

Package: OHPL (via r-universe)

August 20, 2024

Type Package

Title Ordered Homogeneity Pursuit Lasso for Group Variable Selection

Version 1.4.1

Description Ordered homogeneity pursuit lasso (OHPL) algorithm for group variable selection proposed in Lin et al. (2017) <[DOI:10.1016/j.chemolab.2017.07.004](https://doi.org/10.1016/j.chemolab.2017.07.004)>. The OHPL method exploits the homogeneity structure in high-dimensional data and enjoys the grouping effect to select groups of important variables automatically. This feature makes it particularly useful for high-dimensional datasets with strongly correlated variables, such as spectroscopic data.

License GPL-3 | file LICENSE

URL <https://ohpl.io>, <https://ohpl.io/doc/>,
<https://github.com/nanxstats/OHPL>

BugReports <https://github.com/nanxstats/OHPL/issues>

Encoding UTF-8

LazyData true

Depends R (>= 3.0.2)

Imports glmnet, mvtnorm, pls

Roxygen list(markdown = TRUE)

RoxygenNote 7.3.2

Repository <https://nanxstats.r-universe.dev>

RemoteUrl <https://github.com/nanxstats/ohpl>

RemoteRef HEAD

RemoteSha 4046c1273714f98070430538bf69738162806cff

Contents

beer	2
cv.OHPL	3
dlc	4
FOP	5
OHPL	6
OHPL.RMSEP	7
OHPL.sim	8
predict.OHPL	9
proto	10
soil	11
wheat	12
Index	13

beer	<i>The beer dataset</i>
------	-------------------------

Description

The beer dataset contains 60 samples published by Norgaard et al. Recorded with a 30mm quartz cell on the undiluted degassed beer and measured from 1100 to 2250 nm (576 data points) in steps of 2 nm.

Usage

```
data(beer)
```

References

Norgaard, L., Saudland, A., Wagner, J., Nielsen, J. P., Munck, L., & Engelsen, S. B. (2000). Interval partial least-squares regression (iPLS): a comparative chemometric study with an example from near-infrared spectroscopy. *Applied Spectroscopy*, 54(3), 413–419.

Examples

```
data("beer")
x.cal <- beer$xtrain
dim(x.cal)
x.test <- beer$xtest
dim(x.test)
y.cal <- beer$ytrain
dim(y.cal)
y.test <- beer$ytest
dim(y.test)

X <- rbind(x.cal, x.test)
y <- rbind(y.cal, y.test)
n <- nrow(y)
```

```

set.seed(1001)
samp.idx <- sample(1L:n, round(n * 0.7))
X.cal <- X[samp.idx, ]
y.cal <- y[samp.idx]
X.test <- X[-samp.idx, ]
y.test <- y[-samp.idx]

```

cv.OHPL

*Cross-validation for Ordered Homogeneity Pursuit Lasso***Description**

Use cross-validation to help select the optimal number of variable groups and the value of gamma.

Usage

```

cv.OHPL(
  X.cal,
  y.cal,
  maxcomp,
  gamma = seq(0.1, 0.9, 0.1),
  X.test,
  y.test,
  cv.folds = 5L,
  G = 30L,
  type = c("max", "median"),
  scale = TRUE,
  pls.method = "simpls"
)

```

Arguments

X.cal	Predictor matrix (training).
y.cal	Response matrix with one column (training).
maxcomp	Maximum number of components for PLS.
gamma	A vector of the gamma sequence between (0, 1).
X.test	X.test Predictor matrix (test).
y.test	y.test Response matrix with one column (test).
cv.folds	Number of cross-validation folds.
G	Maximum number of variable groups.
type	Find the maximum absolute correlation ("max") or find the median of absolute correlation ("median"). Default is "max".
scale	Should the predictor matrix be scaled? Default is TRUE.
pls.method	Method for fitting the PLS model. Default is "simpls". See the details section in pls::plsr() for all possible options.

Value

A list containing the optimal model, RMSEP, Q2, and other evaluation metrics. Also the optimal number of groups to use in group lasso.

Examples

```
data("wheat")

X <- wheat$x
y <- wheat$protein
n <- nrow(wheat$x)

set.seed(1001)
samp.idx <- sample(1L:n, round(n * 0.7))
X.cal <- X[samp.idx, ]
y.cal <- y[samp.idx]
X.test <- X[-samp.idx, ]
y.test <- y[-samp.idx]

# This could run for a while
## Not run:
cv.fit <- cv.OHPL(
  x, y,
  maxcomp = 6, gamma = seq(0.1, 0.9, 0.1),
  x.test, y.test, cv.folds = 5, G = 30, type = "max"
)
# the optimal G and gamma
cv.fit$opt.G
cv.fit$opt.gamma

## End(Not run)
```

dlc

Compute D, L, and C in the Fisher optimal partitions algorithm

Description

Compute D, L, and C in the Fisher optimal partitions algorithm

Usage

```
dlc(X, maxk)
```

Arguments

X	A set of samples.
maxk	Maximum number of k.

Value

The D, L, and C statistics in the Fisher optimal partitions algorithm.

FOP	<i>Fisher optimal partition</i>
-----	---------------------------------

Description

The Fisher optimal partition algorithm.

Usage

```
FOP(X, k, C)
```

Arguments

X	A set of samples.
k	Number of classes.
C	Statistic from the output of <code>dlc()</code> .

Value

Index vector for each sample's classification.

References

W. D. Fisher (1958). On grouping for maximum homogeneity. *Journal of the American Statistical Association*, vol. 53, pp. 789–798.

Examples

```
X <- matrix(c(
  9.3, 1.8, 1.9, 1.7, 1.5, 1.3,
  1.4, 2.0, 1.9, 2.3, 2.1
))
C <- dlc(X, maxk = 8)$C
F <- FOP(X, 8, C)
```

Description

Fits the ordered homogeneity pursuit lasso (OHPL) model.

Usage

```
OHPL(
  x,
  y,
  maxcomp,
  gamma,
  cv.folds = 5L,
  G = 30L,
  type = c("max", "median"),
  scale = TRUE,
  pls.method = "simpls"
)
```

Arguments

<code>x</code>	Predictor matrix.
<code>y</code>	Response matrix with one column.
<code>maxcomp</code>	Maximum number of components for PLS.
<code>gamma</code>	A number between (0, 1) for generating the gamma sequence. An usual choice for gamma could be $n * 0.05$, where n is a number in 2, 3, ..., 19.
<code>cv.folds</code>	Number of cross-validation folds.
<code>G</code>	Maximum number of variable groups.
<code>type</code>	Find the maximum absolute correlation ("max") or find the median of absolute correlation ("median"). Default is "max".
<code>scale</code>	Should the predictor matrix be scaled? Default is TRUE.
<code>pls.method</code>	Method for fitting the PLS model. Default is "simpls". See the details section in pls::plsr() for all possible options.

Value

A list of fitted OHPL model object with performance metrics.

References

You-Wu Lin, Nan Xiao, Li-Li Wang, Chuan-Quan Li, and Qing-Song Xu (2017). Ordered homogeneity pursuit lasso for group variable selection with applications to spectroscopic data. *Chemometrics and Intelligent Laboratory Systems* 168, 62–71.

Examples

```

# Generate simulation data
dat <- OHPL.sim(
  n = 100, p = 100, rho = 0.8,
  coef = rep(1, 10), snr = 3, p.train = 0.5,
  seed = 1010
)

# Split training and test set
x <- dat$x.tr
y <- dat$y.tr
x.test <- dat$x.te
y.test <- dat$y.te

# Fit the OHPL model
fit <- OHPL(x, y, maxcomp = 3, gamma = 0.5, G = 10, type = "max")

# Selected variables
fit$Vsel

# Make predictions
y.pred <- predict(fit, x.test)

# Compute evaluation metric RMSEP, Q2 and MAE for the test set
perf <- OHPL.RMSEP(fit, x.test, y.test)
perf$RMSEP
perf$Q2
perf$MAE

```

OHPL.RMSEP

Compute RMSEP, MAE, and Q2 for a test set

Description

Makes predictions on new data and computes the performance evaluation metrics RMSEP, MAE, and Q2.

Usage

```
OHPL.RMSEP(object, newx, newy)
```

Arguments

object	Object of class OHPL fitted by OHPL() .
newx	Predictor matrix of the new data.
newy	Response matrix of the new data (matrix with one column).

Value

A list of the performance metrics.

Examples

```
# Generate simulation data
dat <- OHPL.sim(
  n = 100, p = 100, rho = 0.8,
  coef = rep(1, 10), snr = 3, p.train = 0.5,
  seed = 1010
)

# Split training and test set
x <- dat$x.tr
y <- dat$y.tr
x.test <- dat$x.te
y.test <- dat$y.te

# Fit the OHPL model
fit <- OHPL(x, y, maxcomp = 3, gamma = 0.5, G = 10, type = "max")

# Compute evaluation metric RMSEP, Q2 and MAE for the test set
perf <- OHPL.RMSEP(fit, x.test, y.test)
perf$RMSEP
perf$Q2
perf$MAE
```

OHPL.sim

Generate simulation data for benchmarking sparse regressions (Gaussian response)

Description

Generate simulation data (Gaussian case) following the settings in Xiao and Xu (2015).

Usage

```
OHPL.sim(
  n = 100,
  p = 100,
  rho = 0.8,
  coef = rep(1, 10),
  snr = 3,
  p.train = 0.5,
  seed = 1001
)
```


Arguments

n	Number of observations.
p	Number of variables.
rho	Correlation base for generating correlated variables.
coef	Vector of non-zero coefficients.
snr	Signal-to-noise ratio (SNR).
p.train	Percentage of training set.
seed	Random seed for reproducibility.

Value

A list containing x.tr, x.te, y.tr, and y.te.

Author(s)

Nan Xiao

References

Nan Xiao and Qing-Song Xu. (2015). Multi-step adaptive elastic-net: reducing false positives in high-dimensional variable selection. *Journal of Statistical Computation and Simulation* 85(18), 3755–3765.

Examples

```
dat <- OHPL.sim(  
  n = 100, p = 100, rho = 0.8,  
  coef = rep(1, 10), snr = 3, p.train = 0.5,  
  seed = 1010  
)  
  
dim(dat$x.tr)  
dim(dat$x.te)
```

predict.OHPL

Make predictions based on the fitted OHPL model

Description

Make predictions on new data by an OHPL model object.

Usage

```
## S3 method for class 'OHPL'  
predict(object, newx, ncomp = NULL, type = "response", ...)
```

Arguments

object	Object of class OHPL fitted by <code>OHPL()</code> .
newx	Predictor matrix of the new data.
ncomp	Optimal number of components. If is NULL, the optimal number of components stored in the model object will be used.
type	Prediction type.
...	Additional parameters.

Value

Numeric matrix of the predicted values.

Examples

```
# generate simulation data
dat <- OHPL.sim(
  n = 100, p = 100, rho = 0.8,
  coef = rep(1, 10), snr = 3, p.train = 0.5,
  seed = 1010
)

# split training and test set
x <- dat$x.tr
y <- dat$y.tr
x.test <- dat$x.te
y.test <- dat$y.te

# fit the OHPL model
fit <- OHPL(x, y, maxcomp = 3, gamma = 0.5, G = 10, type = "max")
# make predictions
y.pred <- predict(fit, x.test)
y.pred
```

proto

Extract the prototype from each variable group

Description

Extracts the prototype from each variable group.

Usage

```
proto(X, y, groups, type = c("max", "median"), mu = NULL)
```

Arguments

X	Predictor matrix.
y	Response matrix with one column.
groups	An group index vector containing the group number each variable belongs to. For example: <code>c(1, 1, 1, 1, 1, 2, 2, 2, ...)</code> . Variable groups can be generated by the Fisher optimal partition algorithm implemented in <code>FOP()</code> .
type	The rule for extracting the prototype. Possible options are "max" and "median".
mu	The mean value of y for standardization. Default is NULL, which uses the sample mean of y.

Value

The prototypes (variable index) extracted from each group (cluster).

soil	<i>The soil dataset</i>
------	-------------------------

Description

The soil dataset contains 108 sample measurements from the wavelength range of 400–2500 nm (visible and near infrared spectrum) published by Rinnan et al.

Usage

```
data(soil)
```

References

Rinnan, R., & Rinnan, A. (2007). Application of near infrared reflectance (NIR) and fluorescence spectroscopy to analysis of microbiological and chemical properties of arctic soil. *Soil biology and Biochemistry*, 39(7), 1664–1673.

Examples

```
data("soil")

X <- soil$x
y <- soil$som
n <- nrow(soil$x)

set.seed(1001)
samp.idx <- sample(1L:n, round(n * 0.7))
X.cal <- X[samp.idx, ]
y.cal <- y[samp.idx]
X.test <- X[-samp.idx, ]
y.test <- y[-samp.idx]
```

wheat

The wheat dataset

Description

The wheat dataset contains 100 wheat samples with specified protein and moisture content, published by J. Kalivas. Samples were measured by diffuse reflectance as $\log(I/R)$ from 1100 to 2500 nm (701 data points) in 2 nm intervals.

Usage

```
data(wheat)
```

References

Kalivas, J. H. (1997). Two data sets of near infrared spectra. *Chemometrics and Intelligent Laboratory Systems*, 37(2), 255–259.

Examples

```
data("wheat")

X <- wheat$x
y <- wheat$protein
n <- nrow(wheat$x)

set.seed(1001)
samp.idx <- sample(1L:n, round(n * 0.7))
X.cal <- X[samp.idx, ]
y.cal <- y[samp.idx]
X.test <- X[-samp.idx, ]
y.test <- y[-samp.idx]
```

Index

beer, [2](#)

cv.OHPL, [3](#)

dlc, [4](#)

dlc(), [5](#)

FOP, [5](#)

FOP(), [11](#)

OHPL, [6](#)

OHPL(), [7](#), [10](#)

OHPL.RMSEP, [7](#)

OHPL.sim, [8](#)

pls::plsr(), [3](#), [6](#)

predict.OHPL, [9](#)

proto, [10](#)

soil, [11](#)

wheat, [12](#)