
## -----

## first we roughly impute the data
data(housing, package = "randomForestSRC")

## to speed up analysis, convert all factors to real values
iowa <- roughfix(housing)
iowa <- data.frame(data.matrix(iowa))

## canonical call
o <- uvarpro(iowa)

## standard importance
print(importance(o))

## lasso importance
beta <- get.beta.entropy(o)
print(beta)
print(sort(colMeans(beta, na.rm=TRUE), decreasing = TRUE))

## s-dependent graph
sdependent(beta)

## lasso importance without pre-filtering
## beta.nof <- get.beta.entropy(o, pre.filter = FALSE)
## print(beta.nof)
## print(sort(colMeans(beta.nof, na.rm=TRUE), decreasing = TRUE))

## lasso importance with second stage sparsity lasso
## beta.sparse <- get.beta.entropy(o, second.stage = TRUE)
## print(beta.sparse)

## -----
## custom importance
## OPTION 1: use hidden entropy option
## -----

my.entropy <- function(xC, x0, ...) {

  ## xC      x feature data from complementary region

```

```

## x0      x feature data from original region
## ...    used to pass additional options (required)

## custom importance value
wss <- mean(apply(rbind(x0, xC), 2, sd, na.rm = TRUE))
bss <- (mean(apply(xC, 2, sd, na.rm = TRUE)) +
        mean(apply(x0, 2, sd, na.rm = TRUE)))
imp <- 0.5 * bss / wss

## entropy value must contain complementary and original membership
entropy <- list(comp = list(...)$compMembership,
               oob = list(...)$oobMembership)

## return importance and in the second slot the entropy list
list(imp = imp, entropy)

o <- uvarpro(BostonHousing, entropy=my.entropy)
print(importance(o))

## -----
## custom importance
## OPTION 2: direct importance without hidden entropy option
## -----

o <- uvarpro(BostonHousing, ntree=3, max.rules.tree=10)

## convert original/release region into two-class problem
## define importance as the lasso beta values

## For faster performance on Unix systems, consider using:
## library(parallel)
## imp <- do.call(rbind, mclapply(seq_along(o$entropy), function(j) { ... }))

imp <- do.call(rbind, lapply(seq_along(o$entropy), function(j) {
  r0 <- do.call(rbind, lapply(o$entropy[[j]], function(r) {
    xC <- o$x[[r][[1]],names(o$entropy),drop=FALSE]
    x0 <- o$x[[r][[2]],names(o$entropy),drop=FALSE]
    y <- factor(c(rep(0, nrow(xC)), rep(1, nrow(x0))))
    x <- rbind(xC, x0)
    x <- x[, colnames(x) != names(o$entropy)[j]]
    fit <- tryCatch(
      suppressWarnings(glmnet::cv.glmnet(as.matrix(x), y, family = "binomial")),
      error = function(e) NULL
    )
    if (!is.null(fit)) {
      beta <- setNames(rep(0, length(o$entropy)), names(o$entropy))
      bhat <- abs(coef(fit)[-1, 1])
      beta[names(bhat)] <- bhat
      beta
    } else {
      NULL
    }
  })
})

```

```

    )))
  if (!is.null(r0)) {
    val <- colMeans(r0, na.rm = TRUE)
    names(val) <- colnames(r0)
    return(val)
  } else {
    return(NULL)
  }
}) |> setNames(names(o$entropy)))

print(imp)

## -----
## custom importance
## OPTION 3: direct importance using built in lasso beta function
## -----

o <- uvarpro(BostonHousing)
print((get.beta.entropy(o)))

}

```

---

varpro	<i>Model-Independent Variable Selection via the Rule-based Variable Priority (VarPro)</i>
--------	---

---

## Description

Model-Independent Variable Selection via the Rule-based Variable Priority (VarPro) for Regression, Multivariate Regression, Classification and Survival.

## Usage

```

varpro(formula, data, nvar = 30, ntree = 500,
       split.weight = TRUE, split.weight.method = NULL, sparse = TRUE,
       nodesize = NULL, max.rules.tree = 150, max.tree = min(150, ntree),
       parallel = TRUE, cores = get.mc.cores(), verbose = FALSE, seed = NULL, ...)

```

## Arguments

formula	Formula specifying the model to be fit.
data	Data frame containing the training data.
nvar	Maximum number of variables to return.
ntree	Number of trees to grow.

<code>split.weight</code>	Use guided tree-splitting? Candidate variables for splitting are selected with probability proportional to a split-weight, obtained by default from a preliminary lasso+tree step.
<code>split.weight.method</code>	Character string (or vector) specifying method used to generate split-weights. Defaults to lasso+tree. See Details.
<code>sparse</code>	Use sparse split-weights?
<code>nodesize</code>	Minimum terminal node size. If not specified, value is set internally based on sample size and dimension.
<code>max.rules.tree</code>	Maximum number of rules per tree.
<code>max.tree</code>	Maximum number of trees used to extract rules.
<code>parallel</code>	Use parallel execution for lasso folds using <b>doMC</b> .
<code>cores</code>	Number of cores for parallel processing. Defaults to <code>parallel::detectCores()</code> .
<code>verbose</code>	Print verbose output?
<code>seed</code>	Seed for reproducibility.
<code>...</code>	Additional arguments for advanced use.

## Details

Rule-based models, such as decision rules, rule learning, trees, boosted trees, Bayesian additive regression trees, Bayesian forests, and random forests, are widely used for variable selection. These nonparametric methods require no model specification and accommodate various outcomes including regression, classification, survival, and longitudinal data.

Although permutation variable importance (VIMP) and knockoff methods have been extensively studied, their effectiveness can be limited in practice. Both approaches rely on the quality of artificially generated covariates, which may not perform well in complex or high-dimensional settings.

To address these limitations, we introduce a new framework called variable priority (VarPro). Instead of generating synthetic covariates, VarPro constructs *release rules* to assess the impact of each covariate on the response. Neighborhoods of existing data are used for estimation, avoiding the need for artificial data generation. Like VIMP and knockoffs, VarPro imposes no assumptions on the conditional distribution of the response.

The VarPro algorithm proceeds as follows: A forest of `ntree` trees is grown using guided tree-splitting, where candidate variables for node splitting are selected with probability proportional to their split-weights. These split-weights are computed in a preprocessing step. A subset of `max.tree` trees is randomly selected from the forest, and `max.rules.tree` branches are sampled from each selected tree. The resulting rules form the basis of the VarPro importance estimator. The method supports regression, multivariate regression, multiclass classification, and survival outcomes.

Guided tree-splitting encourages rule construction to favor influential features. Thus, `split.weight` should generally remain TRUE, especially for high-dimensional problems. If disabled, it is recommended to increase `nodesize` to improve estimator precision.

By default, split-weights are computed via a lasso-plus-tree strategy. Specifically, the split-weight of a variable is defined as the product of the absolute standardized lasso coefficient and the variable's split frequency from a forest of shallow trees. If sample size and dimension are both moderate, this

may be replaced by the variable's absolute permutation importance. Note: variables are one-hot encoded for use in lasso, and all inputs are converted to numeric values.

To customize split-weight construction, use the `split.weight.method` argument with one or more of the following strings: "lasso", "tree", or "vimp". For example, "lasso" uses only lasso coefficients; "lasso tree" combines lasso and shallow trees; "lasso vimp" combines lasso with permutation importance. See examples below.

Variables are ranked by importance, with higher values indicating greater influence. Cross-validation can be used to determine a cutoff threshold. See `cv.varpro` for details.

Run time can be reduced by using smaller values of `ntree` or larger values of `nodesize`. Additional runtime tuning options are discussed in the examples.

In class-imbalanced two-class settings, the algorithm automatically switches to random forest quantile classification (RFQ; see O'Brien and Ishwaran, 2019) using the geometric mean (`gmean`) metric. This behavior can be overridden via the hidden option `use.rfq`.

## Value

Output containing VarPro estimators used to calculate importance. See `importance.varpro`. Also see `cv.varpro` for automated variable selection.

## Author(s)

Min Lu and Hemant Ishwaran

## References

Ishwaran H. (2025). *Multivariate Statistics: Classical Foundations and Modern Machine Learning*. Chapman and Hall/CRC.

Ishwaran H. and Blackstone E.H. (2025). Harnessing the power of virtual (digital) twins: Graphical causal tools for understanding patient and hospital differences. *Computational and Structural Biotechnology Journal*, 28:312.

Lu, M. and Ishwaran, H. (2024). Model-independent variable selection via the rule-based variable priority. arXiv e-prints, pp.arXiv-2409.

Lu, M. and Ishwaran, H. (2025). Individual variable priority: a model-independent local gradient method for variable importance. *Artificial Intelligence Review*, 58:407.

Zhou L., Lu M. and Ishwaran H. (2026). Variable priority for unsupervised variable selection. *Pattern Recognition*, 172:112727.

## See Also

[alzheimers](#) [cv.varpro](#) [glioma](#) [hrrecov](#) [importance.varpro](#) [partialpro](#) [predict.varpro](#) [isopro](#)  
[ivarpro](#) [outpro](#) [uvarpro](#)

## Examples

```
## -----
## toy example - needed to pass CRAN test
## -----
```

```
## mtcars regression
o <- varpro(mpg ~ ., mtcars, ntree = 1)

## -----
## classification example: iris
## -----

## apply varpro to the iris data
o <- varpro(Species ~ ., iris, max.tree = 5)

## call the importance function and print the results
print(importance(o))

## -----
## regression example: boston housing
## -----

## load the data
data(BostonHousing, package = "mlbench")

## call varpro
o <- varpro(medv~., BostonHousing)

## extract and print importance values
imp <- importance(o)
print(imp)

## another way to extract and print importance values
print(get.vimp(o))
print(get.vimp(o, pretty = FALSE))

## plot importance values
importance(o, plot.it = TRUE)

## -----
## regression example: boston housing illustrating hot-encoding
## -----

## load the data
data(BostonHousing, package = "mlbench")

## convert some of the features to factors
Boston <- BostonHousing
Boston$zn <- factor(Boston$zn)
Boston$chas <- factor(Boston$chas)
Boston$lstat <- factor(round(0.2 * Boston$lstat))
Boston$nox <- factor(round(20 * Boston$nox))
Boston$rm <- factor(round(Boston$rm))

## call varpro and print the importance
```

```

print(importance(o <- varpro(medv~., Boston)))

## get top variables
get.topvars(o)

## map importance values back to original features
print(get.orgvimp(o))

## same as above ... but for all variables
print(get.orgvimp(o, pretty = FALSE))

## -----
## regression example: friedman 1
## -----

o <- varpro(y~., data.frame(mlbench::mlbench.friedman1(1000)))
print(importance(o))

## -----
## example without guided tree-splitting
## -----

o <- varpro(y~., data.frame(mlbench::mlbench.friedman2(1000)),
            nodesize = 10, split.weight = FALSE)
print(importance(o))

## -----
## regression example: all noise
## -----

x <- matrix(rnorm(100 * 50), 100, 50)
y <- rnorm(100)
o <- varpro(y~., data.frame(y = y, x = x))
print(importance(o))

## -----
## multivariate regression example: boston housing
## -----

data(BostonHousing, package = "mlbench")

## using rfsrc multivariate formula call
importance(varpro(Multivar(lstat, nox) ~., BostonHousing))

## using cbind multivariate formula call
importance(varpro(cbind(lstat, nox) ~., BostonHousing))

##-----
## class imbalanced problem
##
## - simulation example using the caret R-package
## - creates imbalanced data by randomly sampling the class 1 values
##

```

```

##-----
if (library("caret", logical.return = TRUE)) {

  ## experimental settings
  n <- 5000
  q <- 20
  ir <- 6
  f <- as.formula(Class ~ .)

  ## simulate the data, create minority class data
  d <- twoClassSim(n, linearVars = 15, noiseVars = q)
  d$Class <- factor(as.numeric(d$Class) - 1)
  idx.0 <- which(d$Class == 0)
  idx.1 <- sample(which(d$Class == 1), sum(d$Class == 1) / ir , replace = FALSE)
  d <- d[c(idx.0,idx.1),, drop = FALSE]
  d <- d[sample(1:nrow(d)), ]

  ## varpro call
  print(importance(varpro(f, d)))

}

## -----
## multiclass example: wine data
## -----
data(wine, package = "randomForestSRC")
wine$quality <- factor(wine$quality)
o <- varpro(quality~., wine)
print(importance(o))

## -----
## survival example: pbc
## -----
data(pbc, package = "randomForestSRC")
o <- varpro(Surv(days, status)~., pbc)
print(importance(o))

## -----
## pbc survival with rmst (restricted mean survival time)
## functional of interest is RMST at 500 days
## -----
data(pbc, package = "randomForestSRC")
o <- varpro(Surv(days, status)~., pbc, rmst = 500)
print(importance(o))

## -----
## pbc survival with rmst vector
## variable importance is a list for each rmst value
## -----
data(pbc, package = "randomForestSRC")
o <- varpro(Surv(days, status)~., pbc, rmst = c(500, 1000))
print(importance(o))

```

```

## -----
## survival example with more variables
## -----
data(peakV02, package = "randomForestSRC")
o <- varpro(Surv(ttodead, died)~., peakV02)
imp <- importance(o, plot.it = TRUE)
print(imp)

## -----
## high dimensional survival example
## -----
data(vdv, package = "randomForestSRC")
o <- varpro(Surv(Time, Censoring)~., vdv)
print(importance(o))

## -----
## high dimensional survival example without sparse option
## -----
data(vdv, package = "randomForestSRC")
o <- varpro(Surv(Time, Censoring)~., vdv, sparse = FALSE)
print(importance(o))

## -----
## high dimensional survival example using different split-weight methods
## -----
data(vdv, package = "randomForestSRC")
f <- as.formula(Surv(Time, Censoring)~.)

## lasso only
print(importance(varpro(f, vdv, split.weight.method = "lasso")))

## lasso and vimp
print(importance(varpro(f, vdv, split.weight.method = "lasso vimp")))

## lasso, vimp and shallow trees
print(importance(varpro(f, vdv, split.weight.method = "lasso vimp tree")))

## -----
## largish data (iowa housing data)
## to speed up calculations convert data to all real
## -----

## first we roughly impute the data
data(housing, package = "randomForestSRC")
dta <- roughfix(housing)
dta <- data.frame(data.matrix(dta))

## varpro call
o <- varpro(SalePrice~., dta)
print(importance(o))

## -----

```

```

## large data: illustrates different ways to improve speed
## -----

n <- 25000
p <- 50
d <- data.frame(y = rnorm(n), x = matrix(rnorm(n * p), n))

## use large nodesize
print(system.time(o <- varpro(y~., d, ntree = 100, nodesize = 200)))
print(importance(o))

## use large nodesize, smaller bootstrap
print(system.time(o <- varpro(y~., d, ntree = 100, nodesize = 200,
                             samplesize = 100)))
print(importance(o))

## -----
## custom split-weights (hidden option)
## -----

## load the data
data(BostonHousing, package = "mlbench")

## make some features into factors
Boston <- BostonHousing
Boston$zn <- factor(Boston$zn)
Boston$chas <- factor(Boston$chas)
Boston$lstat <- factor(round(0.2 * Boston$lstat))
Boston$nox <- factor(round(20 * Boston$nox))
Boston$rm <- factor(round(Boston$rm))

## get default custom split-weights: a named real vector
swt <- get.splitweight.custom(medv~., Boston)

## define custom splits weight
swt <- swt[grepl("crim", names(swt)) |
          grepl("zn", names(swt)) |
          grepl("nox", names(swt)) |
          grepl("rm", names(swt)) |
          grepl("lstat", names(swt))]

swt[grepl("nox", names(swt))] <- 4
swt[grepl("lstat", names(swt))] <- 4

swt <- c(swt, strange=99999)

cat("custom split-weight\n")
print(swt)

## call varpro with the custom split-weights
o <- varpro(medv~., Boston, split.weight.custom=swt, verbose=TRUE, sparse=FALSE)
cat("varpro result\n")

```

```
print(importance(o))
print(get.vimp(o, pretty=FALSE))
print(get.orgvimp(o, pretty=FALSE))
```

---

varpro.news	<i>Show the NEWS file</i>
-------------	---------------------------

---

### Description

Show the NEWS file of the **varPro** package.

### Usage

```
varpro.news(...)
```

### Arguments

... Further arguments passed to or from other methods.

### Value

None.

### Author(s)

Min Lu and Hemant Ishwaran

---

varpro.strength	<i>Obtain Strength Array and Other Values from a VarPro Object</i>
-----------------	--

---

### Description

Used to parse values from a VarPro object.

### Usage

```
varpro.strength(object,
  newdata,
  m.target = NULL,
  max.rules.tree = 150,
  max.tree = 150,
  stat = c("importance", "complement", "oob", "none"),
  membership = FALSE,
  neighbor = 5,
  seed = NULL,
  do.trace = FALSE, ...)
```

**Arguments**

object	rfsrc object.
newdata	Optional test data. If provided, returns branch and complementary branch membership of the training data corresponding to the test cases.
m.target	Character string specifying the target outcome for multivariate families. If unspecified, a default is selected automatically.
max.rules.tree	Maximum number of rules extracted per tree.
max.tree	Maximum number of trees used for rule extraction.
stat	Statistic to output. Options include "importance", "complement mean", and "oob mean".
membership	Return out-of-bag and complementary membership indices for each rule?
neighbor	Nearest neighbor parameter, used only when newdata is specified.
seed	Seed for reproducibility.
do.trace	Enable detailed trace output.
...	Additional arguments.

**Details**

Not intended for direct end-user use; primarily designed for internal package operations.

**Value**

Object coerced so as to work with other functions in the package.

**Examples**

```
## -----
## regression example: boston housing
## -----

## load the data
data(BostonHousing, package = "mlbench")

o <- randomForestSRC::rfsrc(medv~., BostonHousing, ntree=100)

## call varpro.strength
varpro.strength(object = o, max.rules.tree = 10, max.tree = 15)

## call varpro.strength with test data
varpro.strength(object = o, newdata = BostonHousing[1:3,], max.rules.tree = 10, max.tree = 15)
```

# Index

- \* **cv.varpro**
    - cv.varpro, 4
  - \* **datasets**
    - alzheimers, 2
    - gliomas, 8
    - hrrecov, 9
  - \* **documentation**
    - varpro.news, 69
  - \* **hplot**
    - partial.ivarpro, 31
    - plot.ivarpro, 39
  - \* **importance**
    - partial.ivarpro, 31
    - plot.ivarpro, 39
  - \* **individual importance**
    - ivarpro, 22
    - predict.ivarpro, 48
  - \* **outlier**
    - isopro, 18
  - \* **plot**
    - partialpro, 35
    - plot.partialpro, 43
  - \* **predict outlier**
    - predict.isopro, 46
  - \* **predict uvarpro**
    - predict.uvarpro, 51
  - \* **predict varpro**
    - predict.varpro, 53
  - \* **uvarpro**
    - uvarpro, 55
  - \* **varpro.strength**
    - varpro.strength, 69
  - \* **varpro**
    - importance, 12
    - importance.varpro, 15
    - varpro, 61
- alzheimers, 2, 63
- cv.varpro, 4, 13, 16, 63
- glioma, 63
- glioma (gliomas), 8
- gliomas, 8
- hrrecov, 9, 63
- importance, 12
- importance (importance.varpro), 15
- importance.varpro, 5, 15, 63
- isopro, 18, 37, 47, 63
- ivarpro, 22, 32, 33, 40, 42, 49, 50, 63
- outpro, 28, 63
- partial.ivarpro, 31
- partialpro, 35, 45, 63
- plot.ivarpro, 26, 39, 49, 50
- plot.partialpro, 43
- predict.isopro, 20, 46
- predict.ivarpro, 40, 48
- predict.uvarpro, 51
- predict.varpro, 53, 63
- rfsrc, 31
- shap.ivarpro, 33, 42
- uvarpro, 5, 13, 20, 47, 52, 55, 63
- varpro, 5, 13, 16, 20, 26, 31, 33, 37, 42, 47, 50, 54, 57, 61
- varpro.news, 69
- varpro.strength, 69