

CMA

October 5, 2010

CMA-package

Synthesis of microarray-based classification

Description

The aim of the package is to provide a user-friendly environment for the evaluation of classification methods using gene expression data. A strong focus is on combined variable selection, hyperparameter tuning, evaluation, visualization and comparison of (up to now) 21 classification methods from three main fields: Discriminant Analysis, Neural Networks and Machine Learning. Although the package has been created with the intention to be used for Microarray data, it can as well be used in various ($p > n$)-scenarios.

Details

Package: CMA
Type: Package
Version: 1.3.3
Date: 2009-9-14
License: GPL (version 2 or later)

Most Important Steps for the workflow are:

1. Generate evaluation datasets using [GenerateLearningsets](#)
2. (Optionally): Perform variable selection using [GeneSelection](#)
3. (Optionally): Perform hyperparameter tuning using [tune](#)
4. Perform [classification](#) using 1.-3.
5. Repeat 2.-4. based on 1. for several methods: [compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)
6. Evaluate the results from 5. using [evaluation](#) and make a comparison by calling [compare](#)

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

Christoph Bernau <bernau@ibe.med.uni-muenchen.de>

Maintainer: Christoph Bernau <bernau@ibe.med.uni-muenchen.de>.

References

Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

ElasticNetCMA-methods

Classification and variable selection by the ElasticNet

Description

Zou and Hastie (2004) proposed a combined L1/L2 penalty for regularization and variable selection. The Elastic Net penalty encourages a grouping effect, where strongly correlated predictors tend to be in or out of the model together. The computation is done with the function `glmPath` from the package of the same name.

Methods

`X = "matrix", y = "numeric", f = "missing"` signature 1

`X = "matrix", y = "factor", f = "missing"` signature 2

`X = "data.frame", y = "missing", f = "formula"` signature 3

`X = "ExpressionSet", y = "character", f = "missing"` signature 4

For references, further argument and output information, consult [ElasticNetCMA](#)

ElasticNetCMA

Classification and variable selection by the ElasticNet

Description

Zou and Hastie (2004) proposed a combined L1/L2 penalty for regularization and variable selection. The Elastic Net penalty encourages a grouping effect, where strongly correlated predictors tend to be in or out of the model together. The computation is done with the function `glmPath` from the package of the same name.

The method can be used for variable selection alone, s. [GeneSelection](#).

For S4 method information, see `ElasticNetCMA-methods`.

Usage

```
ElasticNetCMA(X, y, f, learnind, norm.fraction = 0.1, alpha=0.5, models=FALSE, .
```

Arguments

<code>X</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>. note: by default, the predictors are scaled to have unit variance and zero mean. Can be changed by passing <code>standardize = FALSE</code> via the <code>...</code> argument.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • missing, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p>
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be missing; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>norm.fraction</code>	L1 Shrinkage intensity, expressed as the fraction of the coefficient L1 norm compared to the maximum possible L1 norm (corresponds to <code>fraction = 1</code>). Lower values correspond to higher shrinkage. Note that the default (0.1) need not produce good results, i.e. tuning of this parameter is recommended.
<code>alpha</code>	The elasticnet mixing parameter, with $0 < \alpha \leq 1$. The penalty is defined as $(1-\alpha)/2 \ \beta\ _2^2 + \alpha \ \beta\ _1$. <code>alpha=1</code> is the lasso penalty; Currently ' <code>alpha<0.01</code> ' not reliable, unless you supply your own <code>lambda</code> sequence
<code>models</code>	a logical value indicating whether the model object shall be returned
<code>...</code>	Further arguments passed to the function <code>glm</code> from the package of the same name.

Value

An object of class `clvareselectoutput`.

Note

For a strongly related method, s. `LassoCMA`.
Up to now, this method can only be applied to binary classification.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

Christoph Bernau <bernau@ibe.med.uni-muenchen.de>

References

- Zhou, H., Hastie, T. (2004).
Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society B*, 67(2),301-320
- Young-Park, M., Hastie, T. (2007)
L1-regularization path algorithm for generalized linear models.
Journal of the Royal Statistical Society B, 69(4), 659-677

See Also

[compBoostCMA](#), [dldaCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression
golubX <- as.matrix(golub[,-1])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run ElasticNet - penalized logistic regression (no tuning)
result <- ElasticNetCMA(X=golubX, y=golubY, learnind=learnind, norm.fraction = 0.2, alpha)
show(result)
ftable(result)
plot(result)
```

GeneSelection-methods

General method for variable selection with various methods

Description

Performs gene selection for the following signatures:

Methods

- X = "matrix", y = "numeric", f = "missing"** signature 1
- X = "matrix", y = "factor", f = "missing"** signature 2
- X = "data.frame", y = "missing", f = "formula"** signature 3
- X = "ExpressionSet", y = "character", f = "missing"** signature 4

For further argument and output information, consult [GeneSelection](#).

GeneSelection

*General method for variable selection with various methods***Description**

For different learning data sets as defined by the argument `learningsets`, this method ranks the genes from the most relevant to the less relevant using one of various 'filter' criteria or provides a sparse collection of variables (Lasso, ElasticNet, Boosting). The results are typically used for variable selection for the classification procedure that follows.

For S4 class information, s. `GeneSelection-methods`.

Usage

```
GeneSelection(X, y, f, learningsets, method = c("t.test", "welch.test", "wilcox.test"))
```

Arguments

- | | |
|---------------------------|--|
| <code>X</code> | Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>. |
| <code>y</code> | Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if <code>X</code> is an <code>ExpressionSet</code>. • missing, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. |
| <code>f</code> | A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables. |
| <code>learningsets</code> | An object of class <code>learningsets</code> . May be missing, then the complete datasets is used as learning set. |
| <code>method</code> | A character specifying the method to be used: <ul style="list-style-type: none"> <code>t.test</code> two-sample t.test (equal variances for both classes assumed). <code>welch.test</code> Welch modification of the t.test (unequal variances for both classes). <code>wilcox.test</code> Wilcoxon rank sum test. <code>f.test</code> F test belonging to the linear hypothesis that the mean is the same for all classes. Usually used for the multiclass scheme, is equivalent to <code>method = t.test</code> in the two-class case. <code>kruskal.test</code> Multi-class generalization of the Wilcoxon rank sum test and the nonparametric pendant to the F test, respectively. <code>limma</code> 'Moderated t' statistic for the two-class case and 'moderated F' statistic for the multiclass case, described in Smyth (2003). Requires the package <code>limma</code>. <code>rfe</code> One-step Recursive Feature Elimination, based on the Support Vector Machine. The method is described in Guyon et al. (2002). Requires the package <code>e1071</code>. Take care that appropriate hyperparameters are passed by the <code>...</code> argument. <code>rf</code> Random Forest Variable Importance Measure. Requires the package <code>randomForest</code> |

	<p><code>lasso</code> L1 penalized logistic regression leads to sparsity with respect to the variables used. Calls the function <code>LassoCMA</code>, which requires the package <code>glmPath</code>. warning: Take care that appropriate hyperparameters are passed by the <code>...</code> argument.</p> <p><code>elasticnet</code> Penalized logistic regression with both L1 and L2 penalty, claimed by Zhou and Hastie (2004) to select 'variable groups'. Calls the function <code>ElasticNetCMA</code>, which requires the package <code>glmPath</code>. warning: Take care that appropriate hyperparameters are passed by the <code>...</code> argument.</p> <p><code>boosting</code> Componentwise boosting (Buehlmann and Yu, 2003) has been shown to mimic the LASSO (Efron et al., 2004; Buehlmann and Yu, 2006). Calls the function <code>compBoostCMA</code>. Take care that appropriate hyperparameters are passed by the <code>...</code> argument.</p> <p><code>golub</code> The (theoretically unfounded) variable selection criterion used by Golub et al. (1999), s. <code>golub</code>.</p>
<code>scheme</code>	The scheme to be used in the case of a non-binary response. Must be one of "pairwise", "one-vs-all" or "multiclass". The last case only makes sense if <code>method</code> is one of <code>f.test</code> , <code>limma</code> , <code>rf</code> , <code>boosting</code> , which can directly be applied to the multi class case.
<code>trace</code>	Should the progress be traced ? Default is TRUE.
<code>...</code>	Further arguments passed to the function performing variable selection, s. <code>method</code> .

Value

An object of class `genesel`.

Note

most of the methods described above are only apt for the binary classification case. The only ones that can be used without restriction in the multiclass case are

- `f.test`
- `kruskal.test`
- `rf`
- `boosting`

For the rest, pairwise or one-vs-all schemes are used.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

Christoph Bernau <bernau@ibe.med.uni-muenchen.de>

References

Smyth, G. K., Yang, Y.-H., Speed, T. P. (2003).
Statistical issues in microarray data analysis.
Methods in Molecular Biology 224, 111-136.

Guyon, I., Weston, J., Barnhill, S., Vapnik, V. (2002).
Gene Selection for Cancer Classification using support vector machines. *Journal of Machine Learning Research*, 46, 389-422

Zhou, H., Hastie, T. (2004).
Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society B*, 67(2),301-320

Buehlmann, P., Yu, B. (2003).
Boosting with the L2 loss: Regression and Classification.
Journal of the American Statistical Association, 98, 324-339

Efron, B., Hastie, T., Johnstone, I., Tibshirani, R. (2004).
Least Angle Regression.
Annals of Statistics, 32:407-499

Buehlmann, P., Yu, B. (2006).
Sparse Boosting.
Journal of Machine Learning Research, 7- 1001:1024

Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

See Also

[filter](#), [GenerateLearningsets](#), [tune](#), [classification](#)

Examples

```
# load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression from first 10 genes
golubX <- as.matrix(golub[,-1])
### Generate five different learningsets
set.seed(111)
five <- GenerateLearningsets(y=golubY, method = "CV", fold = 5, strat = TRUE)
### simple t-test:
selttest <- GeneSelection(golubX, golubY, learningsets = five, method = "t.test")
### show result:
show(selttest)
toplist(selttest, k = 10, iter = 1)
plot(selttest, iter = 1)
```

GenerateLearningsets

Repeated Divisions into learn- and test sets

Description

Due to very small sample sizes, the classical division learnset/testset does not give accurate information about the classification performance. Therefore, several different divisions should be used and aggregated. The implemented methods are discussed in Braga-Neto and Dougherty (2003) and Molinaro et al. (2005) whose terminology is adopted.

This function is usually the basis for all deeper analyses.

Usage

```
GenerateLearningsets(n, y, method = c("LOOCV", "CV", "MCCV", "bootstrap"),
                    fold = NULL, niter = NULL, ntrain = NULL, strat = FALSE)
```

Arguments

n	The total number of observations in the available data set. May be missing if y is provided instead.
y	A vector of class labels, either numeric or a factor. <i>Must</i> be given if strat=TRUE or n is not specified.
method	Which kind of scheme should be used to generate divisions into learning sets and test sets ? Can be one of the following: "LOOCV" Leaving-One-Out Cross Validation. "CV" (Ordinary) Cross-Validation. Note that fold must as well be specified. "MCCV" Monte-Carlo Cross Validation, i.e. random divisions into learning sets with ntrain(s.below) observations and tests sets with ntrain observations. "bootstrap" Learning sets are generated by drawing ntrain times with replacement from all observations. Those not drawn not all form the test set.
fold	Gives the number of CV-groups. Used only when method="CV"
niter	Number of iterations (s.details) .
ntrain	Number of observations in the learning sets. Used only when method="MCCV" or method="bootstrap".
strat	Logical. Should stratified sampling be performed, i.e. the proportion of observations from each class in the learning sets be the same as in the whole data set ? Does not apply for method = "LOOCV".

Details

- When method="CV", niter gives the number of times the whole CV-procedure is repeated. The output matrix has then foldxniter rows. When method="MCCV" or method="bootstrap", niter is simply the number of considered learning sets.
- Note that method="CV", fold=n is equivalent to method="LOOCV".

Value

An object of class `learningsets`

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

Christoph Bernau <bernau@ibe.med.uni-muenchen.de>

References

- Braga-Neto, U.M., Dougherty, E.R. (2003).
Is cross-validation valid for small-sample microarray classification ?
Bioinformatics, 20(3), 374-380
- Molinaro, A.M., Simon, R., Pfeiffer, R.M. (2005).
Prediction error estimation: a comparison of resampling methods.
Bioinformatics, 21(15), 3301-3307
- Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

See Also

[learningsets](#), [GeneSelection](#), [tune](#), [classification](#)

Examples

```
# LOOCV
loo <- GenerateLearningsets(n=40, method="LOOCV")
show(loo)
# five-fold-CV
CV5 <- GenerateLearningsets(n=40, method="CV", fold=5)
show(loo)
# MCCV
mccv <- GenerateLearningsets(n=40, method = "MCCV", niter=3, ntrain=30)
show(mccv)
# Bootstrap
boot <- GenerateLearningsets(n=40, method="bootstrap", niter=3)
# stratified five-fold-CV
set.seed(113)
classlabels <- sample(1:3, size = 50, replace = TRUE, prob = c(0.3, 0.5, 0.2))
CV5strat <- GenerateLearningsets(y = classlabels, method="CV", fold=5, strat = TRUE)
show(CV5strat)
```

LassoCMA-methods *L1 penalized logistic regression*

Description

The Lasso (Tibshirani, 1996) is one of the most popular tools for simultaneous shrinkage and variable selection. Recently, Friedman, Hastie and Tibshirani (2008) have developed an algorithm to compute the entire solution path of the Lasso for an arbitrary generalized linear model, implemented in the package `glmnet`. The method can be used for variable selection alone, s. [GeneSelection](#)

Methods

- X = "matrix", y = "numeric", f = "missing"** signature 1
- X = "matrix", y = "factor", f = "missing"** signature 2
- X = "data.frame", y = "missing", f = "formula"** signature 3
- X = "ExpressionSet", y = "character", f = "missing"** signature 4

For references, further argument and output information, consult [LassoCMA](#).

LassoCMA

*L1 penalized logistic regression***Description**

The Lasso (Tibshirani, 1996) is one of the most popular tools for simultaneous shrinkage and variable selection. Recently, Friedman, Hastie and Tibshirani (2008) have developed an algorithm to compute the entire solution path of the Lasso for an arbitrary generalized linear model, implemented in the package `glmnet`. The method can be used for variable selection alone, s. [GeneSelection](#).

For S4 method information, see `LassoCMA-methods`.

Usage

```
LassoCMA(X, y, f, learnind, norm.fraction = 0.1, models=FALSE, ...)
```

Arguments

<code>X</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>. note: by default, the predictors are scaled to have unit variance and zero mean. Can be changed by passing <code>standardize = FALSE</code> via the <code>...</code> argument.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p>
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be <code>missing</code> ; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>norm.fraction</code>	L1 Shrinkage intensity, expressed as the fraction of the coefficient L1 norm compared to the maximum possible L1 norm (corresponds to <code>fraction = 1</code>). Lower values correspond to higher shrinkage. Note that the default (0.1) need not produce good results, i.e. tuning of this parameter is recommended.
<code>models</code>	a logical value indicating whether the model object shall be returned
<code>...</code>	Further arguments passed to the function <code>glm</code> from the package of the same name.

Value

An object of class `clvargseloutput`.

Note

For a strongly related method, s. [ElasticNetCMA](#).
Up to now, this method can only be applied to binary classification.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>
Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>
Christoph Bernau <bernaue@ibe.med.uni-muenchen.de>

References

Tibshirani, R. (1996)
Regression shrinkage and selection via the lasso.
Journal of the Royal Statistical Society B, 58(1), 267-288

Friedman, J., Hastie, T. and Tibshirani, R. (2008) Regularization
Paths for Generalized Linear Models via Coordinate Descent
<http://www-stat.stanford.edu/~hastie/Papers/glmnet.pdf>

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#),
[ldaCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [pnnCMA](#),
[qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression
golubX <- as.matrix(golub[,-1])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run L1 penalized logistic regression (no tuning)
lassoresult <- LassoCMA(X=golubX, y=golubY, learnind=learnind, norm.fraction = 0.2)
show(lassoresult)
ftable(lassoresult)
plot(lassoresult)
```

Planarplot-methods *Visualize Separability of different classes*

Description

Given two variables, the methods trains a classifier (argument `classifier`) based on these two variables and plots the resulting class regions, learning- and test observations in the plane.

Appropriate variables are usually found by [GeneSelection](#).

Methods

X = "matrix", y = "numeric", f = "missing" signature 1
X = "matrix", y = "factor", f = "missing" signature 2
X = "data.frame", y = "missing", f = "formula" signature 3
X = "ExpressionSet", y = "character", f = "missing" signature 4
 For further argument and output information, consult [Planarplot](#).

 Planarplot

 Visualize Separability of different classes

Description

Given two variables, the methods trains a classifier (argument `classifier`) based on these two variables and plots the resulting class regions, learning- and test observations in the plane.

Appropriate variables are usually found by [GeneSelection](#).

For S4 method information, s. [Planarplot-methods](#).

Usage

```
Planarplot(X, y, f, learnind, predind, classifier, gridsize = 100, ...)
```

Arguments

<code>X</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided.
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be <code>missing</code> ; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>predind</code>	A vector containing <i>exactly</i> two indices that denote the two variables used for classification.
<code>classifier</code>	Name of function ending with <code>CMA</code> indicating the classifier to be used.
<code>gridsize</code>	The <code>gridsize</code> used for two-dimensional plotting. For both variables specified in <code>predind</code> , an equidistant grid of size <code>gridsize</code> is created. The resulting two grids are then combined to obtain <code>gridsize^2</code> points in the real plane which are used to draw the class regions. Defaults to 100 which is usually a reasonable choice, but takes some time.
<code>...</code>	Further argument passed to <code>classifier</code> .

Value

No return.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>. Idea is from the `MLInterfaces` package, contributed by Jess Mar, Robert Gentleman and Vince Carey.

See Also

[GeneSelection](#), [compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### simple linear discrimination for the golub data:
data(golub)
golubY <- golub[,1]
golubX <- as.matrix(golub[,-1])
golubn <- nrow(golubX)
set.seed(111)
learnind <- sample(golubn, size=floor(2/3*golubn))
Planarplot(X=golubX, y=golubY, learnind=learnind, predind=c(2,4),
           classifier=ldaCMA)
```

best

Show best hyperparameter settings

Description

In this package hyperparameter tuning is performed by an inner cross-validation step for each `learningset`. A grid of values is tried and evaluated in terms of the misclassification rate, the results are saved in an object of class `tuningresult`. This method displays (separately for each `learningset`) the hyperparameter/ hyperparameter combination that showed the best results. Note that this must not be unique; in this case, only one combination is displayed.

Usage

```
best(object, ...)
```

Arguments

`object` An object of class `tuningresult`.
`...` Currently unused argument.

Value

A list with elements equal to the number of different `learningsets`. Each element contains the best hyperparameter combination and the corresponding misclassification rate.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

See Also

tune

boxplot

Make a boxplot of the classifier evaluation

Description

This method displays the slot `scores` of performance scores of an object of class `evaloutput`.

Arguments

`x` An object of class `evaloutput`.

`...` Further graphical parameters passed to the classical `boxplot` function.

Value

The only return is a boxplot.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

See Also

[evaluation](#)

 classification-methods

General method for classification with various methods

Description

Perform classification for the following signatures:

Methods

`X = "matrix", y = "numeric", f = "missing"` signature 1

`X = "matrix", y = "factor", f = "missing"` signature 2

`X = "data.frame", y = "missing", f = "formula"` signature 3

`X = "ExpressionSet", y = "character", f = "missing"` signature 4

For further argument and output information, consult [classification](#).

 classification

General method for classification with various methods

Description

Most general function in the package, providing an interface to perform variable selection, hyper-parameter tuning and classification in one step. Alternatively, the first two steps can be performed separately and can then be plugged into this function.

For S4 method information, s. [classification-methods](#).

Usage

```
classification(X, y, f, learningsets, genesel, genesellist = list(), nbgene, cla
```

Arguments

- | | |
|----------------|--|
| <code>X</code> | Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>. |
| <code>y</code> | Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p> |
| <code>f</code> | A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables. |

<code>learningsets</code>	An object of class <code>learningsets</code> . May be missing, then the complete datasets is used as learning set.
<code>genesel</code>	Optional (but usually recommended) object of class <code>genesel</code> containing variable importance information for the argument <code>learningsets</code>
<code>geneselist</code>	In the case that the argument <code>genesel</code> is missing, this is an argument list passed to <code>GeneSelection</code> . If both <code>genesel</code> and <code>geneselist</code> are missing, no variable selection is performed.
<code>nbgene</code>	Number of best genes to be kept for classification, based on either <code>genesel</code> or the call to <code>GeneSelection</code> using <code>geneselist</code> . In the case that both are missing, this argument is not necessary. note: <ul style="list-style-type: none"> • If the gene selection method has been one of "lasso", "elasticnet", "boosting", <code>nbgene</code> will be reset to <code>min(s, nbgene)</code> where <code>s</code> is the number of nonzero coefficients. • if the gene selection scheme has been "one-vs-all", "pairwise" for the multiclass case, there exist several rankings. The top <code>nbgene</code> will be kept of <i>each</i> of them, so the number of effective used genes will sometimes be much larger.
<code>classifier</code>	Name of function ending with <code>CMA</code> indicating the classifier to be used.
<code>tuneres</code>	Analogous to the argument <code>genesel</code> - object of class <code>tuningresult</code> containing information about the best hyperparameter choice for the argument <code>learningsets</code> .
<code>tuninglist</code>	Analogous to the argument <code>geneselist</code> . In the case that the argument <code>tuneres</code> is missing, this in argument list passed to <code>tune</code> . If both <code>tuneres</code> and <code>tuninglist</code> are missing, no variable selection is performed. warning: Note that if a user-defined hyperparameter grid is passed, this will result in a list within a list: <code>tuninglist = list(grid=list(argname = c()), s. example. warning: Contrary to <code>tune</code>, if <code>tuninglist</code> is an empty list (default), no hyperparameter tuning will be performed at all. To use pre-defined hyperparameter grids, the argument is <code>tuninglist = list(grid = list())</code>.</code>
<code>trace</code>	Should progress be traced ? Default is <code>TRUE</code> .
<code>models</code>	a logical value indicating whether the model object shall be returned
<code>...</code>	Further arguments passed to the function <code>classifier</code> .

Details

For details about hyperparameter tuning, consult `tune`.

Value

A list of objects of class `cloutput` and `clvarseloutput`, respectively; its length equals the number of different `learningsets`. The single elements of the list can conveniently be combined using the `join` function. The results can be analyzed and evaluated by various measures using the method `evaluation`.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

Christoph Bernau <bernau@ibe.med.uni-muenchen.de>

References

Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

See Also

[GeneSelection](#), [tune](#), [evaluation](#), [compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### a simple k-nearest neighbour example
### datasets
## Not run: plot(x)
data(golub)
golubY <- golub[,1]
golubX <- as.matrix(golub[,-1])
### learningsets
set.seed(111)
lset <- GenerateLearningsets(y=golubY, method = "CV", fold=5, strat =TRUE)
### 1. GeneSelection
seltest <- GeneSelection(golubX, golubY, learningsets = lset, method = "t.test")
### 2. tuning
tunek <- tune(golubX, golubY, learningsets = lset, genesel = seltest, nbgene = 20, class
### 3. classification
knn1 <- classification(golubX, golubY, learningsets = lset, genesel = seltest,
                      tunerest = tunek, nbgene = 20, classifier = knnCMA)
### steps 1.-3. combined into one step:
knn2 <- classification(golubX, golubY, learningsets = lset,
                      genesellist = list(method = "t.test"), classifier = knnCMA,
                      tuninglist = list(grid = list(k = c(1:8))), nbgene = 20)
### show and analyze results:
knnjoin <- join(knn2)
show(knn2)
eval <- evaluation(knn2, measure = "misclassification")
show(eval)
summary(eval)
boxplot(eval)

## End(Not run)
```

cloutput-class "*cloutput*"

Description

Object returned by one of the classifiers (functions ending with CMA)

Slots

learnind: Vector of indices that indicates which observations were used in the learning set.

y: Actual (true) class labels of predicted observations.

yhat: Predicted class labels by the classifier.

prob: A numeric matrix whose rows equals the number of predicted observations (length of `y/yhat`) and whose columns equal the number of different classes in the learning set. Rows add up to one. Entry `j, k` of this matrix contains the probability for the `j`-th predicted observation to belong to class `k`. Can be a matrix of NAs, if the classifier used does not provide any probabilities

method: Name of the classifier used.

mode: character, one of "binary" (if the number of classes in the learning set is two) or `multiclass` (if it is more than two).

model: List containing the constructed classifiers.

Methods

show Use `show(cloutput-object)` for brief information

ftable Use `ftable(cloutput-object)` to obtain a confusion matrix/cross-tabulation of `y` vs. `yhat`, s. [ftable, cloutput-method](#).

plot Use `plot(cloutput-object)` to generate a probability plot of the matrix `prob` described above, s. [plot, cloutput-method](#)

roc Use `roc(cloutput-object)` to compute the empirical ROC curve and the Area Under the Curve (AUC) based on the predicted probabilities, s. [roc, cloutput-method](#)

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

See Also

`clvareoutput` [compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

```
clvareoutput-class
      "clvareoutput"
```

Description

Object returned by all classifiers that can perform variable selection or compute variable importance. These are:

- Random Forest, s. [rfCMA](#),
- Componentwise Boosting, s. [compBoostCMA](#),
- LASSO-logistic regression, s. [LassoCMA](#),
- ElasticNet-logistic regression, s. [ElasticNetCMA](#)

. Objects of class `clvareoutput` extend both the class `cloutput` and `varsel`, s. below.

Slots

learnind: Vector of indices that indicates which observations were used in the learning set.

y: Actual (true) class labels of predicted observations.

yhat: Predicted class labels by the classifier.

prob: A numeric matrix whose rows equals the number of predicted observations (length of `y/yhat`) and whose columns equal the number of different classes in the learning set. Rows add up to one. Entry `j, k` of this matrix contains the probability for the `j`-th predicted observation to belong to class `k`. Can be a matrix of NAs, if the classifier used does not provide any probabilities

method: Name of the classifier used.

mode: character, one of "binary" (if the number of classes in the learning set is two) or "multiclass" (if it is more than two).

varels: numeric vector of variable importance measures (for Random Forest) or absolute values of regression coefficients (for the other three methods mentioned above) (from which the majority will be zero).

Extends

Class "`cloutput`", directly. Class "`varelooutput`", directly.

Methods

show Use `show(cloutput-object)` for brief information

ftable Use `ftable(cloutput-object)` to obtain a confusion matrix/cross-tabulation of `y` vs. `yhat`, s. [ftable, cloutput-method](#).

plot Use `plot(cloutput-object)` to generate a probability plot of the matrix `prob` described above, s. [plot, cloutput-method](#)

roc Use `roc(cloutput-object)` to compute the empirical ROC curve and the Area Under the Curve (AUC) based on the predicted probabilities, s. [roc, cloutput-method](#)

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

See Also

[rfCMA](#), [compBoostCMA](#), [LassoCMA](#), [ElasticNetCMA](#)

 compBoostCMA-methods

Componentwise Boosting

Description

Roughly speaking, Boosting combines 'weak learners' in a weighted manner in a stronger ensemble.

'Weak learners' here consist of linear functions in one component (variable), as proposed by Buehlmann and Yu (2003).

It also generates sparsity and can as well be as used for variable selection alone. (s. [GeneSelection](#).)

Methods

X = "matrix", y = "numeric", f = "missing" signature 1

X = "matrix", y = "factor", f = "missing" signature 2

X = "data.frame", y = "missing", f = "formula" signature 3

X = "ExpressionSet", y = "character", f = "missing" signature 4

For further argument and output information, consult [compBoostCMA](#).

 compBoostCMA

Componentwise Boosting

Description

Roughly speaking, Boosting combines 'weak learners' in a weighted manner in a stronger ensemble.

'Weak learners' here consist of linear functions in one component (variable), as proposed by Buehlmann and Yu (2003).

It also generates sparsity and can as well be as used for variable selection alone. (s. [GeneSelection](#)).

For S4 method information, see [compBoostCMA-methods](#).

Usage

```
compBoostCMA(X, y, f, learnind, loss = c("binomial", "exp", "quadratic"), mstop
```

Arguments

- | | |
|---|--|
| X | Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>. |
| y | Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. |

- A character if `X` is an `ExpressionSet` that specifies the phenotype variable.
- `missing`, if `X` is a `data.frame` and a proper formula `f` is provided.

WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.

<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be <code>missing</code> ; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>loss</code>	Character specifying the loss function - one of "binomial" (LogitBoost), "exp" (AdaBoost), "quadratic"(L2Boost).
<code>mstop</code>	Number of boosting iterations, i.e. number of updates to perform. The default (100) does not necessarily produce good results, therefore usage of <code>tune</code> for this argument is highly recommended.
<code>nu</code>	Shrinkage factor applied to the update steps, defaults to 0.1. In most cases, it suffices to set <code>nu</code> to a very low value and to concentrate on the optimization of <code>mstop</code> .
<code>models</code>	a logical value indicating whether the model object shall be returned
<code>...</code>	Currently unused arguments.

Details

The method is partly based on code from the package `mboost` from T. Hothorn and P. Buehlmann. The algorithm for the multiclass case is described in Lutz and Buehlmann (2006) as 'rowwise updating'.

Value

An object of class `clvarseloutput`.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Buehlmann, P., Yu, B. (2003).

Boosting with the L2 loss: Regression and Classification.

Journal of the American Statistical Association, 98, 324-339

Buehlmann, P., Hothorn, T.

Boosting: A statistical perspective.

Statistical Science (to appear)

Lutz, R., Buehlmann, P. (2006).

Boosting for high-multivariate responses in high-dimensional linear regression.

Statistica Sinica 16, 471-494.

See Also

[dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression
golubX <- as.matrix(golub[,-1])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run componentwise (logit)-boosting (not tuned)
result <- compBoostCMA(X=golubX, y=golubY, learnind=learnind, mstop = 500)
### show results
show(result)
ftable(result)
plot(result)
### multiclass example:
### load Khan data
data(khan)
### extract class labels
khanY <- khan[,1]
### extract gene expression
khanX <- as.matrix(khan[,-1])
### select learningset
set.seed(111)
learnind <- sample(length(khanY), size=floor(ratio*length(khanY)))
### run componentwise multivariate (logit)-boosting (not tuned)
result <- compBoostCMA(X=khanX, y=khanY, learnind=learnind, mstop = 1000)
### show results
show(result)
ftable(result)
plot(result)
```

compare-methods

Compare different classifiers

Description

Compare different classifiers for the following signatures:

Methods

clresultlist = "list" signature 1

For further argument and output information, consult [compare](#)

compare

*Compare different classifiers***Description**

Classifiers can be evaluated separately using the method `evaluation`. Normally, several classifiers are used for the same dataset and their performance is compared. This comparison procedure is essentially facilitated by this method. For S4 method information, s. `compare-methods`

Usage

```
compare(clresultlist, measure = c("misclassification", "sensitivity",
  "specificity", "average probability", "brier score", "auc"), aggfun =
  meanrm, plot = FALSE, ...)
```

Arguments

- `clresultlist` A list of lists (!) of objects of class `cloutput` or `clvarseloutput`. Each inner list is usually returned by `classification`. Additionally, the different list elements of the outer list should have been created by different classifiers, s. also example below.
- `measure` A character vector containing one or more of the elements listed below. By default, all measures are computed, using `evaluation` with `scheme = "iterationwise"`. Note that "sensitivity", "specificity", "auc" cannot be computed for the multiclass case.
- "misclassification" The missclassification rate.
 - "sensitivity" The sensitivity or 1-false negative rate. Can only be computed for binary classification.
 - "specificity" The specificity or 1-false positive rate. Can only be computed for binary classification.
 - "average probability" The average probability assigned to the correct class. Requirement is that the used classifier provides probability estimations. The optimum performance is 1.
 - "brier score" The Brier Score is generally defined as $\langle \text{sum over all observation } i \rangle \langle \text{sum over all classes } k \rangle (I(y_i=k) - P(k))^2$, with $I()$ denoting the indicator function and $P(k)$ the estimated probability for class k . The optimum performance is 0.
 - "auc" The Area under the Curve (AUC) belonging to the empirical ROC curve computed from the estimated probabilities and the true class labels. Can only be computed for binary classification and if "scheme = iterationwise", s. below. S. also `roc`, `cloutput-method`.
- `aggfun` Function that determines how performance among different iterations are aggregated. Default is `meanrm`, which computes the mean using `na.rm=T`. Other possible choices are quantiles.
- `plot` Should the performance of different classifiers be visualized by a joint boxplot? Default is `FALSE`.
- `...` Further arguments passed to `boxplot` in the case that `plot = TRUE`.

Value

A `data.frame` with rows corresponding to the compared classifiers and columns to the performance measures, aggregated by `aggfun`, s. above.

Note

If more than one measure is computed and `plot = TRUE`, one separate plot is created for each of them.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

Christoph Bernau <bernauc@ibe.med.uni-muenchen.de>

References

- Dudoit, S., Fridlyand, J., Speed, T. P. (2002)
Comparison of discrimination methods for the classification of tumors using gene expression data.
Journal of the American Statistical Association 97, 77-87
- Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

See Also

[classification](#), [evaluation](#)

Examples

```
## Not run:
### compare the performance of several discriminant analysis methods
### for the Khan dataset:
data(khan)
khanX <- as.matrix(khan[,-1])
khanY <- khan[,1]
set.seed(27611)
fiveCV10iter <- GenerateLearningsets(y=khanY, method = "CV", fold = 5, niter = 2, strat =
### candidate methods: DLDA, LDA, QDA, pls_LDA, sclda
class_dllda <- classification(X = khanX, y=khanY, learningsets = fiveCV10iter, classifier =
### perform GeneSlection for LDA, FDA, QDA (using F-Tests):
genesel_da <- GeneSelection(X=khanX, y=khanY, learningsets = fiveCV10iter, method = "f.te
###
class_lda <- classification(X = khanX, y=khanY, learningsets = fiveCV10iter, classifier =
class_qda <- classification(X = khanX, y=khanY, learningsets = fiveCV10iter, classifier =

### We now make a comparison concerning the performance (sev. measures):
### first, collect in a list:
dalike <- list(class_dllda, class_lda, class_qda)
### use pre-defined compare function:
comparison <- compare(dalike, plot = TRUE, measure = c("misclassification", "brier score")
print(comparison)

## End(Not run)
```

dldaCMA-methods *Diagonal Discriminant Analysis*

Description

Performs a diagonal discriminant analysis under the assumption of a multivariate normal distribution in each classes (with equal, diagonally structured) covariance matrices. The method is also known under the name 'naive Bayes' classifier.

Methods

X = "matrix", y = "numeric", f = "missing" signature 1

X = "matrix", y = "factor", f = "missing" signature 2

X = "data.frame", y = "missing", f = "formula" signature 3

X = "ExpressionSet", y = "character", f = "missing" signature 4

For further argument and output information, consult [dldaCMA](#).

dldaCMA *Diagonal Discriminant Analysis*

Description

Performs a diagonal discriminant analysis under the assumption of a multivariate normal distribution in each classes (with equal, diagonally structured) covariance matrices. The method is also known under the name 'naive Bayes' classifier.

For S4 method information, see [dldaCMA-methods](#).

Usage

```
dldaCMA(X, y, f, learnind, models=FALSE, ...)
```

Arguments

- | | |
|---|---|
| X | Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>. |
| y | Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if X is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if X is a <code>data.frame</code> and a proper formula <code>f</code> is provided. |

WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.

f	A two-sided formula, if X is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
learnind	An index vector specifying the observations that belong to the learning set. May be missing; in that case, the learning set consists of all observations and predictions are made on the learning set.
models	a logical value indicating whether the model object shall be returned
...	Currently unused argument.

Value

An object of class `cloutput`.

Note

As opposed to linear or quadratic discriminant analysis, variable selection is not strictly necessary.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

McLachlan, G.J. (1992).

Discriminant Analysis and Statistical Pattern Recognition.

Wiley, New York

See Also

[compBoostCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression
golubX <- as.matrix(golub[,-1])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run DLDA
dldaresult <- dldaCMA(X=golubX, y=golubY, learnind=learnind)
### show results
show(dldaresult)
ftable(dldaresult)
plot(dldaresult)
### multiclass example:
### load Khan data
```

```

data(khan)
### extract class labels
khanY <- khan[,1]
### extract gene expression
khanX <- as.matrix(khan[,-1])
### select learningset
set.seed(111)
learnind <- sample(length(khanY), size=floor(ratio*length(khanY)))
### run LDA
ldaresult <- dldaCMA(X=khanX, y=khanY, learnind=learnind)
### show results
show(dldaresult)
ftable(dldaresult)
plot(dldaresult)

```

```
evaloutput-class    "evaloutput"
```

Description

Object returned by the method [evaluation](#).

Slots

score: A numeric vector of performance scores whose length depends on "scheme", s.below. It equals the number of iterations (number of different datasets) if "scheme = iterationwise" and the number of all observations in the complete dataset otherwise. As not necessarily all observation must be predicted at least one time, score can also contain NAs for those observations not classified at all.

measure: performance measure used, s. [evaluation](#).

scheme: scheme used, s. [evaluation](#)

method: name of the classifier that has been evaluated.

Methods

show Use `show(evaloutput-object)` for brief information.

summary Use `summary(evaloutput-object)` to apply the classic `summary()` function to the slot `score`, s. [summary](#), [evaloutput-method](#)

boxplot Use `boxplot(evaloutput-object)` to display a boxplot of the slot `score`, s. [boxplot](#), [evaloutput-method](#).

obsinfo Use `obsinfo(evaloutput-object, threshold)` to display all observations consistently correctly or incorrectly classified (depending on the value of the argument `threshold`), s. [obsinfo](#).

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

See Also

[evaluation](#)

evaluation-methods *Evaluation of classifiers*

Description

Evaluate classifiers for the following signatures:

Methods

clresult = "list" signature 1

For further argument and output information, consult [evaluation](#).

evaluation *Evaluation of classifiers*

Description

The performance of classifiers can be evaluated by six different measures and two different schemes that are described more precisely below.

For S4 method information, s. [evaluation-methods](#).

Usage

```
evaluation(clresult, cltrain = NULL, cost = NULL, y = NULL, measure = c("misclas
scheme = c("iterationwise", "observationwise", "classwise")
```

Arguments

clresult	A list of objects of class <code>cloutput</code> or <code>clvarseloutput</code>
cltrain	An object of class <code>cloutput</code> in which the <i>whole</i> dataset was used as learning set. Only used if <code>method = "0.632"</code> or <code>method = "0.632+"</code> in order to obtain an estimation for the resubstitution error rate.
cost	An optional cost matrix used if <code>measure = "misclassification"</code> . If it is not specified (default), the cost is the usual indicator loss. Otherwise, entry i, j of <code>cost</code> quantifies the loss when the true class is class $i-1$ and the predicted class is $j-1$, provided the conventional coding $0, \dots, K-1$ in the case of K classes is used. Usually, the matrix contains only non-negative entries with zeros on the diagonal, but this is not obligatory. Make sure that the dimension of the matrix matches the number of classes.
y	A vector containing the true class labels. Only needed if <code>scheme = "classwise"</code> .
measure	Performance measure to be used: "misclassification" The missclassification rate. "sensitivity" The sensitivity or 1-false negative rate. Can only be computed for binary classification. "specificity" The specificity or 1-false positive rate. Can only be computed for binary classification.

	"average probability" The average probability assigned to the correct class. Requirement is that the used classifier provides probability estimations. The optimum performance is 1.
	"brier score" The Brier Score is generally defined as $\langle \text{sum over all observation } i \rangle \langle \text{sum over all classes } k \rangle (I(y_{i=k}) - P(k))^2$, with $I()$ denoting the indicator function and $P(k)$ the estimated probability for class k . The optimum performance is 0.
	"auc" The Area under the Curve (AUC) belonging to the empirical ROC curve computed from the estimated probabilities and the true class labels. Can only be computed for binary classification and if "scheme = iterationwise", s. below. S. also <code>roc</code> , <code>cloutput-method</code> .
	"0.632" The 0.632 estimator (s. reference) for the misclassification rate (applied iteration- or) observationwise, if bootstrap learning sets have been used. Note that <code>cltrain</code> must be provided.
	"0.632+" The 0.632+ estimator (s. reference) for the misclassification rate (applied iteration- or) observationwise, if bootstrap learning sets have been used. Note that <code>cltrain</code> must be provided.
scheme	"iterationwise" The performance measures listed above are computed for each different iteration, i.e. each different <code>learningset</code>
	"observationwise" The performance measures listed above (except for "auc") are computed separately for each observation classified one or several times, depending on the <code>learningset</code> scheme.
	"classwise" The performance measures (exceptions: "auc", "0.632", "0.632+") are computed separately for each class, averaged over both iterations and observations.

Value

An object of class `evaloutput`.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

Christoph Bernau <bernau@ibe.med.uni-muenchen.de>

References

Efron, B. and Tibshirani, R. (1997). Improvements on cross-validation: The .632+ bootstrap method.

Journal of the American Statistical Association, 92, 548-560.

Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

See Also

`evaloutput`, `classification`, `compare`

Examples

```
### simple linear discriminant analysis example using bootstrap datasets:
### datasets:
data(golub)
golubY <- golub[,1]
### extract gene expression from first 10 genes
golubX <- as.matrix(golub[,2:11])
### generate 25 bootstrap datasets
set.seed(333)
bootds <- GenerateLearningsets(y = golubY, method = "bootstrap", ntrain = 30, niter = 10,
### run classification()
ldalist <- classification(X=golubX, y=golubY, learningsets = bootds, classifier=ldaCMA)
### Evaluation:
eval_iter <- evaluation(ldalist, scheme = "iter")
eval_obs <- evaluation(ldalist, scheme = "obs")
show(eval_iter)
show(eval_obs)
summary(eval_iter)
summary(eval_obs)
### auc with boxplot
eval_auc <- evaluation(ldalist, scheme = "iter", measure = "auc")
boxplot(eval_auc)
### which observations have often been misclassified ?
obsinfo(eval_obs, threshold = 0.75)
```

Description

Fisher's Linear Discriminant Analysis constructs a subspace of 'optimal projections' in which classification is performed. The directions of optimal projections are computed by the function `cancor` from the package `stats`. For an exhaustive treatment, see e.g. Ripley (1996).

Methods

X = "matrix", y = "numeric", f = "missing" signature 1

X = "matrix", y = "factor", f = "missing" signature 2

X = "data.frame", y = "missing", f = "formula" signature 3

X = "ExpressionSet", y = "character", f = "missing" signature 4

For references, further argument and output information, consult [fdaCMA](#).

fdaCMA

*Fisher's Linear Discriminant Analysis***Description**

Fisher's Linear Discriminant Analysis constructs a subspace of 'optimal projections' in which classification is performed. The directions of optimal projections are computed by the function `cancor` from the package `stats`. For an exhaustive treatment, see e.g. Ripley (1996).

For S4 method information, see [fdaCMA-methods](#).

Usage

```
fdaCMA(X, y, f, learnind, comp = 1, plot = FALSE, models=FALSE)
```

Arguments

<code>X</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A <code>factor</code>. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be <code>missing</code> ; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>comp</code>	Number of discriminant coordinates (projections) to compute. Default is one, must be smaller than or equal to $K-1$, where K is the number of classes.
<code>plot</code>	Should the projections onto the space spanned by the optimal projection directions be plotted? Default is <code>FALSE</code> .
<code>models</code>	a logical value indicating whether the model object shall be returned

Value

An object of class `cloutput`.

Note

Excessive variable selection has usually to performed before `fdaCMA` can be applied in the $p > n$ setting. Not reducing the number of variables can result in an error message.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Ripley, B.D. (1996)

Pattern Recognition and Neural Networks.

Cambridge University Press

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression from first 10 genes
golubX <- as.matrix(golub[,2:11])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run FDA
fdaresult <- fdaCMA(X=golubX, y=golubY, learnind=learnind, comp = 1, plot = TRUE)
### show results
show(fdaresult)
ftable(fdaresult)
plot(fdaresult)
### multiclass example:
### load Khan data
data(khan)
### extract class labels
khanY <- khan[,1]
### extract gene expression from first 10 genes
khanX <- as.matrix(khan[,2:11])
### select learningset
set.seed(111)
learnind <- sample(length(khanY), size=floor(ratio*length(khanY)))
### run FDA
fdaresult <- fdaCMA(X=khanX, y=khanY, learnind=learnind, comp = 2, plot = TRUE)
### show results
show(fdaresult)
ftable(fdaresult)
plot(fdaresult)
```

`filter`*Filter functions for Gene Selection*

Description

The functions listed above are usually not called by the user but via [GeneSelection](#).

Usage

```
ttest(X, y, learnind, ...)  
welchtest(X, y, learnind, ...)  
ftest(X, y, learnind, ...)  
kruskaltest(X, y, learnind, ...)  
limmatest(X, y, learnind, ...)  
golubcrit(X, y, learnind, ...)  
rfe(X, y, learnind, ...)
```

Arguments

<code>X</code>	A numeric matrix of gene expression values.
<code>y</code>	A numeric vector of class labels.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set.
<code>...</code>	Currently unused argument.

Value

An object of class [varseloutput](#).

References

Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

`flexdaCMA-methods`*Flexible Discriminant Analysis*

Description

This method is experimental.

It is easy to show that, after appropriate scaling of the predictor matrix X , Fisher's Linear Discriminant Analysis is equivalent to Discriminant Analysis in the space of the fitted values from the linear regression of the $n_{\text{learn}} \times K$ indicator matrix of the class labels on X . This gives rise to 'non-linear discriminant analysis' methods that expand X in a suitable, more flexible basis. In order to avoid overfitting, penalization is used. In the implemented version, the linear model is replaced by a generalized additive one, using the package `mgcv`.

Methods

`X = "matrix", y = "numeric", f = "missing"` signature 1
`X = "matrix", y = "factor", f = "missing"` signature 2
`X = "data.frame", y = "missing", f = "formula"` signature 3
`X = "ExpressionSet", y = "character", f = "missing"` signature 4
 For further argument and output information, consult [flexdaCMA](#).

 flexdaCMA

Flexible Discriminant Analysis

Description

This method is experimental.

It is easy to show that, after appropriate scaling of the predictor matrix X , Fisher's Linear Discriminant Analysis is equivalent to Discriminant Analysis in the space of the fitted values from the linear regression of the $n_{\text{learn}} \times K$ indicator matrix of the class labels on X . This gives rise to 'non-linear discriminant analysis' methods that expand X in a suitable, more flexible basis. In order to avoid overfitting, penalization is used. In the implemented version, the linear model is replaced by a generalized additive one, using the package `mgcv`.

For S4 method information, s. [flexdaCMA-methods](#).

Usage

```
flexdaCMA(X, y, f, learnind, comp = 1, plot = FALSE, models=FALSE, ...)
```

Arguments

<code>X</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p>
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be <code>missing</code> ; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>comp</code>	Number of discriminant coordinates (projections) to compute. Default is one, must be smaller than or equal to $K-1$, where K is the number of classes.

plot	Should the projections onto the space spanned by the optimal projection directions be plotted ? Default is FALSE.
models	a logical value indicating whether the model object shall be returned
...	Further arguments passed to the function <code>gam</code> from the package <code>mgcv</code> .

Value

An object of class `cloutput`.

Note

Excessive variable selection has usually to performed before `flexdaCMA` can be applied in the `p > n` setting. Recall that the original predictor dimension is even enlarged, therefore, it should be applied only with very few variables.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Ripley, B.D. (1996)
 Pattern Recognition and Neural Networks.
Cambridge University Press

See Also

[compBoostCMA](#), [dlldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression from first 5 genes
golubX <- as.matrix(golub[,2:6])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run flexible Discriminant Analysis
result <- flexdaCMA(X=golubX, y=golubY, learnind=learnind, comp = 1)
### show results
show(result)
ftable(result)
plot(result)
```

ftable

*Cross-tabulation of predicted and true class labels***Description**

An object of class `cloutput` contains (among others) the slot `y` and `yhat`. The former contains the true, the last the predicted class labels. Both are cross-tabulated in order to obtain a so-called confusion matrix. Counts out of the diagonal are misclassifications.

Arguments

`x` An object of class `cloutput`
`...` Currently unused argument.

Value

No return.

Author(s)

Martin Slawski <martin.slawski@campus.lmu.de>
 Anne-Laure Boulesteix <http://www.slcmsr.net/boulesteix>

See Also

For more advanced evaluation: [evaluation](#)

gbmCMA-methods

*Tree-based Gradient Boosting***Description**

Roughly speaking, Boosting combines 'weak learners' in a weighted manner in a stronger ensemble. This method calls the function `gbm.fit` from the package `gbm`. The 'weak learners' are simple trees that need only very few splits (default: 1).

Methods

`X = "matrix", y = "numeric", f = "missing"` signature 1
`X = "matrix", y = "factor", f = "missing"` signature 2
`X = "data.frame", y = "missing", f = "formula"` signature 3
`X = "ExpressionSet", y = "character", f = "missing"` signature 4
 For further argument and output information, consult [gbmCMA](#).

gbmCMA

*Tree-based Gradient Boosting***Description**

Roughly speaking, Boosting combines 'weak learners' in a weighted manner in a stronger ensemble. This method calls the function `gbm.fit` from the package `gbm`. The 'weak learners' are simple trees that need only very few splits (default: 1).

For S4 method information, see [gbmCMA-methods](#).

Usage

```
gbmCMA(X, y, f, learnind, models=FALSE, ...)
```

Arguments

<code>X</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • missing, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p>
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be missing; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>models</code>	a logical value indicating whether the model object shall be returned
<code>...</code>	Further arguments passed to the function <code>gbm.fit</code> from the package of the same name. Worth mentioning are <ul style="list-style-type: none"> <code>ntrees</code> Number of trees to fit (size of the ensemble), defaults to 100. This parameter should be optimized using tune. <code>shrinkage</code> The learning rate (default is 0.001). Usually fixed to a very low value. <code>distribution</code> Loss function to be used. Default is "bernoulli", i.e. <code>LogitBoost</code>, a (less robust) alternative is "adaboost". <code>interaction.depth</code> Number of splits used by the 'weak learner' (single decision tree). Default is 1.

Value

An object of class `cloutput`.

Note

Up to now, this method can only be applied to binary classification.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Ridgeway, G. (1999).

The state of boosting.

Computing Science and Statistics, 31:172-181

Friedman, J. (2001).

Greedy Function Approximation: A Gradient Boosting Machine.

Annals of Statistics 29(5):1189-1232.

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression
golubX <- as.matrix(golub[,-1])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run tree-based gradient boosting (no tuning)
gbmresult <- gbmCMA(X=golubX, y=golubY, learnind=learnind, n.trees = 500)
show(gbmresult)
ftable(gbmresult)
plot(gbmresult)
```

genesel-class

"genesel"

Description

Object returned from a call to [GeneSelection](#)

Slots

rankings: A list of matrices. For the two-class case and the multi-class case where a genuine multi-class method has been used for variable selection, the length of the list is one. Otherwise, it is named according to the different binary scenarios (e.g. 1 vs 3). Each list element is a matrix with rows corresponding to iterations (different `learningsets`) and columns to variables. Each row thus contains an index vector representing the order of the variables with respect to their variable importance (s. slot `importance`)

importance: A list of matrices, with the same structure as described for the slot `rankings`. Each row of these matrices are ordered according to `rankings` and contain the variable importance measure (absolute value of test statistic or regression coefficient).

method: Name of the method used for variable selection, s. [GeneSelection](#).

scheme: The scheme used in the case of a non-binary response, one of "pairwise", "one-vs-all" or "multiclass".

Methods

show Use `show(genesel-object)` for brief information

toplist Use `toplist(genesel-object, k=10, iter = 1)` to display the top first 10 variables and their variable importance for the first iteration (first `learningset`), s. [toplist](#).

plot Use `plot(genesel-object, k=10, iter=1)` to display a barplot of the variable importance of the top first 10 variables, s. [plot, genesel-method](#)

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

See Also

[GeneSelection](#)

golub

ALL/AML dataset of Golub et al. (1999)

Description

s. below

Usage

```
data(golub)
```

Format

A data frame with 38 observations and 3052 variables. The first column (named `golub.c1`) contains the tumor classes (ALL = acute lymphatic leukaemia, AML = acute myeloid leukaemia).\ `golub.c1`: a factor with levels ALL AML.\ X2-X3051: Gene expression values.

Source

Adopted from the dataset in the package `multtest`.

References

Golub, T., Slonim, D. K., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov, J. P., Coller, H., Loh, M. L., Downing, J., Caligiuri, M. A., Bloomfeld, C. D., Lander, E. S. (1999).
Molecular classification of cancer: class discovery and class prediction by gene expression monitoring.
Science 286, 531-537.

Examples

```
data(golub)
```

`internals`

Internal functions

Description

Not intended to be called directly by the user.

`join-methods`

Combine list elements returned by the method classification

Description

The list of objects of class `cloutput` can be unified into one object for the following signatures:

Methods

`cloutputlist = "list"` signature 1

For further argument and output information, consult `join`.

`join`*Combine list elements returned by the method classification*

Description

The method `classification` returns a list of class `cloutput` or `clvarseloutput`. It is often more convenient to work with an object of class `cloutput` instead with a whole list, e.g. because the convenience method defined for that class can be used.

For S4 method information, s. [join-methods](#)

Usage

```
join(cloutputlist)
```

Arguments

`cloutputlist` A list of objects of classes `cloutput` or `clvarseloutput`, usually that returned by a call to the method `classification`. The only requirement for a succesful join is that the used dataset and classifier are the same for each list element.

Value

An object of class `cloutput`. **warning:**If the elements of `cloutputlist` have originally been of class `clvarseloutput`, the slot `varsel` will be dropped !

Note

The result of the `join` method is incompatible with the methods `evaluation`, `compare`. These require the lists returned by `classification`.

See Also

[classification](#), [evaluation](#)

`khan`*Small blue round cell tumor dataset of Khan et al. (2001)*

Description

s. below

Usage

```
data(khan)
```

Format

A data frame with 63 observations on the following 2309 variables. The first column (named `khanY`) contains the tumor classes (BL = Burkitt Lymphoma, EWS = Ewing Sarcoma, NB = Neuro Blastoma, RMS = Rhabdomyosarcoma).\

`khanY`: a factor with levels BL EWS NB RMS \ X2-X2309: Gene expression values.

Source

Adopted from the dataset in the package `pamr`.

References

Khan, J., Wei, J. S., Ringner, M., Saal, L. H., Ladanyi, M., Westermann, F., Berthold, F., Schwab, M., Antonescu, C. R., Peterson, C., Meltzer, P. S., (2001).

Classification and diagnostic prediction of cancers using gene expression profiling and artificial neural networks.

Nature Medicine 7, 673-679.

Examples

```
data(khan)
```

<code>knnCMA-methods</code>	<i>Nearest Neighbours</i>
-----------------------------	---------------------------

Description

Ordinary `k` nearest neighbours algorithm from the very fast implementation in the package `class`

Methods

X = "matrix", y = "numeric", f = "missing" signature 1

X = "matrix", y = "factor", f = "missing" signature 2

X = "data.frame", y = "missing", f = "formula" signature 3

X = "ExpressionSet", y = "character", f = "missing" signature 4

For further argument and output information, consult [knnCMA](#).

knnCMA	<i>Nearest Neighbours</i>
--------	---------------------------

Description

Ordinary k nearest neighbours algorithm from the very fast implementation in the package `class`.
 For S4 method information, see [knnCMA-methods](#).

Usage

```
knnCMA(X, y, f, learnind, models=FALSE, ...)
```

Arguments

<code>X</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A <code>factor</code>. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p>
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. Must not be missing for this method.
<code>models</code>	a logical value indicating whether the model object shall be returned
<code>...</code>	Further arguments to be passed to <code>knn</code> from the package <code>class</code> , in particular the number of nearest neighbours to use (argument <code>k</code>).

Value

An object of class `cloutput`.

Note

Class probabilities are *not* returned. For a probabilistic variant of `knn`, s. [pknnCMA](#).

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Ripley, B.D. (1996)
 Pattern Recognition and Neural Networks.
Cambridge University Press

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [ldaCMA](#),
[LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfCMA](#), [pnnCMA](#),
[qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression from first 10 genes
golubX <- as.matrix(golub[,-1])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run k-nearest neighbours
result <- knnCMA(X=golubX, y=golubY, learnind=learnind, k = 3)
### show results
show(result)
ftable(result)
### multiclass example:
### load Khan data
data(khan)
### extract class labels
khanY <- khan[,1]
### extract gene expression
khanX <- as.matrix(khan[,-1])
### select learningset
set.seed(111)
learnind <- sample(length(khanY), size=floor(ratio*length(khanY)))
### run knn
result <- knnCMA(X=khanX, y=khanY, learnind=learnind, k = 5)
### show results
show(result)
ftable(result)
```

Description

Performs a linear discriminant analysis for the following signatures:

Methods

`X = "matrix", y = "numeric", f = "missing"` signature 1
`X = "matrix", y = "factor", f = "missing"` signature 2
`X = "data.frame", y = "missing", f = "formula"` signature 3
`X = "ExpressionSet", y = "character", f = "missing"` signature 4
 For further argument and output information, consult [ldaCMA](#).

ldaCMA

*Linear Discriminant Analysis***Description**

Performs a linear discriminant analysis under the assumption of a multivariate normal distribution in each classes (with equal, but generally structured) covariance matrices. The function `lda` from the package `MASS` is called for computation.

For S4 method information, see [ldaCMA-methods](#).

Usage

```
ldaCMA(X, y, f, learnind, models=FALSE, ...)
```

Arguments

<code>X</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p>
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be <code>missing</code> ; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>models</code>	a logical value indicating whether the model object shall be returned
<code>...</code>	Further arguments to be passed to <code>lda</code> from the package <code>MASS</code>

Value

An object of class `cloutput`.

Note

Excessive variable selection has usually to be performed before ldaCMA can be applied in the $p > n$ setting. Not reducing the number of variables can result in an error message.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

McLachlan, G.J. (1992).

Discriminant Analysis and Statistical Pattern Recognition.

Wiley, New York

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfCMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression from first 10 genes
golubX <- as.matrix(golub[,2:11])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run LDA
ldaresult <- ldaCMA(X=golubX, y=golubY, learnind=learnind)
### show results
show(ldaresult)
ftable(ldaresult)
plot(ldaresult)
### multiclass example:
### load Khan data
data(khan)
### extract class labels
khanY <- khan[,1]
### extract gene expression from first 10 genes
khanX <- as.matrix(khan[,2:11])
### select learningset
set.seed(111)
learnind <- sample(length(khanY), size=floor(ratio*length(khanY)))
### run LDA
ldaresult <- ldaCMA(X=khanX, y=khanY, learnind=learnind)
### show results
show(ldaresult)
ftable(ldaresult)
plot(ldaresult)
```

```
learningsets-class "learningsets"
```

Description

An object returned from `GenerateLearningsets` which is usually passed as arguments to `GeneSelection`, `tune` and `classification`.

Slots

`learnmatrix`: A matrix of dimension `niter` x `ntrain`. Each row contains the indices of those observations representing the learningset for one iteration. If `method = CV`, zeros appear due to rounding issues.

`method`: The method used to generate the `learnmatrix`, s.`GenerateLearningsets`

`ntrain`: Number of observations in one learning set. If `method = CV`, this number is not attained for all iterations, due to rounding issues.

`iter`: Number of iterations (different learningsets) that are stored in `learnmatrix`.

Methods

- `show` Use `show(learningsets-object)` for brief information.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

Christoph Bernau <bernau@ibe.med.uni-muenchen.de>

See Also

`GenerateLearningsets`, `GeneSelection`, `tune`, `classification`

```
nnetCMA-methods      Feed-Forward Neural Networks
```

Description

This method provides access to the function `nnet` in the package of the same name that trains Feed-forward Neural Networks with one hidden layer.

Methods

`X = "matrix", y = "numeric", f = "missing"` signature 1

`X = "matrix", y = "factor", f = "missing"` signature 2

`X = "data.frame", y = "missing", f = "formula"` signature 3

`X = "ExpressionSet", y = "character", f = "missing"` signature 4

For further argument and output information, consult `nnetCMA`.

nnetCMA

*Feed-forward Neural Networks***Description**

This method provides access to the function `nnet` in the package of the same name that trains Feed-forward Neural Networks with one hidden layer.
For S4 method information, see [nnetCMA-methods](#)

Usage

```
nnetCMA(X, y, f, learnind, eigengenes = FALSE, models=FALSE, ...)
```

Arguments

- | | |
|------------|--|
| X | Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>. |
| y | Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if X is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if X is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p> |
| f | A two-sided formula, if X is a <code>data.frame</code> . The left part correspond to class labels, the right to variables. |
| learnind | An index vector specifying the observations that belong to the learning set. May be <code>missing</code> ; in that case, the learning set consists of all observations and predictions are made on the learning set. |
| eigengenes | Should the training be performed be in the space of eigengenes obtained from a singular value decomposition of the Gene expression data matrix ? Default is <code>FALSE</code> ; in this case, variable selection is necessary to reduce the number of weights that have to be optimized. |
| models | a logical value indicating whether the model object shall be returned |
| ... | Further arguments passed to the function <code>nnet</code> from the package of the same name.
Important parameters are: <ul style="list-style-type: none"> • "size", i.e. the number of units in the hidden layer • "decay" for weight decay. |

Value

An object of class `cloutput`.

Note

- Excessive variable selection is usually necessary if `eigengenes = FALSE`
- Different runs of this method on the same dataset not necessarily produce the same results due to the fact that optimization for Feed-Forward Neural Networks is rather difficult and depends on the choice of (normally randomly chosen) starting values for the network weights.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

Christoph Bernau <bernau@ibe.med.uni-muenchen.de>

References

Ripley, B.D. (1996)
 Pattern Recognition and Neural Networks.
 Cambridge University Press

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#),
[ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcMA](#),
[pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression from first 10 genes
golubX <- as.matrix(golub[,2:11])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run nnet (not tuned)
nnetresult <- nnetCMA(X=golubX, y=golubY, learnind=learnind, size = 3, decay = 0.01)
### show results
show(nnetresult)
ftable(nnetresult)
plot(nnetresult)
### in the space of eigengenes (not tuned)
golubXfull <- as.matrix(golubX[,-1])
nnetresult <- nnetCMA(X=golubXfull, y=golubY, learnind = learnind, eigengenes = TRUE,
                    size = 3, decay = 0.01)

### show results
show(nnetresult)
ftable(nnetresult)
plot(nnetresult)
```

obsinfo

*Classifiability of observations***Description**

Some observations are harder to classify than others. It is frequently of interest to know which observations are consistently misclassified; these are candidates for outliers or wrong class labels.

Arguments

object	An object of class <code>evaluation</code> , generated with <code>scheme = "observationwise"</code>
threshold	threshold value of (observation-wise) performance measure, s. <code>evaluation</code> that has to be exceeded in order to speak of consistent misclassification. If <code>measure = "average probability"</code> , then values <i>below</i> threshold are regarded as consistent misclassification. Note that the default values 1 is not sensible in that case
show	Should the information be printed ? Default is TRUE.

Details

As not all observation must have been classified at least once, observations not classified at all are also shown.

Value

A list with two components

`misclassification`

A `data.frame` containing the indices of consistently misclassified observations and the corresponding performance measure.

`notclassified`

The indices of those observations not classified at all, s. details.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

See Also

`evaluation`

pknnCMA-methods *Probabilistic nearest neighbours*

Description

Nearest neighbour variant that replaces the simple voting scheme by a weighted one (based on euclidean distances). This is also used to compute class probabilities.

Methods

X = "matrix", **y** = "numeric", **f** = "missing" signature 1

X = "matrix", **y** = "factor", **f** = "missing" signature 2

X = "data.frame", **y** = "missing", **f** = "formula" signature 3

X = "ExpressionSet", **y** = "character", **f** = "missing" signature 4

For further argument and output information, consult [pknnCMA](#).

pknnCMA *Probabilistic Nearest Neighbours*

Description

Nearest neighbour variant that replaces the simple voting scheme by a weighted one (based on euclidean distances). This is also used to compute class probabilities.

For S4 class information, see [pknnCMA-methods](#).

Usage

```
pknnCMA(X, y, f, learnind, beta = 1, k = 1, models=FALSE, ...)
```

Arguments

- | | |
|---|--|
| X | Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>. |
| y | Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if X is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if X is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p> |
| f | A two-sided formula, if X is a <code>data.frame</code> . The left part correspond to class labels, the right to variables. |

learnind	An index vector specifying the observations that belong to the learning set. Must not be missing for this method.
beta	Slope parameter for the logistic function which is used for the computation of class probabilities. The default value (1) need not produce reasonable results and can produce warnings.
k	Number of nearest neighbours to use.
models	a logical value indicating whether the model object shall be returned
...	Currently unused argument.

Details

The algorithm is as follows:

- Determine the k nearest neighbours
- For each class represented among these, compute the average euclidean distance.
- The negative distances are plugged into the logistic function with parameter β .
- Classify into the class with highest probability.

Value

An object of class `cloutput`.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression from first 10 genes
golubX <- as.matrix(golub[,-1])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run probabilistic k-nearest neighbours
result <- pknnCMA(X=golubX, y=golubY, learnind=learnind, k = 3)
### show results
show(result)
ftable(result)
plot(result)
```

`plot`*Probability plot*

Description

A popular way of visualizing the output of classifier is to plot, separately for each class, the predicted probability of each predicted observations for the respective class. For this purpose, the plot area is divided into K parts, where K is the number of classes. Predicted observations are assigned, according to their true class, to one of those parts. Then, for each part and each predicted observation, the predicted probabilities are plotted, displayed by coloured dots, where each colour corresponds to one class.

Arguments

<code>x</code>	An object of class <code>cloutput</code> whose slot <code>probmatrix</code> does not contain any missing value, i.e. probability estimations are provided by the classifier.
<code>main</code>	A title for the plot (character).

Value

No return.

Note

The plot usually only makes sense if a sufficiently large numbers of observations has been classified. This is usually achieved by running the classifier on several `learningsets` with the method `classification`. The output can then be processed via `join` to obtain an object of class `cloutput` to which this method can be applied.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

See Also

`cloutput`

Barplot

Barplot of variable importance

Description

This method can be seen as a visual pendant to [toplist](#). The plot visualizes variable importance by a barplot. The height of the barplots correspond to variable importance. What variable importance exactly means depends on the method chosen when calling [GeneSelection](#), s. [genesel](#).

Arguments

<code>x</code>	An object of class genesel
<code>top</code>	Number of top genes whose variable importance should be displayed. Defaults to 10.
<code>iter</code>	Iteration number (<code>learningset</code>) for which variable importance should be displayed.
<code>...</code>	Further graphical options passed to <code>barplot</code> .

Value

No return.

Note

Note the following

- If `scheme = "multiclass"`, only one plot will be made. Otherwise, one plot will be made for each binary scenario (depending on whether `scheme` is `"one-vs-all"` or `"pairwise"`).
- Variable importance do not make sense for variable selection (ranking) methods that are essentially discrete, such as the Wilcoxon-Rank sum statistic or the Kruskal-Wallis statistic.
- For the methods `"lasso"`, `"elasticnet"`, `"boosting"` the number of nonzero coefficients can be very small, resulting in bars of height zero if `top` has been chosen too large.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

See Also

[genesel](#), [GeneSelection](#), [toplist](#)

plot tuningresult *Visualize results of tuning*

Description

After hyperparameter tuning using `tune` it is useful to see which choice of hyperparameters is suitable and how good the performance is.

Arguments

<code>x</code>	An object of class <code>tuningresult</code> .
<code>iter</code>	Iteration number (<code>learningset</code>) for which tuning results should be displayed.
<code>which</code>	Character vector (maximum length is two) naming the arguments for which tuning results should be display. Default is <code>NULL</code> ; if the number of tuned hyperparameter is less or equal than two, then the results for these hyperparameters will be plotted. If this number is two, then a <code>contour</code> plot will be made, otherwise a simple line segment plot. If the number of tuned hyperparameters exceeds two, then <code>which</code> may not be <code>NULL</code> .
<code>...</code>	Further graphical options passed either to <code>plot</code> or <code>contour</code> .

Value

no return.

Note

Frequently, several hyperparameter (combinations) perform "best", s. also the remark in `best`.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

See Also

`tune`, `tuningresult`

plrCMA-methods *L2 penalized logistic regression*

Description

High dimensional logistic regression combined with an L2-type (Ridge-)penalty. Multiclass case is also possible.

Methods

X = "matrix", y = "numeric", f = "missing" signature 1

X = "matrix", y = "factor", f = "missing" signature 2

X = "data.frame", y = "missing", f = "formula" signature 3

X = "ExpressionSet", y = "character", f = "missing" signature 4

For further argument and output information, consult [plrCMA](#).

plrCMA *L2 penalized logistic regression*

Description

High dimensional logistic regression combined with an L2-type (Ridge-)penalty. Multiclass case is also possible. For S4 method information, see [plrCMA-methods](#)

Usage

```
plrCMA(X, y, f, learnind, lambda = 0.01, scale = TRUE, models=FALSE,...)
```

Arguments

- | | |
|---|--|
| X | Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <i>matrix</i>. Rows correspond to observations, columns to variables. • A <i>data.frame</i>, when <i>f</i> is <i>not</i> missing (s. below). • An object of class <i>ExpressionSet</i>. |
| y | Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if X is an <i>ExpressionSet</i> that specifies the phenotype variable. • <i>missing</i>, if X is a <i>data.frame</i> and a proper formula <i>f</i> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p> |
| f | A two-sided formula, if X is a <i>data.frame</i> . The left part correspond to class labels, the right to variables. |

learnind	An index vector specifying the observations that belong to the learning set. May be missing; in that case, the learning set consists of all observations and predictions are made on the learning set.
lambda	Parameter governing the amount of penalization. This hyperparameter should be tuned .
scale	Scale the predictors as specified by X to have unit variance and zero mean.
models	a logical value indicating whether the model object shall be returned
...	Currently unused argument.

Value

An object of class [cloutput](#).

Author(s)

Special thanks go to

Ji Zhu (University of Ann Arbor, Michigan)

Trevor Hastie (Stanford University)

who provided the basic code that was then adapted by

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>.

References

Zhu, J., Hastie, T. (2004). Classification of gene microarrays by penalized logistic regression. *Biostatistics* 5:427-443.

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression from first 10 genes
golubX <- as.matrix(golub[,-1])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run penalized logistic regression (no tuning)
plrresult <- plrCMA(X=golubX, y=golubY, learnind=learnind)
### show results
show(plrresult)
ftable(plrresult)
plot(plrresult)
### multiclass example:
```

```

### load Khan data
data(khan)
### extract class labels
khanY <- khan[,1]
### extract gene expression from first 10 genes
khanX <- as.matrix(khan[,-1])
### select learningset
set.seed(111)
learnind <- sample(length(khanY), size=floor(ratio*length(khanY)))
### run penalized logistic regression (no tuning)
plrresult <- plrCMA(X=khanX, y=khanY, learnind=learnind)
### show results
show(plrresult)
ftable(plrresult)
plot(plrresult)

```

pls_ldaCMA-methods *Partial Least Squares combined with Linear Discriminant Analysis*

Description

-This method constructs a classifier that extracts Partial Least Squares components that are plugged into Linear Discriminant Analysis. The Partial Least Squares components are computed by the package `pls-genomics`.

Methods

X = "matrix", y = "numeric", f = "missing" signature 1

X = "matrix", y = "factor", f = "missing" signature 2

X = "data.frame", y = "missing", f = "formula" signature 3

X = "ExpressionSet", y = "character", f = "missing" signature 4

For further argument and output information, consult [pls_ldaCMA](#).

pls_ldaCMA

Partial Least Squares combined with Linear Discriminant Analysis

Description

This method constructs a classifier that extracts Partial Least Squares components that are plugged into Linear Discriminant Analysis. The Partial Least Squares components are computed by the package `pls-genomics`.

For S4 method information, see [pls_ldaCMA-methods](#).

Usage

```
pls_ldaCMA(X, y, f, learnind, comp = 2, plot = FALSE, models=FALSE)
```

Arguments

<code>X</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p>
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be <code>missing</code> ; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>comp</code>	Number of Partial Least Squares components to extract. Default is 2 which can be suboptimal, depending on the particular dataset. Can be optimized using tune .
<code>plot</code>	If <code>comp</code> \leq 2, should the classification space of the Partial Least Squares components be plotted? Default is <code>FALSE</code> .
<code>models</code>	a logical value indicating whether the model object shall be returned

Value

An object of class `cloutput`.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Nguyen, D., Rocke, D. M., (2002).

Tumor classification by partial least squares using microarray gene expression data.

Bioinformatics 18, 39-50

Boulesteix, A.L., Strimmer, K. (2007).

Partial least squares: a versatile tool for the analysis of high-dimensional genomic data.

Briefings in Bioinformatics 7:32-44.

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```

### load Khan data
data(khan)
### extract class labels
khanY <- khan[,1]
### extract gene expression
khanX <- as.matrix(khan[,-1])
### select learningset
set.seed(111)
learnind <- sample(length(khanY), size=floor(2/3*length(khanY)))
### run Shrunken Centroids classifier, without tuning
plsresult <- pls_ldaCMA(X=khanX, y=khanY, learnind=learnind, comp = 4)
### show results
show(plsresult)
ftable(plsresult)
plot(plsresult)

```

pls_lrCMA-methods *Partial Least Squares followed by logistic regression*

Description

This method constructs a classifier that extracts Partial Least Squares components that form the the covariates in a binary logistic regression model. The Partial Least Squares components are computed by the package `plsgenomics`.

Methods

X = "matrix", y = "numeric", f = "missing" signature 1
X = "matrix", y = "factor", f = "missing" signature 2
X = "data.frame", y = "missing", f = "formula" signature 3
X = "ExpressionSet", y = "character", f = "missing" signature 4
 For further argument and output information, consult [pls_lrCMA](#)

pls_lrCMA *Partial Least Squares followed by logistic regression*

Description

This method constructs a classifier that extracts Partial Least Squares components that form the the covariates in a binary logistic regression model. The Partial Least Squares components are computed by the package `plsgenomics`.

For S4 method information, see [pls_lrCMA-methods](#).

Usage

```
pls_lrCMA(X, y, f, learnind, comp = 2, lambda = 1e-4, plot = FALSE, models=FALSE)
```

Arguments

<code>x</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p>
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be <code>missing</code> ; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>comp</code>	Number of Partial Least Squares components to extract. Default is 2 which can be suboptimal, depending on the particular dataset. Can be optimized using tune .
<code>lambda</code>	Parameter controlling the amount of L2 penalization for logistic regression, usually taken to be a small value in order to stabilize estimation in the case of separable data.
<code>plot</code>	If <code>comp</code> \leq 2, should the classification space of the Partial Least Squares components be plotted? Default is <code>FALSE</code> .
<code>models</code>	a logical value indicating whether the model object shall be returned

Value

An object of class `cloutput`.

Note

Up to now, only the two-class case is supported.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Boulesteix, A.L., Strimmer, K. (2007).

Partial least squares: a versatile tool for the analysis of high-dimensional genomic data.

Briefings in Bioinformatics 7:32-44.

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_rfcMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression
golubX <- as.matrix(golub[,-1])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run PLS, combined with logistic regression
result <- pls_lrCMA(X=golubX, y=golubY, learnind=learnind)
### show results
show(result)
ftable(result)
plot(result)
```

pls_rfcMA-methods *Partial Least Squares followed by random forests*

Description

This method constructs a classifier that extracts Partial Least Squares components used to generate Random Forests, s. [rfCMA](#). The Partial Least Squares components are computed by the package [plsgenomics](#).

Methods

X = "matrix", y = "numeric", f = "missing" signature 1

X = "matrix", y = "factor", f = "missing" signature 2

X = "data.frame", y = "missing", f = "formula" signature 3

X = "ExpressionSet", y = "character", f = "missing" signature 4

For further argument and output information, consult [pls_rfcMA](#).

pls_rfCMA

*Partial Least Squares followed by random forests***Description**

This method constructs a classifier that extracts Partial Least Squares components used to generate Random Forests, s. [rfCMA](#).

For S4 method information, see [pls_rfCMA-methods](#).

Usage

```
pls_rfCMA(X, y, f, learnind, comp = 2 * nlevels(as.factor(y)), seed = 111, models
```

Arguments

X	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>.
y	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if X is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if X is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p>
f	A two-sided formula, if X is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
learnind	An index vector specifying the observations that belong to the learning set. May be <code>missing</code> ; in that case, the learning set consists of all observations and predictions are made on the learning set.
comp	Number of Partial Least Squares components to extract. Default ist two times the number of different classes.
seed	Fix Random number generator seed to <code>seed</code> . This is useful to guarantee reproducibility of the results, due to the random component in the random Forest.
models	a logical value indicating whether the model object shall be returned
...	Further arguments to be passed to <code>randomForests</code> from the package of the same name.

Value

An object of class `cloutput`.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Boulesteix, A.L., Strimmer, K. (2007).

Partial least squares: a versatile tool for the analysis of high-dimensional genomic data.

Briefings in Bioinformatics 7:32-44.

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression
golubX <- as.matrix(golub[,-1])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run PLS, combined with Random Forest
result <- pls_rfCMA(X=golubX, y=golubY, learnind=learnind)
### show results
show(result)
ftable(result)
plot(result)
```

pnnCMA-methods

Probabilistic Neural Networks

Description

Probabilistic Neural Networks is the term Specht (1990) used for a Gaussian kernel estimator for the conditional class densities.

Methods

X = "matrix", **y** = "numeric", **f** = "missing" signature 1

X = "matrix", **y** = "factor", **f** = "missing" signature 2

X = "data.frame", **y** = "missing", **f** = "formula" signature 3

X = "ExpressionSet", **y** = "character", **f** = "missing" signature 4

For references, further argument and output information, consult [pnnCMA](#).

pnnCMA

*Probabilistic Neural Networks***Description**

Probabilistic Neural Networks is the term Specht (1990) used for a Gaussian kernel estimator for the conditional class densities.

For S4 method information, see [pnnCMA-methods](#).

Usage

```
pnnCMA(X, y, f, learnind, sigma = 1, models=FALSE)
```

Arguments

X	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A matrix. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>. Each variable (gene) will be scaled for unit variance and zero mean.
y	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if X is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if X is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p>
f	A two-sided formula, if X is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
learnind	An index vector specifying the observations that belong to the learning set. For this method, this must <i>not</i> be <code>missing</code> .
sigma	Standard deviation of the Gaussian Kernel used. This hyperparameter should be tuned, s. tune . The default is 1, but this generally does not lead to good results. Actually, this method reacts very sensitively to the value of sigma. Take care if warnings appear related to the particular choice.
models	a logical value indicating whether the model object shall be returned

Value

An object of class `cloutput`.

Note

There is actually no strong relation of this method to Feed-Forward Neural Networks, s. [nnetCMA](#).

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Specht, D.F. (1990).

Probabilistic Neural Networks. *Neural Networks*, 3, 109-118.

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression from first 10 genes
golubX <- as.matrix(golub[,2:11])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run PNN
pnnresult <- pnnCMA(X=golubX, y=golubY, learnind=learnind, sigma = 3)
### show results
show(pnnresult)
ftable(pnnresult)
plot(pnnresult)
```

prediction-methods *General method for predicting class labels of new observations*

Description

Perform prediction signatures:

Methods

X.tr = "matrix", X.new="matrix", y.tr='any', f = "missing" signature 1

X.tr = "data.frame", X.new="data.frame", y.tr = "missing", f = "formula" signature 2

X.tr = "ExpressionSet", X.new = "ExpressionSet", y.tr = "character", f = "missing" signature 3

For further argument and output information, consult [classification](#).

prediction

General method for predicting classes of new observations

Description

This method constructs the given classifier using the specified training data, gene selection and tuning results.. Subsequently, class labels are predicted for new observations.

For S4 method information, s. [classification-methods](#).

Usage

```
prediction(X.tr, y.tr, X.new, f, classifier, genesel, models=F, nbgene, tuner, ...)
```

Arguments

- | | |
|----------------------|--|
| <code>X.tr</code> | <p>Training gene expression data. Can be one of the following:</p> <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>. |
| <code>X.new</code> | <p>gene expression data. Can be one of the following:</p> <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>. |
| <code>y.tr</code> | <p>Class labels of training observation. Can be one of the following:</p> <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • missing, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded for classifier construction to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p> |
| <code>f</code> | <p>A two-sided formula, if <code>X</code> is a <code>data.frame</code>. The left part correspond to class labels, the right to variables.</p> |
| <code>genesel</code> | <p>Optional (but usually recommended) object of class <code>genesel</code> containing variable importance information for the argument <code>learningsets</code>. In this case the object contains a single variable selection. Appropriate <code>genesel</code>-objects can be obtained using the function <code>genesel</code> without <code>learningset</code> and setting <code>X=X.tr</code> and <code>y=y.tr</code> (i.e. corresponding to the training data of this function).</p> |
| <code>nbgene</code> | <p>Number of best genes to be kept for classification, based on either <code>genesel</code> or the call to <code>GeneSelection</code> using <code>genesellist</code>. In the case that both are missing, this argument is not necessary. note:</p> <ul style="list-style-type: none"> • If the gene selection method has been one of "lasso", "elasticnet", "boosting", <code>nbgene</code> will be reset to $\min(s, \text{nbgene})$ where <code>s</code> is the number of nonzero coefficients. |

- if the gene selection scheme has been "one-vs-all", "pairwise" for the multiclass case, there exist several rankings. The top `nbgene` will be kept of *each* of them, so the number of effective used genes will sometimes be much larger.

<code>classifier</code>	Name of function ending with <code>CMA</code> indicating the classifier to be used.
<code>tuneres</code>	Analogous to the argument <code>genesel</code> - object of class <code>tuningresult</code> containing information about the best hyperparameter choice for the argument <code>learningsets</code> . Appropriate tuning-objects can be obtained using the function <code>tune</code> without <code>learningsets</code> and setting parameters <code>X=X.tr</code> , <code>y=y.tr</code> and <code>genesel=genesel</code> (i.e. using the same training data and gene selection as in this function)
<code>models</code>	a logical value indicating whether the model object shall be returned
<code>...</code>	Further arguments passed to the function <code>classifier</code> .

Details

This function builds the specified classifier and predicts the class labels of new observations. Hence, its usage differs from those of most other prediction functions in R.

Value

A object of class `predoutput-class`; Predicted classes can be seen by `show(predoutput)`

Author(s)

Christoph Bernau <berнау@ibe.med.uni-muenchen.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

See Also

[GeneSelection](#), [tune](#), [evaluation](#), [compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMAclassification](#)

Examples

```
### a simple k-nearest neighbour example
### datasets
## Not run: plot(x)
data(golub)
golubY <- golub[,1]
golubX <- as.matrix(golub[,-1])
### Splitting data into training and test set
X.tr<-golubX[1:30]
X.new<-golubX[31:39]
y.tr<-golubY[1:30]
### 1. GeneSelection
seltest <- GeneSelection(X=X.tr, y=y.tr, method = "t.test")
### 2. tuning
```

```
tunek <- tune(X.tr, y.tr, genesel = selttest, nbgene = 20, classifier = knnCMA)
### 3. classification
pred <- prediction(X.tr=X.tr,y.tr=y.tr,X.new=X.new, genesel = selttest,
                  tunereres = tunek, nbgene = 20, classifier = knnCMA)
### show and analyze results:
show(pred)

## End(Not run)
```

predoutput-class *"predoutput"*

Description

Object returned by the function `prediction`

Slots

`Xnew`: Gene Expression matrix of new observations

`yhat`: Predicted class labels for the new data.

`model`: List containing the constructed classifier.

Methods

show Returns predicted class labels for the new data.

Author(s)

Christoph Bernau <berнау@ibe.med.uni-muenchen.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#),
[ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcMA](#),
[pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

qdaCMA-methods

Quadratic Discriminant Analysis

Description

Performs a quadratic discriminant analysis under the assumption of a multivariate normal distribution in each classes without restriction concerning the covariance matrices. The function `qda` from the package `MASS` is called for computation.

Methods

`X = "matrix", y = "numeric", f = "missing"` signature 1
`X = "matrix", y = "factor", f = "missing"` signature 2
`X = "data.frame", y = "missing", f = "formula"` signature 3
`X = "ExpressionSet", y = "character", f = "missing"` signature 4
 For further argument and output information, consult [qdaCMA](#).

 qdaCMA

Quadratic Discriminant Analysis

Description

Performs a quadratic discriminant analysis under the assumption of a multivariate normal distribution in each classes without restriction concerning the covariance matrices. The function `qda` from the package `MASS` is called for computation.

For S4 method information, see [qdaCMA-methods](#).

Usage

```
qdaCMA(X, y, f, learnind, models=FALSE, ...)
```

Arguments

<code>X</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A <code>factor</code>. • A <code>character</code> if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p>
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be <code>missing</code> ; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>models</code>	a logical value indicating whether the model object shall be returned
<code>...</code>	Further arguments to be passed to <code>qda</code> from the package <code>MASS</code>

Value

An object of class `cloutput`.

Note

Excessive variable selection has usually to be performed before qdaCMA can be applied in the $p > n$ setting. Not reducing the number of variables can result in an error message.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

McLachlan, G.J. (1992).

Discriminant Analysis and Statistical Pattern Recognition.

Wiley, New York

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [pnnCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression from first 3 genes
golubX <- as.matrix(golub[,2:4])
### select learningset
ratio <- 2/3
set.seed(112)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run QDA
qdaresult <- qdaCMA(X=golubX, y=golubY, learnind=learnind)
### show results
show(qdaresult)
ftable(qdaresult)
plot(qdaresult)
### multiclass example:
### load Khan data
data(khan)
### extract class labels
khanY <- khan[,1]
### extract gene expression from first 4 genes
khanX <- as.matrix(khan[,2:5])
### select learningset
set.seed(111)
learnind <- sample(length(khanY), size=floor(ratio*length(khanY)))
### run QDA
qdaresult <- qdaCMA(X=khanX, y=khanY, learnind=learnind)
### show results
show(qdaresult)
ftable(qdaresult)
plot(qdaresult)
```

Description

Random Forests were proposed by Breiman (2001) and are implemented in the package `randomForest`. In this package, they can as well be used to rank variables according to their importance, s. `GeneSelection`.

Methods

X = "matrix", y = "numeric", f = "missing" signature 1

X = "matrix", y = "factor", f = "missing" signature 2

X = "data.frame", y = "missing", f = "formula" signature 3

X = "ExpressionSet", y = "character", f = "missing" signature 4

For references, further argument and output information, consult [rfCMA](#)

Description

Random Forests were proposed by Breiman (2001) and are implemented in the package `randomForest`. In this package, they can as well be used to rank variables according to their importance, s. `GeneSelection`. For S4 method information, see [rfCMA-methods](#)

Usage

```
rfCMA(X, y, f, learnind, varimp = TRUE, seed = 111, models=FALSE, ...)
```

Arguments

- | | |
|---|--|
| X | Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>. |
| y | Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if X is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if X is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p> |
| f | A two-sided formula, if X is a <code>data.frame</code> . The left part correspond to class labels, the right to variables. |

<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be missing; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>varimp</code>	Should variable importance measures be computed ? Defaults to TRUE.
<code>seed</code>	Fix Random number generator seed to <code>seed</code> . This is useful to guarantee reproducibility of the results.
<code>models</code>	a logical value indicating whether the model object shall be returned
<code>...</code>	Further arguments to be passed to <code>randomForest</code> from the package of the same name.

Value

If `varimp`, then an object of class `clvarseloutput` is returned, otherwise an object of class `cloutput`

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Breiman, L. (2001)
 Random Forest.
Machine Learning, 45:5-32.

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfCMA](#), [pnnCMA](#), [qdaCMA](#), [scdaCMA](#), [shrinkldaCMA](#), [svmCMA](#)

Examples

```
### load Khan data
data(khan)
### extract class labels
khanY <- khan[,1]
### extract gene expression
khanX <- as.matrix(khan[,-1])
### select learningset
set.seed(111)
learnind <- sample(length(khanY), size=floor(2/3*length(khanY)))
### run random Forest
rfresult <- rfCMA(X=khanX, y=khanY, learnind=learnind, varimp = FALSE)
### show results
show(rfresult)
ftable(rfresult)
plot(rfresult)
```

`roc`*Receiver Operator Characteristic*

Description

The empirical Receiver Operator Characteristic (ROC) is widely used for the evaluation of diagnostic tests, but also for the evaluation of classifiers. In this implementation, it can only be used for the binary classification case. The input are a numeric vector of class probabilities (which play the role of a test result) and the true class labels. Note that misclassification performance can (partly widely) differ from the Area under the ROC (AUC). This is due to the fact that misclassification rates are always computed for the threshold 'probability = 0.5'.

Arguments

<code>object</code>	An object of <code>cloutput</code> .
<code>plot</code>	Should the ROC curve be plotted ? Default is TRUE.
<code>...</code>	Argument to specify further graphical options.

Value

The empirical area under the curve (AUC).

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

See Also

[evaluation](#)

`scdaCMA-methods`*Shrunken Centroids Discriminant Analysis*

Description

The nearest shrunken centroid classification algorithm is detailly described in Tibshirani et al. (2002).

It is widely known under the name PAM (prediction analysis for microarrays), which can also be found in the package `pamr`.

Methods

`X = "matrix", y = "numeric", f = "missing"` signature 1

`X = "matrix", y = "factor", f = "missing"` signature 2

`X = "data.frame", y = "missing", f = "formula"` signature 3

`X = "ExpressionSet", y = "character", f = "missing"` signature 4

For references, further argument and output information, consult [scdaCMA](#).

scdaCMA

Shrunken Centroids Discriminant Analysis

Description

The nearest shrunken centroid classification algorithm is detailly described in Tibshirani et al. (2002).

It is widely known under the name PAM (prediction analysis for microarrays), which can also be found in the package `pamr`.

For S4 method information, see [scdaCMA-methods](#).

Usage

```
scdaCMA(X, y, f, learnind, delta = 0.5, models=FALSE, ...)
```

Arguments

<code>X</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A <code>factor</code>. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p>
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be <code>missing</code> ; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>delta</code>	The shrinkage intensity for the class centroids - a hyperparameter that must be tuned. The default <code>0.5</code> not necessarily produces good results.
<code>models</code>	a logical value indicating whether the model object shall be returned
<code>...</code>	Currently unused argument.

Value

An object of class `cloutput`.

Note

The results can differ from those obtained by using the package `pamr`.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Tibshirani, R., Hastie, T., Narasimhan, B., and Chu, G., (2003).

Class prediction by nearest shrunken centroids with applications to DNA microarrays.

Statistical Science, 18, 104-117

See Also

`compBoostCMA`, `dldaCMA`, `ElasticNetCMA`, `fdaCMA`, `flexdaCMA`, `gbmCMA`, `knnCMA`,
`ldaCMA`, `LassoCMA`, `nnetCMA`, `pknnCMA`, `plrCMA`, `pls_ldaCMA`, `pls_lrCMA`, `pls_rfcCMA`,
`pnnCMA`, `qdaCMA`, `rfCMA`, `shrinkldaCMA`, `svmCMA`

Examples

```
### load Khan data
data(khan)
### extract class labels
khanY <- khan[,1]
### extract gene expression
khanX <- as.matrix(khan[,-1])
### select learningset
set.seed(111)
learnind <- sample(length(khanY), size=floor(2/3*length(khanY)))
### run Shrunken Centroids classifier, without tuning
scdaresult <- scdaCMA(X=khanX, y=khanY, learnind=learnind)
### show results
show(scdaresult)
ftable(scdaresult)
plot(scdaresult)
```

shrinkldaCMA-methods

Shrinkage linear discriminant analysis

Description

Linear Discriminant Analysis combined with the James-Stein-Shrinkage approach of Schaefer and Strimmer (2005) for the covariance matrix.

Currently still an experimental version. For S4 method information, see [shrinkldaCMA-methods](#)

Methods

`X = "matrix", y = "numeric", f = "missing"` signature 1
`X = "matrix", y = "factor", f = "missing"` signature 2
`X = "data.frame", y = "missing", f = "formula"` signature 3
`X = "ExpressionSet", y = "character", f = "missing"` signature 4
 For further argument and output information, consult [shrinkldaCMA](#).

shrinkldaCMA

Shrinkage linear discriminant analysis

Description

Linear Discriminant Analysis combined with the James-Stein-Shrinkage approach of Schaefer and Strimmer (2005) for the covariance matrix.

Currently still an experimental version.

For S4 method information, see [shrinkldaCMA-methods](#)

Usage

```
shrinkldaCMA(X, y, f, learnind, models=FALSE, ...)
```

Arguments

<code>X</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p>
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be <code>missing</code> ; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>models</code>	a logical value indicating whether the model object shall be returned
<code>...</code>	Further arguments to be passed to <code>cov.shrink</code> from the package <code>corpcor</code>

Value

An object of class `cloutput`.

Note

This is still an experimental version.

Covariance shrinkage is performed by calling functions from the package `corpcor`.

Variable selection is *not* necessary.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Schaefer, J., Strimmer, K. (2005).

A shrinkage approach to large-scale covariance estimation and implications for functional genomics.

Statistical Applications in Genetics and Molecular Biology, 4:32.

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [svmCMA](#).

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression
golubX <- as.matrix(golub[,-1])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run shrinkage-LDA
result <- shrinkldaCMA(X=golubX, y=golubY, learnind=learnind)
### show results
show(result)
ftable(result)
plot(result)
```

summary

Summarize classifier evaluation

Description

This method principally does nothing more than applying the pre-implemented `summary()` function to the slot `score` of an object of class `evaloutput`. One then obtains the usual five-point-summary, consisting of minimum and maximum, lower and upper quartile and the median. Additionally, the mean is also shown.

Arguments

`object` An object of class `evaloutput`.
`...` Further arguments passed to the pre-implemented `summary` function.

Value

No return.

Note

That the results normally differ for different evaluation schemes ("iterationwise" or "observation-wise").

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

See Also

[evaluation](#), [compare](#), [obsinfo](#).

svmCMA-methods

Support Vector Machine

Description

Calls the function `svm` from the package `e1071` that provides an interface to the award-winning LIBSVM routines.

Methods

`X = "matrix", y = "numeric", f = "missing"` signature 1

`X = "matrix", y = "factor", f = "missing"` signature 2

`X = "data.frame", y = "missing", f = "formula"` signature 3

`X = "ExpressionSet", y = "character", f = "missing"` signature 4

For further argument and output information, consult [svmCMA](#).

Description

Calls the function `svm` from the package `e1071` that provides an interface to the award-winning LIBSVM routines. For S4 method information, see [svmCMA-methods](#)

Usage

```
svmCMA(X, y, f, learnind, probability, models=FALSE, ...)
```

Arguments

<code>X</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided. <p>WARNING: The class labels will be re-coded to range from 0 to $K-1$, where K is the total number of different classes in the learning set.</p>
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learnind</code>	An index vector specifying the observations that belong to the learning set. May be <code>missing</code> ; in that case, the learning set consists of all observations and predictions are made on the learning set.
<code>probability</code>	logical indicating whether the model should allow for probability predictions.
<code>models</code>	a logical value indicating whether the model object shall be returned
<code>...</code>	Further arguments to be passed to <code>svm</code> from the package <code>e1071</code>

Value

An object of class `cloutput`.

Note

Contrary to the default settings in `e1071:::svm`, the used kernel is a linear kernel which has turned to be out a better default setting in the small sample, large number of predictors - situation, because additional nonlinearity is mostly not necessary there. It additionally avoids the tuning of a further kernel parameter `gamma`, s. help of the package `e1071` for details.

Nevertheless, hyperparameter tuning concerning the parameter `cost` must usually be performed to obtain reasonable results, s. [tune](#).

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

Christoph Bernau <bernau@ibe.med.uni-muenchen.de>

References

Boser, B., Guyon, I., Vapnik, V. (1992)

A training algorithm for optimal margin classifiers.

Proceedings of the fifth annual workshop on Computational learning theory, pages 144-152, ACM Press.

Chang, Chih-Chung and Lin, Chih-Jen : LIBSVM: a library for Support Vector Machines <http://www.csie.ntu.edu.tw/~cjlin/libsvm>

Schoelkopf, B., Smola, A.J. (2002)

Learning with kernels. *MIT Press, Cambridge, MA.*

See Also

[compBoostCMA](#), [dldaCMA](#), [ElasticNetCMA](#), [fdaCMA](#), [flexdaCMA](#), [gbmCMA](#), [knnCMA](#), [ldaCMA](#), [LassoCMA](#), [nnetCMA](#), [pknnCMA](#), [plrCMA](#), [pls_ldaCMA](#), [pls_lrCMA](#), [pls_rfcCMA](#), [pnnCMA](#), [qdaCMA](#), [rfCMA](#), [scdaCMA](#), [shrinkldaCMA](#)

Examples

```
### load Golub AML/ALL data
data(golub)
### extract class labels
golubY <- golub[,1]
### extract gene expression
golubX <- as.matrix(golub[,-1])
### select learningset
ratio <- 2/3
set.seed(111)
learnind <- sample(length(golubY), size=floor(ratio*length(golubY)))
### run _untuned_linear SVM
svmresult <- svmCMA(X=golubX, y=golubY, learnind=learnind,probability=TRUE)
### show results
show(svmresult)
ftable(svmresult)
plot(svmresult)
```

toplist

Display 'top' variables

Description

This is a convenient method to get quick access to the most important variables, based on the result of call to [GeneSelection](#).

Usage

```
toplist(object, k = 10, iter = 1, show = TRUE, ...)
```

Arguments

<code>object</code>	An object of <code>genesel</code> .
<code>k</code>	Number of top genes for which information should be displayed. Defaults to 10.
<code>iter</code>	iteration number (<code>learningset</code>) for which tuning results should be displayed.
<code>show</code>	Should the results be printed ? Default is <code>TRUE</code> .
<code>...</code>	Currently unused argument.

Value

The type of output depends on the gene selection scheme. For the multiclass case, if gene selection has been run with the "pairwise" or "one-vs-all" scheme, then the output will be a list of `data.frames`, each containing the gene indices plus variable importance for the top `k` genes. The list elements are named according to the binary scenarios (e.g., 1 vs. 3). Otherwise, a single `data.frame` is returned.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

References

Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

See Also

`genesel`, `GeneSelection`, `plot`, `genesel-method`

tune-methods

Hyperparameter tuning for classifiers

Description

Performs hyperparameter tuning for the following signatures:

Methods

`X = "matrix", y = "numeric", f = "missing"` signature 1

`X = "matrix", y = "factor", f = "missing"` signature 2

`X = "data.frame", y = "missing", f = "formula"` signature 3

`X = "ExpressionSet", y = "character", f = "missing"` signature 4

For further argument and output information, consult `tune`.

Description

Most classifiers implemented in this package depend on one or even several hyperparameters (s. details) that should be optimized to obtain good (and comparable !) results. As tuning scheme, we propose three fold Cross-Validation on each `learningset` (for fixed selected variables). Note that `learningsets` usually do not contain the complete dataset, so tuning involves a second level of splitting the dataset. Increasing the number of folds leads to larger datasets (and possibly to higher accuracy), but also to higher computing times.

For S4 method information, s. `link{tune-methods}`

Usage

```
tune(X, y, f, learningsets, genesel, genesellist = list(), nbgene, classifier, f
```

Arguments

<code>X</code>	Gene expression data. Can be one of the following: <ul style="list-style-type: none"> • A <code>matrix</code>. Rows correspond to observations, columns to variables. • A <code>data.frame</code>, when <code>f</code> is <i>not</i> missing (s. below). • An object of class <code>ExpressionSet</code>.
<code>y</code>	Class labels. Can be one of the following: <ul style="list-style-type: none"> • A numeric vector. • A factor. • A character if <code>X</code> is an <code>ExpressionSet</code> that specifies the phenotype variable. • <code>missing</code>, if <code>X</code> is a <code>data.frame</code> and a proper formula <code>f</code> is provided.
<code>f</code>	A two-sided formula, if <code>X</code> is a <code>data.frame</code> . The left part correspond to class labels, the right to variables.
<code>learningsets</code>	An object of class <code>learningsets</code> . May be missing, then the complete datasets is used as learning set.
<code>genesel</code>	Optional (but usually recommended) object of class <code>genesel</code> containing variable importance information for the argument <code>learningsets</code>
<code>genesellist</code>	In the case that the argument <code>genesel</code> is missing, this is an argument list passed to <code>GeneSelection</code> . If both <code>genesel</code> and <code>genesellist</code> are missing, no variable selection is performed.
<code>nbgene</code>	Number of best genes to be kept for classification, based on either <code>genesel</code> or the call to <code>GeneSelection</code> using <code>genesellist</code> . In the case that both are missing, this argument is not necessary. note: <ul style="list-style-type: none"> • If the gene selection method has been one of "lasso", "elasticnet", "boosting", <code>nbgene</code> will be reset to <code>min(s, nbgene)</code> where <code>s</code> is the number of nonzero coefficients. • if the gene selection scheme has been "one-vs-all", "pairwise" for the multiclass case, there exist several rankings. The top <code>nbgene</code> will be kept of <i>each</i> of them, so the number of effective used genes will sometimes be much larger.

<code>classifier</code>	Name of function ending with <code>CMA</code> indicating the classifier to be used.
<code>fold</code>	The number of cross-validation folds used within each <code>learningset</code> . Default is 3. Increasing <code>fold</code> will lead to higher computing times.
<code>strat</code>	Should stratified cross-validation according to the class proportions in the complete dataset be used ? Default is <code>FALSE</code> .
<code>grids</code>	A named list. The names correspond to the arguments to be tuned, e.g. <code>k</code> (the number of nearest neighbours) for <code>knnCMA</code> , or <code>cost</code> for <code>svmCMA</code> . Each element is a numeric vector defining the grid of candidate values. Of course, several hyperparameters can be tuned simultaneously (though requiring much time). By default, <code>grids</code> is an empty list. In that case, a pre-defined list will be used, s. details.
<code>trace</code>	Should progress be traced ? Default is <code>TRUE</code> .
<code>...</code>	Further arguments to be passed to <code>classifier</code> , of course not one of the arguments to be tuned (!).

Details

The following default settings are used, if the arguments `grids` is an empty list:

```

gbmCMA n.trees = c(50, 100, 200, 500, 1000)
compBoostCMA mstop = c(50, 100, 200, 500, 1000)
LassoCMA norm.fraction = seq(from=0.1, to=0.9, length=9)
ElasticNetCMA norm.fraction = seq(from=0.1, to=0.9, length=5), lambda2
  = 2^{-(5:1)}
plrCMA lambda = 2^{-4:4}
pls_ldaCMA comp = 1:10
pls_lrCMA comp = 1:10
pls_rfcCMA comp = 1:10
rfCMA mtry = ceiling(c(0.1, 0.25, 0.5, 1, 2)*sqrt(ncol(X))), nodesize
  = c(1,2,3)
knnCMA k=1:10
pknnCMA k = 1:10
scdaCMA delta = c(0.1, 0.25, 0.5, 1, 2, 5)
pnnCMA sigma = c(2^{-2:2}),
nnetCMA size = 1:5, decay = c(0, 2^{-(4:1)})
svmCMA,kernel = "linear" cost = c(0.1, 1, 5, 10, 50, 100, 500)
svmCMA,kernel = "radial" cost = c(0.1, 1, 5, 10, 50, 100, 500), gamma
  = 2^{-2:2}
svmCMA,kernel = "polynomial" cost = c(0.1, 1, 5, 10, 50, 100, 500),
  degree = 2:4

```

Value

An object of class `tuningresult`

Note

The computation time can be enormously high. Note that for each different learningset, the classifier must be trained `foldtimes` number of possible different hyperparameter combinations times. E.g. if the number of the learningsets is fifty, `fold = 3` and two hyperparameters (each with 5 candidate values) are tuned, $50 \times 3 \times 25 = 3750$ training iterations are necessary !

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

Christoph Bernau <bernau@ibe.med.uni-muenchen.de>

References

Slawski, M. Daumer, M. Boulesteix, A.-L. (2008) CMA - A comprehensive Bioconductor package for supervised classification with high dimensional data. *BMC Bioinformatics* 9: 439

See Also

[tuningresult](#), [GeneSelection](#), [classification](#)

Examples

```
## Not run:
### simple example for a one-dimensional grid, using compBoostCMA.
### dataset
data(golub)
golubY <- golub[,1]
golubX <- as.matrix(golub[,-1])
### learningsets
set.seed(111)
lset <- GenerateLearningsets(y=golubY, method = "CV", fold=5, strat =TRUE)
### tuning after gene selection with the t.test
tunerres <- tune(X = golubX, y = golubY, learningsets = lset,
                genesellist = list(method = "t.test"),
                classifier=compBoostCMA, nbgene = 100,
                grids = list(mstop = c(50, 100, 250, 500, 1000)))
### inspect results
show(tunerres)
best(tunerres)
plot(tunerres, iter = 3)

## End(Not run)
```

tuningresult-class "*tuningresult*"

Description

Object returned by the function [tune](#)

Slots

hypergrid: A `data.frame` representing the grid of values that were tried and evaluated. The number of columns equals the number of tuned hyperparameters and the number rows equals the number of all possible combinations of the discrete grids.

tuneres: A list whose lengths equals the number of different `learningsets` for which tuning has been performed and whose elements are numeric vectors with length equal to the number of rows of `hypergrid` (s.above), containing the misclassification rate belonging to the respective hyperparameter/hyperparameter combination. In order to get an overview about the best hyperparameter/hyperparameter combination, use the convenience method `best`

method: Name of the classifier that has been tuned.

fold: Number of cross-validation fold used for tuning, s. argument of the same name in `tune`

Methods

show Use `show(tuninresult-object)` for brief information.

best Use `best(tuningresult-object)` to see which hyperparameter/hyperparameter combination has performed best in terms of the misclassification rate, s. `best,tuningresult-method`

plot Use `plot(tuningresult-object, iter, which)` to display the performance of hyperparameter/hyperparameter combinations graphically, either as one-dimensional or as two-dimensional (contour) plot, s. `plot,tuningresult-method`

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

See Also

`tune`

`varseloutput-class "varseloutput"`

Description

An object returned by the functions described in `filter`, usually not created directly by the user.

Slots

varsel: numeric vector of variable importance measures, e.g. absolute of genewise statistics.

Methods

No methods are currently defined.

Author(s)

Martin Slawski <ms@cs.uni-sb.de>

Anne-Laure Boulesteix <boulesteix@ibe.med.uni-muenchen.de>

See Also

[filter](#), [clvarseloutput](#)

Index

*Topic **datasets**

golub, 39
khan, 41

*Topic **multivariate**

Barplot, 54
best, 13
boxplot, 14
classification, 15
classification-methods, 15
cloutput-class, 17
clvarseloutput-class, 18
CMA-package, 1
compare, 23
compare-methods, 22
compBoostCMA, 20
compBoostCMA-methods, 20
dldaCMA, 25
dldaCMA-methods, 25
ElasticNetCMA, 2
ElasticNetCMA-methods, 2
evaloutput-class, 27
evaluation, 28
evaluation-methods, 28
fdaCMA, 31
fdaCMA-methods, 30
filter, 33
flexdaCMA, 34
flexdaCMA-methods, 33
ftable, 36
gbmCMA, 37
gbmCMA-methods, 36
GenerateLearningsets, 7
genesel-class, 38
GeneSelection, 5
GeneSelection-methods, 4
internals, 40
join, 41
join-methods, 40
knnCMA, 43
knnCMA-methods, 42
LassoCMA, 10
LassoCMA-methods, 9
ldaCMA, 45

ldaCMA-methods, 44
learningsets-class, 47
nnetCMA, 48
nnetCMA-methods, 47
obsinfo, 50
pknnCMA, 51
pknnCMA-methods, 51
Planarplot, 12
Planarplot-methods, 11
plot, 53
plot tuningresult, 55
plrCMA, 56
plrCMA-methods, 56
pls_ldaCMA, 58
pls_ldaCMA-methods, 58
pls_lrcCMA, 60
pls_lrcCMA-methods, 60
pls_rfcCMA, 63
pls_rfcCMA-methods, 62
pnnCMA, 65
pnnCMA-methods, 64
prediction, 67
prediction-methods, 66
predoutput-class, 69
qdaCMA, 70
qdaCMA-methods, 69
rfCMA, 72
rfCMA-methods, 72
roc, 74
scdaCMA, 75
scdaCMA-methods, 74
shrinkldaCMA, 77
shrinkldaCMA-methods, 76
summary, 78
svmCMA, 80
svmCMA-methods, 79
toplist, 81
tune, 83
tune-methods, 82
tuningresult-class, 85
varseloutput-class, 86

Barplot, 54
best, 13, 55, 86

- best, tuningresult-method, 86
- best, tuningresult-method (*best*), 13
- bklr (*internals*), 40
- bkreg (*internals*), 40
- boxplot, 14
- boxplot, evaloutput-method, 27
- boxplot, evaloutput-method (*boxplot*), 14

- care.dev (*internals*), 40
- care.exp (*internals*), 40
- characterplot (*internals*), 40
- classification, 1, 7, 9, 15, 15, 23, 24, 29, 41, 47, 53, 66, 68, 85
- classification, data.frame, missing, formula-method (*classification-methods*), 15
- classification, ExpressionSet, character, missing-method (*classification-methods*), 15
- classification, matrix, factor, missing-method (*classification-methods*), 15
- classification, matrix, numeric, missing-method (*classification-methods*), 15
- classification-methods, 15, 67
- classification-methods, 15
- cloutput, 16, 19, 23, 26, 28, 31, 35–37, 40, 41, 43, 45, 48, 52, 53, 57, 59, 61, 63, 65, 70, 73, 74, 76, 77, 80
- cloutput (*cloutput-class*), 17
- cloutput-class, 17
- clvarelooutput, 3, 16, 21, 23, 28, 41, 73, 87
- clvarelooutput (*clvarelooutput-class*), 18
- clvarelooutput-class, 18
- CMA (*CMA-package*), 1
- CMA-package, 1
- compare, 1, 22, 23, 29, 41, 79
- compare, list-method (*compare-methods*), 22
- compare-methods, 23
- compare-methods, 22
- compBoostCMA, 1, 4, 6, 11, 13, 17–19, 20, 20, 26, 32, 35, 38, 44, 46, 49, 52, 57, 59, 62, 64, 66, 68, 69, 71, 73, 76, 78, 81, 84
- compBoostCMA, data.frame, missing, formula-method (*compBoostCMA-methods*), 20
- compBoostCMA, ExpressionSet, character, missing-method (*compBoostCMA-methods*), 20
- compBoostCMA, matrix, factor, missing-method (*compBoostCMA-methods*), 20
- compBoostCMA, matrix, numeric, missing-method (*compBoostCMA-methods*), 20
- compBoostCMA-methods, 20
- compBoostCMA-methods, 20

- dldaCMA, 1, 4, 11, 13, 17, 18, 22, 25, 25, 32, 35, 38, 44, 46, 49, 52, 57, 59, 62, 64, 66, 68, 69, 71, 73, 76, 78, 81
- dldaCMA, data.frame, missing, formula-method (*dldaCMA-methods*), 25
- dldaCMA, ExpressionSet, character, missing-method (*dldaCMA-methods*), 25
- dldaCMA, matrix, factor, missing-method (*dldaCMA-methods*), 25
- dldaCMA, matrix, numeric, missing-method (*dldaCMA-methods*), 25
- dldaCMA-methods, 25, 25

- ElasticNetCMA, 1, 2, 2, 6, 11, 13, 17–19, 22, 26, 32, 35, 38, 44, 46, 49, 52, 57, 59, 62, 64, 66, 68, 69, 71, 73, 76, 78, 81, 84
- ElasticNetCMA, data.frame, missing, formula-method (*ElasticNetCMA-methods*), 2
- ElasticNetCMA, ExpressionSet, character, missing-method (*ElasticNetCMA-methods*), 2
- ElasticNetCMA, matrix, factor, missing-method (*ElasticNetCMA-methods*), 2
- ElasticNetCMA, matrix, numeric, missing-method (*ElasticNetCMA-methods*), 2
- ElasticNetCMA-methods, 2
- evaloutput, 14, 29, 78, 79
- evaloutput (*evaloutput-class*), 27
- evaloutput-class, 27
- evaluation, 1, 14, 16, 17, 23, 24, 27, 28, 28, 36, 41, 50, 68, 74, 79
- evaluation, list-method (*evaluation-methods*), 28
- evaluation-methods, 28
- evaluation-methods, 28

- fdaCMA, 1, 4, 11, 13, 17, 18, 22, 26, 30, 31, 32, 35, 38, 44, 46, 49, 52, 57, 59, 62, 64, 66, 68, 69, 71, 73, 76, 78, 81
- fdaCMA, data.frame, missing, formula-method (*fdaCMA-methods*), 30
- fdaCMA, ExpressionSet, character, missing-method (*fdaCMA-methods*), 30

- learningsets-class, 47
- limmatest (*filter*), 33
- mklr (*internals*), 40
- mkreg (*internals*), 40
- my.care.exp (*internals*), 40
- nnetCMA, 1, 4, 11, 13, 17, 18, 22, 26, 32, 35, 38, 44, 46, 47, 48, 49, 52, 57, 59, 62, 64–66, 68, 69, 71, 73, 76, 78, 81, 84
- nnetCMA, data.frame, missing, formula-method (*nnetCMA-methods*), 47
- nnetCMA, ExpressionSet, character, missing-method (*nnetCMA-methods*), 47
- nnetCMA, matrix, factor, missing-method (*nnetCMA-methods*), 47
- nnetCMA, matrix, numeric, missing-method (*nnetCMA-methods*), 47
- nnetCMA-methods, 47, 48
- obsinfo, 27, 50, 79
- obsinfo, evaloutput-method (*evaloutput-class*), 27
- pknnCMA, 1, 4, 11, 13, 17, 18, 22, 26, 32, 35, 38, 43, 44, 46, 49, 51, 51, 57, 59, 62, 64, 66, 68, 69, 71, 73, 76, 78, 81, 84
- pknnCMA, data.frame, missing, formula-method (*pknnCMA-methods*), 51
- pknnCMA, ExpressionSet, character, missing-method (*pknnCMA-methods*), 51
- pknnCMA, matrix, factor, missing-method (*pknnCMA-methods*), 51
- pknnCMA, matrix, numeric, missing-method (*pknnCMA-methods*), 51
- pknnCMA-methods, 51, 51
- Planarplot, 12, 12
- Planarplot, data.frame, missing, formula-method (*Planarplot-methods*), 11
- Planarplot, ExpressionSet, character, missing-method (*Planarplot-methods*), 11
- Planarplot, matrix, factor, missing-method (*Planarplot-methods*), 11
- Planarplot, matrix, numeric, missing-method (*Planarplot-methods*), 11
- Planarplot-methods, 12
- Planarplot-methods, 11
- plot, 53
- plot tuningresult, 55
- plot, cloutput, missing-method (*plot*), 53
- plot, cloutput-method, 18, 19
- plot, cloutput-method (*plot*), 53
- plot, genesel, missing-method (*Barplot*), 54
- plot, genesel-method, 39, 82
- plot, genesel-method (*Barplot*), 54
- plot, tuningresult, missing-method (*plot tuningresult*), 55
- plot, tuningresult-method, 86
- plot, tuningresult-method (*plot tuningresult*), 55
- plotprob (*internals*), 40
- plrCMA, 1, 4, 11, 13, 17, 18, 22, 26, 32, 35, 38, 44, 46, 49, 52, 56, 56, 59, 62, 64, 66, 68, 69, 71, 73, 76, 78, 81, 84
- plrCMA, data.frame, missing, formula-method (*plrCMA-methods*), 56
- plrCMA, ExpressionSet, character, missing-method (*plrCMA-methods*), 56
- plrCMA, matrix, factor, missing-method (*plrCMA-methods*), 56
- plrCMA, matrix, numeric, missing-method (*plrCMA-methods*), 56
- plrCMA-methods, 56, 56
- pls_ldaCMA, 1, 4, 11, 13, 17, 18, 22, 26, 32, 35, 38, 44, 46, 49, 52, 57, 58, 58, 59, 62, 64, 66, 68, 69, 71, 73, 76, 78, 81, 84
- pls_ldaCMA, data.frame, missing, formula-method (*pls_ldaCMA-methods*), 58
- pls_ldaCMA, ExpressionSet, character, missing-method (*pls_ldaCMA-methods*), 58
- pls_ldaCMA, matrix, factor, missing-method (*pls_ldaCMA-methods*), 58
- pls_ldaCMA, matrix, numeric, missing-method (*pls_ldaCMA-methods*), 58
- pls_ldaCMA-methods, 58
- pls_ldaCMA-methods, 58
- pls_lrCMA, 1, 4, 11, 13, 17, 18, 22, 26, 32, 35, 38, 44, 46, 49, 52, 57, 59, 60, 60, 64, 66, 68, 69, 71, 73, 76, 78, 81, 84
- pls_lrCMA, data.frame, missing, formula-method (*pls_lrCMA-methods*), 60
- pls_lrCMA, ExpressionSet, character, missing-method (*pls_lrCMA-methods*), 60
- pls_lrCMA, matrix, factor, missing-method (*pls_lrCMA-methods*), 60
- pls_lrCMA, matrix, numeric, missing-method (*pls_lrCMA-methods*), 60
- pls_lrCMA-methods, 60
- pls_lrCMA-methods, 60
- pls_rfcMA, 1, 4, 11, 13, 17, 18, 22, 26, 32, 35, 38, 44, 46, 49, 52, 57, 59, 62, 63, 66, 68, 69, 71, 73, 76, 78, 81, 84

- pls_rfcCMA, data.frame, missing, formula-method
 (*pls_rfcCMA-methods*), 62
- pls_rfcCMA, ExpressionSet, character, missing-method
 (*pls_rfcCMA-methods*), 62
- pls_rfcCMA, matrix, factor, missing-method
 (*pls_rfcCMA-methods*), 62
- pls_rfcCMA, matrix, numeric, missing-method
 (*pls_rfcCMA-methods*), 62
- pls_rfcCMA-methods, 63
- pls_rfcCMA-methods, 62
- pnnCMA, 1, 4, 11, 13, 17, 18, 22, 26, 32, 35,
 38, 44, 46, 49, 52, 57, 59, 62, 64, 65,
 68, 69, 71, 73, 76, 78, 81, 84
- pnnCMA, data.frame, missing, formula-method
 (*pnnCMA-methods*), 64
- pnnCMA, ExpressionSet, character, missing-method
 (*pnnCMA-methods*), 64
- pnnCMA, matrix, factor, missing-method
 (*pnnCMA-methods*), 64
- pnnCMA, matrix, numeric, missing-method
 (*pnnCMA-methods*), 64
- pnnCMA-methods, 64, 65
- prediction, 67
- prediction, data.frame, missing, data.frame, formula-method
 (*prediction-methods*), 66
- prediction, ExpressionSet, character, ExpressionSet, missing-method
 (*prediction-methods*), 66
- prediction, matrix, ANY, matrix, missing-method
 (*prediction-methods*), 66
- prediction-methods, 66
- predoutput (*predoutput-class*), 69
- predoutput-class, 68
- predoutput-class, 69
- qdaCMA, 1, 4, 11, 13, 17, 18, 22, 26, 32, 35,
 38, 44, 46, 49, 52, 57, 59, 62, 64, 66,
 68, 69, 70, 70, 73, 76, 78, 81
- qdaCMA, data.frame, missing, formula-method
 (*qdaCMA-methods*), 69
- qdaCMA, ExpressionSet, character, missing-method
 (*qdaCMA-methods*), 69
- qdaCMA, matrix, factor, missing-method
 (*qdaCMA-methods*), 69
- qdaCMA, matrix, numeric, missing-method
 (*qdaCMA-methods*), 69
- qdaCMA-methods, 69, 70
- rfcCMA, 1, 4, 11, 13, 17–19, 22, 26, 32, 35, 38,
 44, 46, 49, 52, 57, 59, 62–64, 66, 68,
 69, 71, 72, 72, 76, 78, 81, 84
- rfcCMA, data.frame, missing, formula-method
 (*rfcCMA-methods*), 72
- rfcCMA, ExpressionSet, character, missing-method
 (*rfcCMA-methods*), 72
- rfcCMA, matrix, factor, missing-method
 (*rfcCMA-methods*), 72
- rfcCMA, matrix, numeric, missing-method
 (*rfcCMA-methods*), 72
- rfcCMA-methods, 72, 72
- rfe (*filter*), 33
- roc, 74
- roc, cloutput-method, 18, 19
- roc, cloutput-method (*roc*), 74
- ROCinternal (*internals*), 40
- roundvector (*internals*), 40
- Rowswaps (*internals*), 40
- safeexp (*internals*), 40
- scdaCMA, 1, 4, 11, 13, 17, 18, 22, 26, 32, 35,
 38, 44, 46, 49, 52, 57, 59, 62, 64, 66,
 68, 69, 71, 73, 75, 75, 78, 81, 84
- scdaCMA, data.frame, missing, formula-method
 (*scdaCMA-methods*), 74
- scdaCMA, ExpressionSet, character, missing-method
 (*scdaCMA-methods*), 74
- scdaCMA, matrix, factor, missing-method
 (*scdaCMA-methods*), 74
- scdaCMA, matrix, numeric, missing-method
 (*scdaCMA-methods*), 74
- scdaCMA-methods, 74, 75
- show, cloutput-method
 (*cloutput-class*), 17
- show, evaloutput-method
 (*evaloutput-class*), 27
- show, genesel-method
 (*genesel-class*), 38
- show, learningsets-method
 (*learningsets-class*), 47
- show, predoutput-method
 (*predoutput-class*), 69
- show, tuningresult-method
 (*tuningresult-class*), 85
- shrinkldaCMA, 1, 4, 11, 13, 17, 18, 22, 26,
 32, 35, 38, 44, 46, 49, 52, 57, 59, 62,
 64, 66, 68, 69, 71, 73, 76, 77, 77, 81
- shrinkldaCMA, data.frame, missing, formula-method
 (*shrinkldaCMA-methods*), 76
- shrinkldaCMA, ExpressionSet, character, missing-method
 (*shrinkldaCMA-methods*), 76
- shrinkldaCMA, matrix, factor, missing-method
 (*shrinkldaCMA-methods*), 76
- shrinkldaCMA, matrix, numeric, missing-method
 (*shrinkldaCMA-methods*), 76
- shrinkldaCMA-methods, 76, 76, 77
- summary, 78

summary, evaloutput-method, 27
summary, evaloutput-method
 (summary), 78
svmCMA, 1, 4, 11, 13, 17, 18, 22, 26, 32, 35,
 38, 44, 46, 49, 52, 57, 59, 62, 64, 66,
 68, 69, 71, 73, 76, 78, 79, 80, 84
svmCMA, data.frame, missing, formula-method
 (svmCMA-methods), 79
svmCMA, ExpressionSet, character, missing-method
 (svmCMA-methods), 79
svmCMA, matrix, factor, missing-method
 (svmCMA-methods), 79
svmCMA, matrix, numeric, missing-method
 (svmCMA-methods), 79
svmCMA-methods, 79, 80

toplist, 39, 54, 81
toplist, genesel-method (toplist),
 81
ttest (filter), 33
tune, 1, 7, 9, 16, 17, 21, 37, 47, 55, 57, 59,
 61, 65, 68, 80, 82, 83, 85, 86
tune, data.frame, missing, formula-method
 (tune-methods), 82
tune, ExpressionSet, character, missing-method
 (tune-methods), 82
tune, matrix, factor, missing-method
 (tune-methods), 82
tune, matrix, numeric, missing-method
 (tune-methods), 82
tune-methods, 82
tuningresult, 13, 16, 55, 68, 84, 85
tuningresult
 (tuningresult-class), 85
tuningresult-class, 85

varseloutput, 19, 33
varseloutput
 (varseloutput-class), 86
varseloutput-class, 86

welchtest (filter), 33
wilcoxtest (filter), 33