

Package ‘orthoDr’

March 26, 2018

Type Package

Title An Orthogonality Constrained Optimization Approach for
Semi-Parametric Dimension Reduction Problems

Version 0.5.1

Author Ruilin Zhao, Jiyang Zhang and Ruoqing Zhu

Maintainer Ruoqing Zhu <teazrq@gmail.com>

Description Utilize an orthogonality constrained optimization algorithm of Wen & Yin (2013) <DOI:10.1007/s10107-012-0584-1> to solve a variety of dimension reduction problems in the semiparametric framework, such as Ma & Zhu (2012) <DOI:10.1080/01621459.2011.646925>, Ma & Zhu (2013) <DOI:10.1214/12-AOS1072>, and Sun, Zhu, Wang & Zeng (2017) <arXiv:1704.05046>. It also serves as a general purpose optimization solver for problems with orthogonality constraints. Parallel computing is enabled through ‘OpenMP’.

License GPL (>= 2)

Encoding UTF-8

LazyData TRUE

RoxygenNote 6.0.1

NeedsCompilation yes

Repository CRAN

Imports Rcpp (>= 0.12.12), survival, dr, pracma, plot3D, rgl, MASS

LinkingTo Rcpp, RcppArmadillo

SystemRequirements GNU make

Date/Publication 2018-03-25 23:07:00 UTC

R topics documented:

CP_SIR	2
distance	3
hMave	4
orthoDr_reg	5
orthoDr_surv	6

ortho_optim	8
predict	9
silverman	10
skcm.clinical	11
skcm.melgene	11
view_dr_surv	12
Index	13

CP_SIR	<i>Counting process based sliced inverse regression model</i>
--------	---

Description

The CP-SIR model for right-censored survival outcome. This model is correct only under very strong assumptions, however, since it only requires an SVD, the solution is used as the initial value in the orthoDr optimization.

Usage

```
CP_SIR(x, y, censor, bw = silverman(1, length(y)))
```

Arguments

x	A matrix for features (continuous only).
y	A vector of observed time.
censor	A vector of censoring indicator.
bw	Kernel bandwidth for nonparametric estimations (one-dimensional), the default is using Silverman's formula.

Value

A list consisting of

values	The eigenvalues of the estimation matrix
vectors	The estimated directions, ordered by eigenvalues

References

Sun, Q., Zhu, R., Wang, T. and Zeng, D. "Counting Process Based Dimension Reduction Method for Censored Outcomes." (2017) <https://arxiv.org/abs/1704.05046> .

Examples

```
# This is setting 1 in Sun et. al. (2017) with reduced sample size
library(MASS)
set.seed(1)
N = 200; P = 6
V=0.5^abs(outer(1:P, 1:P, "-"))
dataX = as.matrix(mvrnorm(N, mu=rep(0,P), Sigma=V))
failedR = as.matrix(c(1, 0.5, 0, 0, 0, rep(0, P-5)))
censorEDR = as.matrix(c(0, 0, 0, 1, 1, rep(0, P-5)))
T = rexp(N, exp(dataX %*% failedR))
C = rexp(N, exp(dataX %*% censorEDR - 1))
ndr = 1
Y = pmin(T, C)
Censor = (T < C)

# fit the model
cpsir.fit = CP_SIR(dataX, Y, Censor)
distance(failedR, cpsir.fit$vectors[, 1:ndr, drop = FALSE], "dist")
```

distance

distance correlation

Description

Calculate the distance correlation between two linear spaces

Usage

```
distance(s1, s2, type = "dist", x = NULL)
```

Arguments

s1	first space
s2	second space
type	type of distance measures: "dist" (default), "trace" or "canonical"
x	the covariate values, for canonical correlation only

Value

The distance between s1 and s2.

Examples

```
# two spaces
failedR = as.matrix(cbind(c(1, 1, 0, 0, 0, 0),
                          c(0, 0, 1, -1, 0, 0)))
B = as.matrix(cbind(c(0.1, 1.1, 0, 0, 0, 0),
                    c(0, 0, 1.1, -0.9, 0, 0)))
```

```
distance(failedR, B, "dist")
distance(failedR, B, "trace")

N=300
P=6
dataX = matrix(rnorm(N*P), N, P)
distance(failedR, B, "canonical", dataX)
```

hMave

Hazard Mave for Censored Survival Data

Description

This is an almost direct R translation of Xia, Zhang & Xu's (2010) hMave Matlab code. We implemented further options for setting a different initial value. The computational algorithm does not utilize the orthogonality constrained optimization.

Usage

```
hMave(x, y, censor, m0, B0 = NULL)
```

Arguments

x	A matrix for features.
y	A vector of observed time.
censor	A vector of censoring indicator.
m0	number of dimensions to use
B0	initial value of B. This is a feature we implemented.

Value

A list consisting of

B	The estimated B matrix
cv	Leave one out cross-validation error

References

Xia, Y., Zhang, D., & Xu, J. (2010). Dimension reduction and semiparametric estimation of survival models. *Journal of the American Statistical Association*, 105(489), 278-290. <http://dx.doi.org/10.1198/jasa.2009.tm09372>.

Examples

```
# generate some survival data
set.seed(1)
P = 7
N = 150
dataX = matrix(runif(N*P), N, P)
failedR = as.matrix(cbind(c(1, 1.3, -1.3, 1, -0.5, 0.5, -0.5, rep(0, P-7))))
T = exp(dataX %>% failedR + rnorm(N))
C = runif(N, 0, 15)
Y = pmin(T, C)
Censor = (T < C)

# fit the model
hMave.fit = hMave(dataX, Y, Censor, 1)
```

orthoDr_reg

orthoDr_reg

Description

The semiparametric dimension reduction method from Ma & Zhu (2012).

Usage

```
orthoDr_reg(x, y, method = "sir", ndr = 2, B.initial = NULL, bw = NULL,
  keep.data = FALSE, control = list(), maxitr = 500, verbose = FALSE,
  ncore = 0)
```

Arguments

x	A matrix or data.frame for features (continuous only). The algorithm will not scale the columns to unit variance
y	A vector of continuous outcome
method	Dimension reduction methods (semi-): sir, save, phd, local or seff. Currently only sir and phd are available.
ndr	The number of directions
B.initial	Initial B values. If specified, must be a matrix with ncol(x) rows and ndr columns. Will be processed by Gram-Schmidt if not orthogonal. If the initial value is not given, three initial values (sir, save and phd) using the traditional method will be tested. The one with smallest l2 norm of the estimating equation will be used.
bw	A Kernel bandwidth, assuming each variables have unit variance
keep.data	Should the original data be kept for prediction. Default is FALSE
control	A list of tuning variables for optimization. epsilon is the size for numerically approximating the gradient. For others, see Wen and Yin (2013).

maxitr	Maximum number of iterations
verbose	Should information be displayed
ncore	Number of cores for parallel computing. The default is the maximum number of threads.

Value

	A orthoDr object; a list consisting of
B	The optimal B value
fn	The final functional value
itr	The number of iterations
converge	convergence code

References

- Ma, Y., & Zhu, L. (2012). A semiparametric approach to dimension reduction. *Journal of the American Statistical Association*, 107(497), 168-179. DOI: <https://doi.org/10.1080/01621459.2011.646925>.
- Ma, Y., & Zhu, L. (2013). Efficient estimation in sufficient dimension reduction. *Annals of statistics*, 41(1), 250. DOI: 10.1214/12-AOS1072 <https://projecteuclid.org/euclid.aos/1364302742>
- Wen, Z. and Yin, W., "A feasible method for optimization with orthogonality constraints." *Mathematical Programming* 142.1-2 (2013): 397-434. DOI: <https://doi.org/10.1007/s10107-012-0584-1>.

Examples

```
# generate some regression data
set.seed(1)
N = 100; P = 4; dataX = matrix(rnorm(N*P), N, P)
Y = -1 + dataX[,1] + rnorm(N)
# fit the semi-sir model
orthoDr_reg(dataX, Y, ndr = 1, method = "sir")
# fit the semi-phd model
Y = -1 + dataX[,1]^2 + rnorm(N)
orthoDr_reg(dataX, Y, ndr = 1, method = "phd")
```

orthoDr_surv

IR-CP model

Description

The counting process based semiparametric dimension reduction (IR-CP) model for right censored survival outcome.

Usage

```
orthoDr_surv(x, y, censor, method = "dm", ndr = ifelse(method == "forward",
  1, 2), B.initial = NULL, bw = NULL, keep.data = FALSE,
  control = list(), maxitr = 500, verbose = FALSE, ncore = 0)
```

Arguments

x	A matrix or data.frame for features. The algorithm will not scale the columns to unit variance
y	A vector of observed time
censor	A vector of censoring indicator
method	Which estimating equation to use: should be forward (1-d model), dn (counting process) or dm (martingale)
ndr	The number of directions
B.initial	Initial B values. Will use the counting process based SIR model CP_SIR as the initial if leaving as NULL. If specified, must be a matrix with <code>ncol(x)</code> rows and <code>ndr</code> columns. Will be processed by Gram-Schmidt if not orthogonal
bw	A Kernel bandwidth, assuming each variables have unit variance
keep.data	Should the original data be kept for prediction. Default is FALSE
control	A list of tuning variables for optimization. <code>epsilon</code> is the size for numerically approximating the gradient. For others, see Wen and Yin (2013).
maxitr	Maximum number of iterations
verbose	Should information be displayed
ncore	Number of cores for parallel computing. The default is the maximum number of threads.

Value

A orthoDr object; a list consisting of	
B	The optimal B value
fn	The final functional value
itr	The number of iterations
converge	convergence code

References

- Sun, Q., Zhu, R., Wang, T. and Zeng, D. "Counting Process Based Dimension Reduction Method for Censored Outcomes." (2017) DOI: <https://arxiv.org/abs/1704.05046>.
- Wen, Z. and Yin, W., "A feasible method for optimization with orthogonality constraints." *Mathematical Programming* 142.1-2 (2013): 397-434. DOI: <https://doi.org/10.1007/s10107-012-0584-1>

Examples

```

# This is setting 1 in Sun et. al. (2017) with reduced sample size
library(MASS)
set.seed(1)
N = 200; P = 6
V=0.5^abs(outer(1:P, 1:P, "-"))
dataX = as.matrix(mvrnorm(N, mu=rep(0,P), Sigma=V))
failedR = as.matrix(c(1, 0.5, 0, 0, 0, rep(0, P-5)))
censorEDR = as.matrix(c(0, 0, 0, 1, 1, rep(0, P-5)))
T = rexp(N, exp(dataX %*% failedR))
C = rexp(N, exp(dataX %*% censorEDR - 1))
ndr = 1
Y = pmin(T, C)
Censor = (T < C)

# fit the model
forward.fit = orthoDr_surv(dataX, Y, Censor, method = "forward")
distance(failedR, forward.fit$B, "dist")

dn.fit = orthoDr_surv(dataX, Y, Censor, method = "dn", ndr = ndr)
distance(failedR, dn.fit$B, "dist")

dm.fit = orthoDr_surv(dataX, Y, Censor, method = "dm", ndr = ndr)
distance(failedR, dm.fit$B, "dist")

```

ortho_optim

*Orthogonality constrained optimization***Description**

A general purpose optimization solver with orthogonality constraint. The orthogonality constrained optimization method is a nearly direct translation from Wen and Yin (2010)'s Matlab code.

Usage

```
ortho_optim(B, fn, grad = NULL, ..., maximize = FALSE, control = list(),
           maxitr = 500, verbose = FALSE)
```

Arguments

B	Initial B values. Must be a matrix, and the columns are subject to the orthogonality constrains. Will be processed by Gram-Schmidt if not orthogonal
fn	A function that calculate the objective function value. The first argument should be B. Returns a single value.
grad	A function that calculate the gradient. The first argument should be B. Returns a matrix with the same dimension as B. If not specified, then numerical approximation is used.
...	Arguments passed to fn and grad

maximize	By default, the solver will try to minimize the objective function unless maximize = TRUE
control	A list of tuning variables for optimization. epsilon is the size for numerically approximating the gradient. For others, see Wen and Yin (2013).
maxitr	Maximum number of iterations
verbose	Should information be displayed

Value

A orthoDr object; a list consisting of	
B	The optimal B value
fn	The final functional value
itr	The number of iterations
converge	convergence code

References

Wen, Z. and Yin, W., "A feasible method for optimization with orthogonality constraints." *Mathematical Programming* 142.1-2 (2013): 397-434. DOI: <https://doi.org/10.1007/s10107-012-0584-1>

Examples

```
# an eigen value problem
library(pracma)
set.seed(1)
n = 100; k = 6
A = matrix(rnorm(n*n), n, n)
A = t(A) %*% A
B = gramSchmidt(matrix(rnorm(n*k), n, k))$Q
fx <- function(B, A) -0.5 * sum(diag(t(B) %*% A %*% B ))
gx <- function(B, A) -A %*% B
fit = ortho_optim(B, fx, gx, A = A)
fx(fit$B, A)

# compare with the solution from the eigen function
sol = eigen(A)$vectors[, 1:k]
fx(sol, A)
```

predict

predict.orthoDr

Description

The prediction function for orthoDr fitted models

Usage

```
## S3 method for class 'orthoDr'
predict(object, testx, ...)
```

Arguments

object	A fitted orthoDr object
testx	Testing data
...	...

Value

The predicted object

Examples

```
# generate some survival data
N = 100; P = 4; dataX = matrix(rnorm(N*P), N, P)
Y = exp(-1 + dataX[,1] + rnorm(N))
Censor = rbinom(N, 1, 0.8)

# fit the model with keep.data = TRUE
orthoDr.fit = orthoDr_surv(dataX, Y, Censor, ndr = 1, method = "dm", keep.data = TRUE)

#predict 10 new observations
predict(orthoDr.fit, matrix(rnorm(10*P), 10, P))
```

silverman

A simple Silverman bandwidth formula

Description

Silverman bandwidth

Usage

```
silverman(d, n)
```

Arguments

d	Number of dimension
n	Number of observation

Value

A simple bandwidth choice

Examples

```
silverman(1, 300)
```

```
skcm.clinical
```

```
skcm.clinical
```

Description

The clinical variables of the SKCM dataset. The original data was obtained from The Cancer Genome Atlas (TCGA).

Usage

```
skcm.clinical
```

Format

Contains 469 subjects with 156 failures. Each row contains one subject, subject ID is indicated by row name. Variables include Time, Censor, Gender and Age. Age has 8 missing values.

References

<https://cancergenome.nih.gov/>

```
skcm.melgene
```

```
skcm.melgene
```

Description

The expression of top 20 genes of cutaneous melanoma literature based on the MelGene Database.

Usage

```
skcm.melgene
```

Format

Each row contains one subject, subject ID is indicated by row name. Gene names in the columns. The columns are scaled.

References

Chatzinasiou, Foteini, Christina M. Lill, Katerina Kypreou, Irene Stefanaki, Vasiliki Nicolaou, George Spyrou, Evangelos Evangelou et al. "Comprehensive field synopsis and systematic meta-analyses of genetic association studies in cutaneous melanoma." *Journal of the National Cancer Institute* 103, no. 16 (2011): 1227-1235.

<http://bioserver-3.bioacademy.gr/Bioserver/MelGene/>

<https://cancergenome.nih.gov/>

 view_dr_surv

2D or 2D view of survival data on reduced dimension

Description

Produce 2D or 3D plots of right censored survival data based on a given dimension reduction space

Usage

```
view_dr_surv(x, y, censor, B = NULL, bw = NULL, FUN = "log",
  type = "2D", legend.add = TRUE, xlab = "Reduced Direction",
  ylab = "Time", zlab = "Survival")
```

Arguments

x	A matrix or data.frame for features (continuous only). The algorithm will not scale the columns to unit variance
y	A vector of observed time
censor	A vector of censoring indicator
B	The dimension reduction subspace, can only be 1 dimensional
bw	A Kernel bandwidth (3D plot only) for approximating the survival function, default is the Silverman's formula
FUN	A scaling function applied to the time points y. Default is "log".
type	2D or 3D plot
legend.add	Should legend be added (2D plot only)
xlab	x axis label
ylab	y axis label
zlab	z axis label

References

Sun, Q., Zhu, R., Wang, T. and Zeng, D. "Counting Process Based Dimension Reduction Method for Censored Outcomes." (2017) <https://arxiv.org/abs/1704.05046>.

Examples

```
# generate some survival data
N = 100; P = 4; dataX = matrix(rnorm(N*P), N, P)
Y = exp(-1 + dataX[,1] + rnorm(N))
Censor = rbinom(N, 1, 0.8)

orthoDr.fit = orthoDr_surv(dataX, Y, Censor, ndr = 1, method = "dm")
view_dr_surv(dataX, Y, Censor, orthoDr.fit$B)
```

Index

*Topic **skcm.clinical**
skcm.clinical, [11](#)

*Topic **skcm.melgene**
skcm.melgene, [11](#)

CP_SIR, [2](#), [7](#)

distance, [3](#)

hMave, [4](#)

ortho_optim, [8](#)
orthoDr_reg, [5](#)
orthoDr_surv, [6](#)

predict, [9](#)

silverman, [10](#)
skcm.clinical, [11](#)
skcm.melgene, [11](#)

view_dr_surv, [12](#)