

Package ‘SPCAvRP’

January 16, 2018

Type Package

Title Sparse Principal Component Analysis via Random Projections
(SPCAvRP)

Version 0.2

Date 2018-01-16

Author Milana Gataric, Tengyao Wang and Richard J. Samworth

Maintainer Milana Gataric <m.gataric@statslab.cam.ac.uk>

Description Implements the SPCA_vRP algorithm, developed and analysed in “Sparse principal component analysis via random projections” Gataric, M., Wang, T. and Samworth, R. J. (2017) <arXiv:1712.05630>. The algorithm is based on the aggregation of eigenvector information from carefully-selected random projections of the sample covariance matrix.

Depends R (>= 3.0.0), parallel, MASS

License GPL-3

URL <https://arxiv.org/abs/1712.05630>

NeedsCompilation no

Repository CRAN

RoxygenNote 6.0.1

Date/Publication 2018-01-16 12:51:09 UTC

R topics documented:

final_estimator	2
project_covariance	3
select_projection	4
select_projections_subspace	5
SPCAvRP	6
SPCAvRP_deflation	8
SPCAvRP_parallel	10
SPCAvRP_ranking	12
SPCAvRP_subspace	13

Index	16
--------------	-----------

final_estimator	<i>Computes the leading eigenvector from its support</i>
-----------------	--

Description

Computes the leading eigenvector of the sample covariance matrix given the indices of variables ranked by their importance and desired sparsity level.

Usage

```
final_estimator(data, cov, l, ranking)
```

Arguments

data	Either the data matrix or the sample covariance matrix.
cov	TRUE if data is given as a sample covariance matrix.
l	Desired sparsity of the final estimator (see Details).
ranking	Original variables ranked by their importance.

Details

If true sparsity level k is known use $l = k$. If k is unknown, l can be an array of different values and then the eigenvectors of the corresponding sparsity levels are returned.

Value

Returns a list of two elements:

vector	A vector or a matrix with $\text{length}(l)$ columns as the estimated eigenvectors of sparsity level l .
value	An array with $\text{length}(l)$ estimated eigenvalues.

Author(s)

Milana Gataric, Tengyao Wang and Richard J. Samworth

References

Milana Gataric, Tengyao Wang and Richard J. Samworth (2017) Sparse principal component analysis via random projections <https://arxiv.org/abs/1712.05630>

Examples

```
p <- 80
k <- 8
n <- 1000
v1 <- c(rep(1/sqrt(k), k), rep(0,p-k))
Sigma <- 2*tcrossprod(v1) + diag(p)
mu <- rep(0, p)
X <- mvrnorm(n, mu, Sigma)
Sigma_hat <- 1/n*crossprod(X)

A <- 200
B <- 100
d <- 10

rand_ind <- matrix(replicate(A*B,sample.int(p,d)), nrow = A*B, byrow = TRUE)
cov_projections <- project_covariance(data = Sigma_hat, cov = TRUE, rand_ind)
ranking <- SPCAVRP_ranking(cov_projections, rand_ind, p, A)

output <- final_estimator(data = Sigma_hat, cov = TRUE, l = (5:11), ranking)
df <- data.frame(5:11,output$value); colnames(df) <- c('l','eigenvalue')
print(df)
```

project_covariance *Projects the sample covariance*

Description

Projects the sample covariance matrix along given axis-aligned projections.

Usage

```
project_covariance(data, cov, rand_ind)
```

Arguments

data	Either the data matrix or the sample covariance matrix (see Details).
cov	TRUE if data is given as a sample covariance matrix.
rand_ind	Matrix whose rows are the indices of non-zero entries of axis-aligned projections.

Details

If the dimension of data is very large, it might be faster if the data matrix is provided as the input.

Value

Returns a list of the sample covariance projections:

```

projections[[1]]
    projected sample covariance along indices of rand_ind[1,]

...

projections[[nrow(rand_ind)]]
    projected sample covariance along indices of rand_ind[nrow(rand_ind),]

```

Author(s)

Milana Gataric, Tengyao Wang and Richard J. Samworth

References

Milana Gataric, Tengyao Wang and Richard J. Samworth (2017) Sparse principal component analysis via random projections <https://arxiv.org/abs/1712.05630>

Examples

```

p <- 50 # dimension of data
k <- 5 # sparsity level
n <- 1000 # number of observations
v1 <- c(rep(1/sqrt(k), k), rep(0,p-k))
Sigma_hat <- 1/n*crossprod(mvrnorm(n, rep(0,p), tcrossprod(v1)+diag(p)))

N <- 1000 # number of projections
d <- k # dimension of projections
rand_ind <- matrix(replicate(N,sample.int(p,d)), nrow = N, byrow = TRUE) # axis-aligned projections

cov_projections <- project_covariance(data = Sigma_hat, cov = TRUE, rand_ind)

```

select_projection	<i>Selects the best projection</i>
-------------------	------------------------------------

Description

Selects the projection yielding the largest eigenvalue among one group of B different d-dimensional axis-aligned projections generated uniformly at random.

Usage

```
select_projection(data, cov = TRUE, p, d, B)
```

Arguments

data	Either the data matrix or the sample covariance matrix.
cov	FALSE if data is given as a data matrix.
p	Original dimension of the data.
d	Dimension of the random projections.
B	Number of random projections to generate and select from.

Details

If p is very large, it might be faster if the data matrix is provided as the input.

Value

Returns eigenvector v_hat_star corresponding to the projected covariance yielding the largest eigenvalue among B different d -dimensional axis-aligned random projections.

Author(s)

Milana Gataric, Tengyao Wang and Richard J. Samworth

References

Milana Gataric, Tengyao Wang and Richard J. Samworth (2017) Sparse principal component analysis via random projections <https://arxiv.org/abs/1712.05630>

Examples

```
p <- 100
k <- 10
n <- 1000
v1 <- c(rep(1/sqrt(k), k), rep(0,p-k))
Sigma <- 2*tcrossprod(v1) + diag(p)
mu <- rep(0, p)

X <- mvrnorm(n, mu, Sigma)

v_hat_star <- select_projection(data = 1/n*crossprod(X), cov = TRUE, p, d = k, B = 100)
```

select_projections_subspace

Selects the best projections for the subspace estimation

Description

Selects A projections yielding the largest r -th eigenvalue among B different random projections, for $r=1:s$.

Usage

```
select_projections_subspace(data, rand_ind, s, p, A)
```

Arguments

data	A list of projected covariances to select from (in the form of the output of project_covariance).
rand_ind	Corresponding projections used to generate data.
s	The number of eigenvalues to estimate.
p	The original dimension of samples.
A	The number of projections to select.

Value

Returns matrix `v_hat_stars` with `A` columns that correspond to `s` eigenvectors yielding the largest eigenvalues.

Author(s)

Milana Gataric, Tengyao Wang and Richard J. Samworth

References

Milana Gataric, Tengyao Wang and Richard J. Samworth (2017) Sparse principal component analysis via random projections <https://arxiv.org/abs/1712.05630>

SPCAvRP

Computes the leading eigenvector using the SPCA_vRP algorithm

Description

Computes l -sparse leading eigenvector of the sample covariance matrix, using A groups of B random axis-aligned projections of dimension d .

Usage

```
SPCAvRP(data, cov = FALSE, l, d = 10, A = 300, B = 100,  
datascaling = TRUE, center = TRUE, scale = TRUE)
```

Arguments

<code>data</code>	Either the data matrix or the sample covariance matrix.
<code>cov</code>	TRUE if data is given as a sample covariance matrix.
<code>l</code>	Desired sparsity level in the final estimator (see Details).
<code>d</code>	The dimension of the random projections.
<code>A</code>	Number of projections over which to aggregate.
<code>B</code>	Number of projections in a group from which to select.
<code>datascaling</code>	TRUE if function <code>scale()</code> should be applied to the data matrix with its arguments <code>center</code> and <code>scale</code> .
<code>center</code>	If <code>datascaling == TRUE</code> , this is the argument of function <code>scale()</code> .
<code>scale</code>	If <code>datascaling == TRUE</code> , this is the argument of function <code>scale()</code> .

Details

This function implements the SPCAvRP algorithm.

If the true sparsity level k is known, use $d = k$ and $l = k$. If k is unknown, the default choice for d is 10, while l can take an array of different values and then the estimators of the corresponding sparsity levels are computed.

It is desirable to choose A as big as possible subject to the computational budget. In general, we suggest using $A = 300$ and $B = 100$ when the dimension of data is a few hundreds, while $A = 600$ and $B = 200$ when the dimension is on order of 1000.

If `datascaling == TRUE` and data is given as a data matrix, the first step is to scale it by executing `scale(data, center, scale)`. By default `center == TRUE`, which means that centering is done by subtracting the column means of data from their corresponding columns; also `scale == TRUE`, which means that the scaling is done by dividing the (centered) columns of data by their standard deviations in case `center == TRUE`, and the root mean square otherwise.

Value

Returns a list of two elements:

<code>vector</code>	A vector or a matrix with <code>length(l)</code> columns as the estimated eigenvectors of sparsity level l .
<code>value</code>	An array with <code>length(l)</code> estimated eigenvalues.

Author(s)

Milana Gataric, Tengyao Wang and Richard J. Samworth

References

Milana Gataric, Tengyao Wang and Richard J. Samworth (2017) Sparse principal component analysis via random projections <https://arxiv.org/abs/1712.05630>

See Also

[SPCAvRP_parallel](#), [SPCAvRP_ranking](#)

Examples

```
p <- 100 # dimension of data
k <- 10 # true sparsity level
n <- 1000 # number of observations
v1 <- c(rep(1/sqrt(k), k), rep(0,p-k)) # leading eigenvector
Sigma <- 2*tcrossprod(v1) + diag(p) # population covariance
mu <- rep(0, p) # population mean
X <- mvrnorm(n, mu, Sigma) # data matrix

spca <- SPCAvRP(data = X, cov = FALSE, l = k, d = k, A = 200, B = 70, datascaling = FALSE)
spca$vector
spca$value
```

SPCAvRP_deflation	<i>Computes the leading eigenvectors using the modified deflation scheme</i>
-------------------	--

Description

Computes s leading eigenvectors of the sample covariance matrix which are sparse and orthogonal, using the modified deflation scheme in conjunction with the SPCAvRP algorithm.

Usage

```
SPCAvRP_deflation(data, cov = FALSE, s, l, d = 10,
  A = 300, B = 100, datascaling = TRUE, center = TRUE, scale = TRUE)
```

Arguments

<code>data</code>	Either the data matrix or the sample covariance matrix.
<code>cov</code>	TRUE if data is given as a sample covariance matrix.
<code>s</code>	The number of eigenvectors to compute.
<code>l</code>	The array of length s with the desired sparsity levels in the final estimators.
<code>d</code>	The dimension of the random projections.
<code>A</code>	Number of projections over which to aggregate.
<code>B</code>	Number of projections in a group from which to select.
<code>datascaling</code>	TRUE if function <code>scale()</code> should be applied to the data matrix with its arguments <code>center</code> and <code>scale</code> .
<code>center</code>	If <code>datascaling == TRUE</code> , this is the argument of function <code>scale()</code> .
<code>scale</code>	If <code>datascaling == TRUE</code> , this is the argument of function <code>scale()</code> .

Details

This function implements the modified deflation scheme in conjunction with the SPCAvRP in order to compute s sparse eigenvectors that are orthogonal. If possible, use `SPCAvRP_subspace` instead.

If the true sparsity level is known and for each component is equal to k , use $d = k$ and $l = \text{rep}(k, s)$. Sparsity levels of different components may take different values. If k is unknown, appropriate k could be chosen from an array of different values by inspecting the explained variance for one component at the time and by using SPCAvRP in a combination with the deflation scheme implemented in `SPCAvRP_deflation`.

It is desirable to choose A as big as possible subject to the computational budget. In general, we suggest using $A = 300$ and $B = 100$ when the dimension of data is a few hundreds, while $A = 600$ and $B = 200$ when the dimension is on order of 1000 .

If `datascaling == TRUE` and `data` is given as a data matrix, the first step is to scale it by executing `scale(data, center, scale)`. By default `center == TRUE`, which means that centering is done by subtracting the column means of `data` from their corresponding columns; also `scale == TRUE`, which means that the scaling is done by dividing the (centered) columns of `data` by their standard deviations in case `center == TRUE`, and the root mean square otherwise.

Value

Returns a list of two elements:

vector	A matrix whose s columns are the estimated eigenvectors.
value	An array with s estimated eigenvalues.

Author(s)

Milana Gataric, Tengyao Wang and Richard J. Samworth

References

Milana Gataric, Tengyao Wang and Richard J. Samworth (2017) Sparse principal component analysis via random projections <https://arxiv.org/abs/1712.05630>

See Also

`SPCAvRP`, `SPCAvRP_subspace`

Examples

```
p <- 50
k <- 8
theta <- 40
v1 <- c(rep(1/sqrt(k), k), rep(0, p-k))
theta2 <- 20
v2 <- c(rep(0,4), 1/sqrt(k), -1/sqrt(k), 1/sqrt(k), -1/sqrt(k), rep(1/sqrt(k),4), rep(0,p-12))
theta3 <- 5
v3 <- c(rep(0,6), 1/sqrt(k), -rep(1/sqrt(k),4), rep(1/sqrt(k),3), rep(0,p-14))
Sigma <- diag(p) + theta*tcrossprod(v1) + (theta2)*tcrossprod(v2) + (theta3)*tcrossprod(v3)
mu <- rep(0, p)
```

```
n <- 2000
X <- mvrnorm(n, mu, Sigma)

spcarp <- SPCAvRP_deflation(data = X, cov = FALSE, s = 1, l = k, d = k,
                             A = 300, B = 100, datascaling = FALSE)
```

SPCAvRP_parallel

Parallel implementation of the SPCAvRP algorithm

Description

Computes l -sparse leading eigenvector of the sample covariance matrix, by parallel selection of A projections, where each projection is selected from a group of B random projections of dimension d .

Usage

```
SPCAvRP_parallel(data, cov = FALSE, l, d = 10, A = 300, B = 100,
                 datascaling = TRUE, center = TRUE, scale = TRUE,
                 cluster_type = "PSOCK", cores = 1, machine_names = NULL)
```

Arguments

<code>data</code>	Either the data matrix or the sample covariance matrix.
<code>cov</code>	TRUE if data is given as a sample covariance matrix.
<code>l</code>	Desired sparsity level in the final estimator (see Details).
<code>d</code>	The dimension of the random projections.
<code>A</code>	Number of projections over which to aggregate.
<code>B</code>	Number of projections in a group from which to select.
<code>datascaling</code>	TRUE if function <code>scale()</code> should be applied to the data matrix with its arguments <code>center</code> and <code>scale</code> .
<code>center</code>	If <code>datascaling == TRUE</code> , this is the argument of function <code>scale()</code> .
<code>scale</code>	If <code>datascaling == TRUE</code> , this is the argument of function <code>scale()</code> .
<code>cluster_type</code>	Can be "PSOCK" or "FORK" (cf. package "parallel").
<code>cores</code>	Number of cores to use if <code>clustertype=="FORK"</code> .
<code>machine_names</code>	Names of computers on the network if <code>clustertype=="PSOCK"</code> .

Details

This function implements the parallelised SPCAvRP algorithm, by calling 'select_projection' A times in parallel. We recommend to use this function if p , A and B are large; otherwise use [SPCAvRP](#).

If the true sparsity level k is known, use $d = k$ and $l = k$. If k is unknown, the default choice for d is 10, while l can take an array of different values and then the estimators of the corresponding sparsity levels are computed.

It is desirable to choose A as big as possible subject to the computational budget. In general, we suggest using A = 300 and B = 100 when the dimension of data is a few hundreds, while A = 600 and B = 200 when the dimension is on order of 1000.

If `datascaling == TRUE` and data is given as a data matrix, the first step is to scale it by executing `scale(data, center, scale)`. By default `center == TRUE`, which means that centering is done by subtracting the column means of data from their corresponding columns; also `scale == TRUE`, which means that the scaling is done by dividing the (centered) columns of data by their standard deviations in case `center == TRUE`, and the root mean square otherwise.

Value

Returns a list of two elements:

vector	A vector or a matrix with <code>length(1)</code> columns as the estimated eigenvectors of sparsity level <code>l</code> .
value	An array with <code>length(1)</code> estimated eigenvalues.

Author(s)

Milana Gataric, Tengyao Wang and Richard J. Samworth

References

Milana Gataric, Tengyao Wang and Richard J. Samworth (2017) Sparse principal component analysis via random projections <https://arxiv.org/abs/1712.05630>

See Also

[SPCAvRP](#)

Examples

```
p <- 100 # dimension of data
k <- 10 # true sparsity level
n <- 1000 # number of observations
v1 <- c(rep(1/sqrt(k), k), rep(0,p-k)) # leading eigenvector
Sigma <- 2*tcrossprod(v1) + diag(p) # population covariance
mu <- rep(0, p) # population mean
X <- mvrnorm(n, mu, Sigma) # data matrix

spca <- SPCAvRP_parallel(data = X, cov = FALSE, l = k, d = k, A = 200, B = 70,
                        datascaling = FALSE, cluster_type = "PSOCK")
```

SPCAvRP_ranking	<i>Ranks the variables</i>
-----------------	----------------------------

Description

Ranks the original variables according to their importance in maximising the explained variance in the data by aggregating over selected projections of the sample covariance matrix.

Usage

```
SPCAvRP_ranking(data, rand_ind, p, A)
```

Arguments

data	A list of projected covariances (see Details).
rand_ind	Corresponding axis-aligned projections (see Details).
p	The dimension of the data.
A	Number of projections over which to aggregate.

Details

This function divides given projections into A groups and selects the best one from each group. This is then followed by an aggregation step. Data is given as a list of d-dimensional projections of p-dimensional sample covariance matrix, as generated by the function [project_covariance](#). Corresponding axis-aligned projections are given as a matrix whose rows are the indices of the projection's non-zero entries.

Value

Returns the vector of indices corresponding to variables ranked by their importance in maximising the explained variance.

Author(s)

Milana Gataric, Tengyao Wang and Richard J. Samworth

References

Milana Gataric, Tengyao Wang and Richard J. Samworth (2017) Sparse principal component analysis via random projections <https://arxiv.org/abs/1712.05630>

See Also

[SPCAvRP](#)

Examples

```

p <- 50 # dimension of data
k <- 5 # true sparsity level
n <- 1000 # number of observations
v1 <- c(rep(1/sqrt(k), k), rep(0,p-k)) # leading eigenvector
Sigma <- tcrossprod(v1) + diag(p) # population covariance
mu <- rep(0, p) # population mean
X <- mvrnorm(n, mu, Sigma) # data matrix
Sigma_hat <- 1/n*crossprod(X) # sample covariance matrix

A <- 200 # number of projections over which to aggregate
B <- 100 # number of projections in a group from which to select
d <- k # dimension of projections
rand_ind <- matrix(replicate(A*B,sample.int(p,d)), nrow = A*B, byrow = TRUE) # random projections

cov_projections <- project_covariance(data = Sigma_hat, cov = TRUE, rand_ind)

ranking <- SPCAvRP_ranking(cov_projections, rand_ind, p, A)
print(ranking)

```

SPCAvRP_subspace	<i>Computes the leading eigenspace using the SPCAvRP algorithm for eigenspace estimation</i>
------------------	--

Description

Computes s leading eigenvectors of the sample covariance matrix which are sparse and orthogonal, using A groups of B random axis-aligned projections of dimension d .

Usage

```

SPCAvRP_subspace(data, cov = FALSE, s, l, d = 10,
A = 300, B = 100, datascaling = TRUE, center = TRUE, scale = TRUE)

```

Arguments

<code>data</code>	Either the data matrix or the sample covariance matrix.
<code>cov</code>	TRUE if data is given as a sample covariance matrix.
<code>s</code>	The dimension of the eigenspace, i.e the number of principal components.
<code>l</code>	The array of length s with the desired sparsity levels in the final estimators.
<code>d</code>	The dimension of the random projections.
<code>A</code>	Number of projections over which to aggregate.
<code>B</code>	Number of projections in a group from which to select.
<code>datascaling</code>	TRUE if function <code>scale()</code> should be applied to the data matrix with its arguments <code>center</code> and <code>scale</code> .
<code>center</code>	If <code>datascaling == TRUE</code> , this is the argument of function <code>scale()</code> .
<code>scale</code>	If <code>datascaling == TRUE</code> , this is the argument of function <code>scale()</code> .

Details

This function implements the SPCA_vRP algorithm for eigenspace estimation.

If the true sparsity level is known and for each component is equal to k , use $d = k$ and $l = \text{rep}(k, s)$. Sparsity levels of different components may take different values. If k is unknown, appropriate k could be chosen from an array of different values by inspecting the explained variance for one component at the time and by using SPCA_vRP in a combination with SPCA_vRP_deflation.

It is desirable to choose A as big as possible subject to the computational budget. In general, we suggest using $A = 300$ and $B = 100$ when the dimension of data is a few hundreds, while $A = 600$ and $B = 200$ when the dimension is on order of 1000 .

If `datascaling == TRUE` and `data` is given as a data matrix, the first step is to scale it by executing `scale(data, center, scale)`. By default `center == TRUE`, which means that centering is done by subtracting the column means of `data` from their corresponding columns; also `scale == TRUE`, which means that the scaling is done by dividing the (centered) columns of `data` by their standard deviations in case `center == TRUE`, and the root mean square otherwise.

Value

Returns a list of two elements:

vector	A matrix whose s columns are the estimated eigenvectors.
value	An array with s estimated eigenvalues.

Author(s)

Milana Gataric, Tengyao Wang and Richard J. Samworth

References

Milana Gataric, Tengyao Wang and Richard J. Samworth (2017) Sparse principal component analysis via random projections <https://arxiv.org/abs/1712.05630>

See Also

SPCA_vRP, SPCA_vRP_deflation

Examples

```
p <- 50
k <- 8
theta <- 40
v1 <- c(rep(1/sqrt(k), k), rep(0, p-k))
theta2 <- 20
v2 <- c(rep(0,4), 1/sqrt(k), -1/sqrt(k), 1/sqrt(k), -1/sqrt(k), rep(1/sqrt(k),4), rep(0,p-12))
theta3 <- 5
v3 <- c(rep(0,6), 1/sqrt(k), -rep(1/sqrt(k),4), rep(1/sqrt(k),3), rep(0,p-14))
Sigma <- diag(p) + theta*tcrossprod(v1) + (theta2)*tcrossprod(v2) + (theta3)*tcrossprod(v3)
mu <- rep(0, p)
n <- 2000
X <- mvrnorm(n, mu, Sigma)
```

```
loss = function(u,v){
  sqrt(abs(1-sum(v*u)^2))
}
loss_sub = function(U,V){
  V<-qr.Q(qr(V))
  norm(tcrossprod(U)-tcrossprod(V),"2")
}

s <- 2
spcarp <- SPCAvRP_subspace(data = X, cov = FALSE, s, l = rep(k,s), d = k,
  A = 200, B = 100, datascaling = FALSE)

subspace_estimation<-data.frame(
  loss(spcarp$vector[,1],v1),
  loss(spcarp$vector[,2],v2),
  loss_sub(matrix(c(v1,v2),ncol=s),spcarp$vector),
  crossprod(spcarp$vector[,1],spcarp$vector[,2]))
colnames(subspace_estimation)<-c("loss_1","loss_2","loss_sub","inp")
rownames(subspace_estimation)<-c("")
subspace_estimation
```

Index

final_estimator, 2

project_covariance, 3, 6, 12

select_projection, 4

select_projections_subspace, 5

SPCAvRP, 6, 9–12, 14

SPCAvRP_deflation, 8, 14

SPCAvRP_parallel, 8, 10

SPCAvRP_ranking, 8, 12

SPCAvRP_subspace, 9, 13