

Package ‘smart’

July 2, 2014

Type Package

Title Sparse Multivariate Analysis via Rank Transformation

Version 1.0.1

Date 2012-08-17

Author Fang Han, Han Liu

Maintainer Fang Han <fhan@jhspk.edu>

Depends R (>= 2.10)

Imports Matrix, gplots, gtools,PMA, elasticnet, pcaPP, igraph

Suggests huge

Description The package “smart” provides a general framework for analyzing (including estimation, feature selection and prediction) and visualize big data. It integrates several novel, efficient and robust data analysis tools, including Transelliptical Component Analysis (TCA), Transelliptical Correlation Estimation (TCE) and Group Nearest Shrunken Centroids (gnsc). We target on high dimensional data analysis(usually $d \gg n$), and exploit computationally efficiently approaches. Results are organized to be visualized properly for users.

License GPL-2

Repository CRAN

Date/Publication 2013-10-22 10:50:06

NeedsCompilation no

R topics documented:

smart-package	2
gnsc.cv	3
gnsc.train	5
plot.gnsc	6
plot.gnscsv	7
plot.TCA	7
plot.TCE	8
print.gnsc	9
print.gnscsv	9
print.TCA	10
print.TCE	11
smart-internal	11
TCA	12
TCE	13
Index	15

smart-package	<i>Sparse Multivariate Analysis via Rank Transformation</i>
---------------	---

Description

A package for Sparse Multivariate Analysis via Rank Transformation

Details

Package:	smart
Type:	Package
Version:	1.0.0
Date:	2012-08-17
License:	GPL-2
LazyLoad:	yes

The package "smart" contain three main functions:

- (i) "gnsc.train" and "gnsc.cv" for conducting Group Nearest Shrunken Centroids.
- (ii) "TCA" for conducting Transeeliptical Component Analysis.
- (iii) "TCE" for conducting Transelliptical Correlation Estimation.

Author(s)

Fang Han, Han Liu
 Maintainer: Fang Han<fhan@jhsph.edu>

References

1. Robert Tibshirani, Trevor Hastie, Balasubramanian Narasimhan, and Gilbert Chu. Diagnosis of multiple cancer types by shrunken centroids of gene expression *PNAS*, 99: 6567-6572.
2. Han Liu, Fang Han, Ming Yuan, John Lafferty, Larry Wasserman. High dimensional semiparametric gaussian copula graphical models. *Annals of Statistics*, to appear.
3. Witten, D., Tibshirani, R., and Hastie, T., A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics*
4. Juemin Yang, Fang Han, Rafa Irizarry, and Han Liu. Gene Context Analysis on Large-scale Genomic Data. *Technical Report*, Johns Hopkins University, 2012
5. Yuan, X. and Zhang, T. (2011). Truncated power method for sparse eigenvalue problems. *Technical Report*, Rutgers, 2011.
6. Tuo Zhao and Han Liu. HUGE: A Package for High-dimensional Undirected Graph Estimation. *Technical Report*, Carnegie Mellon University, 2010
7. Zou, H., Hastie, T. and Tibshirani, R. Sparse principal component analysis. *JCGS*, 2006.

See Also

[TCA](#), [TCE](#), [gns.train](#), and [gns.c.v](#)

gns.c.v	<i>A function to cross-validate the Group Nearest Shrunken Centroid Classifier</i>
---------	--

Description

A function to cross-validate the Group Nearest Shrunken Centroid Classifier produced by `gns.train`

Usage

```
gns.c.v(fit, x, y = NULL, z = NULL, nfold = NULL, folds = NULL, verbose = T)
```

Arguments

<code>fit</code>	The result of a call to <code>gns.train</code>
<code>x</code>	The test data matrix (variables in the rows, samples in the columns).
<code>y</code>	The test class labels for samples, must have the same length as the column length of <code>x</code> .
<code>z</code>	The test class labels for variables, must have the same length as the row length of <code>x</code> .
<code>nfold</code>	Number of cross-validation folds. The default value is the smallest class size.
<code>folds</code>	The fold labels for each sample, must have the same length as <code>y</code> and $\max(\text{folds}) = \text{nfold}$. The default value is <code>sample(1:nfold, n, replace=T)</code> , here <code>n</code> is the sample size.
<code>verbose</code>	If <code>verbose = FALSE</code> , tracing information printing is disabled. The default value is <code>TRUE</code> .

Details

gnsc.cv carries out a cross-validation for Group Nearest Shrunken Centroid Classifier.

Value

An object with S3 class "gnscv" is returned:

lambda	A vector of the thresholds tried in the shrinkage
nlambda	The number of thresholds tried in the shrinkage
lambda.min	The index of the threshold which achieves the lowest cross-validation error
errors	The number of cross-validation errors for each threshold value
nonzero	The number of variables that survived the thresholding
Thresh.mat	A list of estimated $\tilde{\mu}_{\{mk\}}$. See Yang, et.al (2012) for details

Author(s)

Fang Han, Han Liu
Maintainer: Fang Han<fhan@jhsph.edu>

References

1. Juemin Yang, Fang Han, Rafa Irizarry, and Han Liu. Gene Context Analysis on Large-scale Genomic Data. *Technical Report*, Johns Hopkins University, 2012
2. Robert Tibshirani, Trevor Hastie, Balasubramanian Narasimhan, and Gilbert Chu. Diagnosis of multiple cancer types by shrunken centroids of gene expression *PNAS*, 99: 6567-6572.

See Also

[gnsc.train](#)

Examples

```
set.seed(120)
x <- matrix(rnorm(1000*20), ncol=20)
y <- sample(c(1:4), size=20, replace=TRUE)
z <- sample(c(1:10), size=1000, replace=TRUE)
fit=gnsc.train(x, col.struc=y, row.struc=z, lambda.max=5, nlambda=20)
fit
plot(fit)
fit.cv=gnsc.cv(fit,x,y,z)
fit.cv
plot(fit.cv)
```

gns.train	<i>gns.train</i>
-----------	------------------

Description

A function to conduct the Group Nearest Shrunken Centroid Classifier

Usage

```
gns.train(x, col.struc = NULL, row.struc = NULL, standardize = T,
          nlambda = NULL, lambda.max = 10, lambda = NULL, verbose = TRUE)
```

Arguments

<code>x</code>	The train data matrix (variables in the rows, samples in the columns).
<code>col.struc</code>	The train class labels for samples, must have the same length as the column length of <code>x</code> .
<code>row.struc</code>	The train class labels for variables, must have the same length as the row length of <code>x</code> .
<code>standardize</code>	Logical value to determine whether to standardize the data. The default value is TRUE.
<code>nlambda</code>	The number of thresholding parameters. The default value is 10.
<code>lambda.max</code>	The largest lambda value, given the thresholding parameters lambda is not provided by the user.
<code>lambda</code>	A sequence of positive numbers to control to determine the thresholding level.
<code>verbose</code>	If <code>verbose = FALSE</code> , tracing information printing is disabled. The default value is TRUE.

Details

`gns.train` conducts a Group Nearest Shrunken Centroid Classifier.

Value

An object with S3 class "gns" is returned:

<code>lambda</code>	A vector of the thresholds tried in the shrinkage
<code>nlambda</code>	The number of thresholds tried in the shrinkage
<code>yhat</code>	A matrix with the estimated sample labels for each thresholding level in each column
<code>errors</code>	The number of estimated errors for each threshold value
<code>nonzero</code>	The number of variables that survived the thresholding for each thresholding value
<code>...</code>	System reserved (No specific usage)

Author(s)

Fang Han, Han Liu
 Maintainer: Fang Han<fhan@jhsp.edu>

References

1. Juemin Yang, Fang Han, Rafa Irizarry, and Han Liu. Gene Context Analysis on Large-scale Genomic Data. *Technical Report*, Johns Hopkins University, 2012
2. Robert Tibshirani, Trevor Hastie, Balasubramanian Narasimhan, and Gilbert Chu. Diagnosis of multiple cancer types by shrunken centroids of gene expression *PNAS*, 99: 6567-6572.

See Also

[gnsc.cv](#)

Examples

```
set.seed(120)
x <- matrix(rnorm(1000*20), ncol=20)
y <- sample(c(1:4), size=20, replace=TRUE)
z <- sample(c(1:10), size=1000, replace=TRUE)
fit=gnsc.train(x, col.struc=y, row.struc=z, lambda.max=5, nlambda=20)
fit
plot(fit)
```

plot.gnsc

Plot function for S3 class "gnsc"

Description

Plot sparsity level information and heatmap from the gnsc.train

Usage

```
## S3 method for class 'gnsc'
plot(x, which.lambda = NULL, ...)
```

Arguments

x	An object with S3 class "gnsc"
which.lambda	pick one lambda to visualize the intensity heatmap. The default value is median(lambda).
...	System reserved (No specific usage)

Author(s)

Fang Han, Han Liu
 Maintainer: Fang Han<fhan@jhsp.edu>

See Also

[gnsccv.train](#)

plot.gnsccv *Plot function for S3 class "gnsccv"*

Description

Plot sparsity level information and heatmap from the gnsccv

Usage

```
## S3 method for class 'gnsccv'  
plot(x, ...)
```

Arguments

x An object with S3 class "gnsccv"
... System reserved (No specific usage)

Author(s)

Fang Han, Han Liu
Maintainer: Fang Han<fhan@jhsph.edu>

See Also

[gnsccv.cv](#)

plot.TCA *Plot function for S3 class "TCA"*

Description

Plot sparsity level information and projected PCs from the TCA

Usage

```
## S3 method for class 'TCA'  
plot(x, pc = NULL, ...)
```

Arguments

x	An object with S3 class "TCA"
pc	pick two PCs to visualize the projected principal components. The default value is c(1,2).
...	System reserved (No specific usage)

Author(s)

Fang Han, Han Liu
Maintainer: Fang Han<fhan@jhsph.edu>

See Also

[TCA](#)

plot.TCE

Plot function for S3 class "TCE"

Description

Plot sparsity level information and 3 typical sparse graphs from the correlation graph path

Usage

```
## S3 method for class 'TCE'  
plot(x, align = FALSE, ...)
```

Arguments

x	An object with S3 class "TCE"
align	If align = FALSE, 3 plotted graphs are aligned
...	System reserved (No specific usage)

Author(s)

Fang Han, Han Liu
Maintainer: Fang Han<fhan@jhsph.edu>

See Also

[TCE](#)

print.gnsc	<i>Print function for S3 class "gnsc"</i>
------------	---

Description

Print the information about threshold values, selected variables, fitted errors, predicted values

Usage

```
## S3 method for class 'gnsc'  
print(x, ...)
```

Arguments

x	An object with S3 class "gnsc"
...	System reserved (No specific usage)

Author(s)

Fang Han, Han Liu
Maintainer: Fang Han<fhan@jhsph.edu>

See Also

[gnsc.train](#)

print.gnscv	<i>Print function for S3 class "gnscv"</i>
-------------	--

Description

Print the information about threshold values, selected variables, cross-validation errors

Usage

```
## S3 method for class 'gnscv'  
print(x, ...)
```

Arguments

x	An object with S3 class "gnscv"
...	System reserved (No specific usage)

Author(s)

Fang Han, Han Liu
Maintainer: Fang Han<fhan@jhsph.edu>

See Also

[gnsc.cv](#)

print.TCA	<i>Print function for S3 class "TCA"</i>
-----------	--

Description

Print the information about model, algorithm and results

Usage

```
## S3 method for class 'TCA'  
print(x, ...)
```

Arguments

x	An object with S3 class "TCA"
...	System reserved (No specific usage)

Author(s)

Fang Han, Han Liu
Maintainer: Fang Han<fhan@jhsph.edu>

See Also

[TCA](#)

print.TCE *Print function for S3 class "TCE"*

Description

Print the information about the model usage, the graph path length, graph dimension, sparsity level

Usage

```
## S3 method for class 'TCE'  
print(x, ...)
```

Arguments

x	An object with S3 class "TCE"
...	System reserved (No specific usage)

Author(s)

Fang Han, Han Liu
Maintainer: Fang Han<fhan@jhsph.edu>

See Also

[TCE](#)

smart-internal *Internal smart functions*

Description

Internal smart functions

Details

These are not intended for use by users. Please refer to `gnsc.train()`, `gnsc.cv()`, `TCA()`, `TCE()`.

Author(s)

Fang Han, Han Liu
Maintainer: Fang Han<fhan@jhsph.edu>

TCA

*Transelliptical Component Analysis***Description**

A function to conduct Transelliptical Component Analysis

Usage

```
TCA(x, K, para, method = "kendall", algorithm = "tp", max.iter = 200,
    verbose = TRUE, eps.conv = 0.001)
```

Arguments

x	The n by d data matrix or d by d covariance matrix from the input
K	Number of components
para	A vector of length K, indicating the number of sparse loadings.
method	Method to be used to estimating the correlation matrix with 5 options: pearson, ns, npn, spearman and kendall. kendall as default.
algorithm	Algorithm to be used to obtain sparse loadings with 3 options: sp, spca and pmd. tp as default.
max.iter	Maximum number of iterations.
verbose	If verbose = FALSE, tracing information printing is disabled. The default value is TRUE.
eps.conv	Convergence criterion.

Details

PCA and Sparse PCA is sensitive to modeling assumption, outliers, missing values data dependency. We propose an alternative way using rank-based methods including ns, npn, spearman and kendall to approximate the correlation matrix. Details are referred to Han,F. and Liu,H. (2012). Three sparse PCA algorithms are used: truncated power (Yuan, X. and Zhang, T. (2011)), spca(Zou,H., Hastie, T., and Tibshirani, R. (2006)) and pmd (Witten, D., Tibshirani, R., and Hastie, T. (2009)).

Value

cov.input	An indicator of the sample covariance.
loadings	The loadings of the sparse PCs.
pev	An indicator of the sample covariance.
PC	Projected PCs. existing if cov.input=TRUE.
method	The method used in estimating the correlation matrix.
algorithm	The algorithm used in obtaining the sparse loadings.
K	The number of components.

Author(s)

Fang Han, Han Liu
 Maintainer: Fang Han<fhan@jhsph.edu>

References

1. Witten, D., Tibshirani, R., and Hastie, T., A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics*
2. Yuan, X. and Zhang, T. (2011). Truncated power method for sparse eigenvalue problems. *Technical Report*, Rutgers, 2011.
3. Zou, H., Hastie, T. and Tibshirani, R. Sparse principal component analysis. *JCGS*, 2006.

Examples

```
x=matrix(rnorm(20000),100)
fit=TCA(x,K=6, para=c(10,10,10,5,5,5))
fit
plot(fit)
```

TCE

Transelliptical Correlation Estimation

Description

A function to conduct Transelliptical Correlation Estimation

Usage

```
TCE(x, method, nlambda = NULL, lambda.min.ratio = NULL, lambda = NULL, verbose = TRUE)
```

Arguments

x	The n by d data matrix or d by d covariance matrix from the input
method	Method to be used to estimating the correlation matrix with 5 options: pearson, ns, npr, spearman and kendall. kendall as default.
nlambda	The number of regularization/thresholding paramters. The default value is 20.
lambda.min.ratio	The largest sparsity level for the estimated graphs. The default valye is 0.05.
lambda	A sequence of positive numbers for conducting thresholding.
verbose	If verbose = FALSE, tracing information printing is disabled. The default value is TRUE.

Details

The correlation graph is estimated by correlation cut-off based on the given thresholding level.

Value

An object with S3 class "TCE" is returned:

<code>cov.input</code>	An indicator of the sample covariance.
<code>path</code>	A list of k by k adjacency matrices of estimated graphs as a graph path corresponding to <code>lambda</code> .
<code>sparsity</code>	The sparsity levels of the graph path.
<code>method</code>	The method used in the correlation graph estimation stage.
<code>lambda</code>	The sequence of thresholding parameters used.

Author(s)

Fang Han, Han Liu
Maintainer: Fang Han<fhan@jhsph.edu>

References

- 1.Han Liu, Fang Han, Ming Yuan, John Lafferty, Larry Wasserman. High dimensional semiparametric gaussian copula graphical models. *Annals of Statistics*, to appear.
- 2.Tuo Zhao and Han Liu. HUGE: A Package for High-dimensional Undirected Graph Estimation. *Technical Report*, Carnegie Mellon University, 2010

Examples

```
require(huge)
L = huge.generator(n = 200, d = 80, graph = "hub")
out = TCE(L$data,method="kendall")
out
plot(out)
```

Index

`gnsc.cv`, [3](#), [3](#), [6](#), [7](#), [10](#)
`gnsc.heatmap` (`smart-internal`), [11](#)
`gnsc.icov` (`smart-internal`), [11](#)
`gnsc.predict` (`smart-internal`), [11](#)
`gnsc.restruc` (`smart-internal`), [11](#)
`gnsc.train`, [3](#), [4](#), [5](#), [7](#), [9](#)

`minid` (`smart-internal`), [11](#)

`plot.gnsc`, [6](#)
`plot.gnsccv`, [7](#)
`plot.TCA`, [7](#)
`plot.TCE`, [8](#)
`print.gnsc`, [9](#)
`print.gnsccv`, [9](#)
`print.TCA`, [10](#)
`print.TCE`, [11](#)

`smart` (`smart-package`), [2](#)
`smart-internal`, [11](#)
`smart-package`, [2](#)
`smart.npn` (`smart-internal`), [11](#)

TCA, [3](#), [8](#), [10](#), [12](#)
TCE, [3](#), [8](#), [11](#), [13](#)