

Package ‘abess’

May 7, 2026

Type Package

Title Fast Best Subset Selection

Version 0.4.11

Date 2025-12-07

Maintainer Jin Zhu <zhuji1jqx@gmail.com>

Description Extremely efficient toolkit for solving the best subset selection problem <<https://www.jmlr.org/papers/v23/21-1060.html>>. This package is its R interface. The package implements and generalizes algorithms designed in <[doi:10.1073/pnas.2014241117](https://doi.org/10.1073/pnas.2014241117)> that exploits a novel sequencing-and-splicing technique to guarantee exact support recovery and globally optimal solution in polynomial times for linear model. It also supports best subset selection for logistic regression, Poisson regression, Cox proportional hazard model, Gamma regression, multiple-response regression, multinomial logistic regression, ordinal regression, Ising model reconstruction <[doi:10.1080/01621459.2025.2571245](https://doi.org/10.1080/01621459.2025.2571245)>, (sequential) principal component analysis, and robust principal component analysis. The other valuable features such as the best subset of group selection <[doi:10.1287/ijoc.2022.1241](https://doi.org/10.1287/ijoc.2022.1241)> and sure independence screening <[doi:10.1111/j.1467-9868.2008.00674.x](https://doi.org/10.1111/j.1467-9868.2008.00674.x)> are also provided.

License GPL (>= 3) | file LICENSE

Encoding UTF-8

LazyData true

Depends R (>= 3.1.0)

Imports Rcpp, MASS, methods, Matrix, igraph

LinkingTo Rcpp, RcppEigen

RoxygenNote 7.3.3

Suggests testthat, knitr, rmarkdown

VignetteBuilder knitr

URL <https://github.com/abess-team/abess>,
<https://abess-team.github.io/abess/>,
<https://abess.readthedocs.io>

BugReports <https://github.com/abess-team/abess/issues>

NeedsCompilation yes

Author Jin Zhu [aut, cre] (ORCID: <<https://orcid.org/0000-0001-8550-5822>>),
 Zezhi Wang [aut],
 Liyuan Hu [aut],
 Junhao Huang [aut],
 Kangkang Jiang [aut],
 Yanhang Zhang [aut],
 Borui Tang [aut],
 Shiyun Lin [aut],
 Junxian Zhu [aut],
 Canhong Wen [aut],
 Heping Zhang [aut] (ORCID: <<https://orcid.org/0000-0002-0688-4076>>),
 Xueqin Wang [aut] (ORCID: <<https://orcid.org/0000-0001-5205-9950>>),
 spectra contributors [cph] (Spectra implementation)

Repository CRAN

Date/Publication 2025-12-08 09:50:15 UTC

Contents

abess.default	3
abesspca	10
abessrpca	15
coef.abess	18
coef.abesspca	19
coef.abessrpca	20
deviance.abess	20
extract	21
generate.bmn.data	22
generate.data	23
generate.matrix	26
generate.spc.matrix	27
plot.abess	29
plot.abesspca	30
plot.abessrpca	31
predict.abess	31
print.abess	32
print.abesspca	33
print.abessrpca	34
slide	34
trim32	36

Index

37

abess.default	<i>Adaptive best subset selection (for generalized linear model)</i>
---------------	--

Description

Adaptive best-subset selection for regression, (multi-class) classification, counting-response, censored-response, positive response, multi-response modeling in polynomial times.

Usage

```
## Default S3 method:
abess(
  x,
  y,
  family = c("gaussian", "binomial", "poisson", "cox", "mgaussian", "multinomial",
    "gamma", "ordinal"),
  tune.path = c("sequence", "gsection"),
  tune.type = c("gic", "ebic", "bic", "aic", "cv"),
  weight = NULL,
  normalize = NULL,
  fit.intercept = TRUE,
  beta.low = -.Machine$double.xmax,
  beta.high = .Machine$double.xmax,
  c.max = 2,
  support.size = NULL,
  gs.range = NULL,
  lambda = 0,
  always.include = NULL,
  group.index = NULL,
  init.active.set = NULL,
  splicing.type = 2,
  max.splicing.iter = 20,
  screening.num = NULL,
  important.search = NULL,
  warm.start = TRUE,
  nfolds = 5,
  foldid = NULL,
  cov.update = FALSE,
  newton = c("exact", "approx"),
  newton.thresh = 1e-06,
  max.newton.iter = NULL,
  early.stop = FALSE,
  ic.scale = 1,
  num.threads = 0,
  seed = 1,
  ...
)
```

```
## S3 method for class 'formula'
abess(formula, data, subset, na.action, ...)
```

Arguments

<code>x</code>	Input matrix, of dimension $n \times p$; each row is an observation vector and each column is a predictor/feature/variable. Can be in sparse matrix format (inherit from class "dgCMatrix" in package Matrix).
<code>y</code>	The response variable, of n observations. For <code>family = "binomial"</code> should have two levels. For <code>family = "poisson"</code> , <code>y</code> should be a vector with positive integer. For <code>family = "cox"</code> , <code>y</code> should be a Surv object returned by the survival package (recommended) or a two-column matrix with columns named "time" and "status". For <code>family = "mgaussian"</code> , <code>y</code> should be a matrix of quantitative responses. For <code>family = "multinomial"</code> or "ordinal", <code>y</code> should be a factor of at least three levels. Note that, for either "binomial", "ordinal" or "multinomial", if <code>y</code> is presented as a numerical vector, it will be coerced into a factor.
<code>family</code>	One of the following models: "gaussian" (continuous response), "binomial" (binary response), "poisson" (non-negative count), "cox" (left-censored response), "mgaussian" (multivariate continuous response), "multinomial" (multi-class response), "ordinal" (multi-class ordinal response), "gamma" (positive continuous response). Depending on the response. Any unambiguous substring can be given.
<code>tune.path</code>	The method to be used to select the optimal support size. For <code>tune.path = "sequence"</code> , we solve the best subset selection problem for each size in <code>support.size</code> . For <code>tune.path = "gsection"</code> , we solve the best subset selection problem with support size ranged in <code>gs.range</code> , where the specific support size to be considered is determined by golden section.
<code>tune.type</code>	The type of criterion for choosing the support size. Available options are "gic", "ebic", "bic", "aic" and "cv". Default is "gic".
<code>weight</code>	Observation weights. When <code>weight = NULL</code> , we set <code>weight = 1</code> for each observation as default.
<code>normalize</code>	Options for normalization. <code>normalize = 0</code> for no normalization. <code>normalize = 1</code> for subtracting the means of the columns of <code>x</code> and <code>y</code> , and also normalizing the columns of <code>x</code> to have \sqrt{n} norm. <code>normalize = 2</code> for subtracting the mean of columns of <code>x</code> and scaling the columns of <code>x</code> to have \sqrt{n} norm. <code>normalize = 3</code> for scaling the columns of <code>x</code> to have \sqrt{n} norm. If <code>normalize = NULL</code> , <code>normalize</code> will be set 1 for "gaussian" and "mgaussian", 3 for "cox". Default is <code>normalize = NULL</code> .
<code>fit.intercept</code>	A boolean value indicating whether to fit an intercept. We assume the data has been centered if <code>fit.intercept = FALSE</code> . Default: <code>fit.intercept = TRUE</code> .
<code>beta.low</code>	A single value specifying the lower bound of β . Default is <code>.Machine\$double.xmax</code> .
<code>beta.high</code>	A single value specifying the upper bound of β . Default is <code>.Machine\$double.xmax</code> .
<code>c.max</code>	an integer splicing size. Default is: <code>c.max = 2</code> .

<code>support.size</code>	An integer vector representing the alternative support sizes. Only used for <code>tune.path = "sequence"</code> . Default is $0:\min(n, \text{round}(n/(\log(\log(n))\log(p))))$.
<code>gs.range</code>	A integer vector with two elements. The first element is the minimum model size considered by golden-section, the later one is the maximum one. Default is <code>gs.range = c(1, min(n, round(n/(log(log(n))log(p)))))</code> .
<code>lambda</code>	A single lambda value for regularized best subset selection. Default is 0.
<code>always.include</code>	An integer vector containing the indexes of variables that should always be included in the model.
<code>group.index</code>	A vector of integers indicating the which group each variable is in. For variables in the same group, they should be located in adjacent columns of <code>x</code> and their corresponding index in <code>group.index</code> should be the same. Denote the first group as 1, the second 2, etc. If you do not fit a model with a group structure, please set <code>group.index = NULL</code> (the default).
<code>init.active.set</code>	A vector of integers indicating the initial active set. Default: <code>init.active.set = NULL</code> .
<code>splicing.type</code>	Optional type for splicing. If <code>splicing.type = 1</code> , the number of variables to be spliced is <code>c.max, ..., 1</code> ; if <code>splicing.type = 2</code> , the number of variables to be spliced is <code>c.max, c.max/2, ..., 1</code> . (Default: <code>splicing.type = 2</code> .)
<code>max.splicing.iter</code>	The maximum number of performing splicing algorithm. In most of the case, only a few times of splicing iteration can guarantee the convergence. Default is <code>max.splicing.iter = 20</code> .
<code>screening.num</code>	An integer number. Preserve <code>screening.num</code> number of predictors with the largest marginal maximum likelihood estimator before running algorithm.
<code>important.search</code>	An integer number indicating the number of important variables to be splicing. When <code>important.search</code> \ll <code>p</code> variables, it would greatly reduce runtimes. Default: <code>important.search = 128</code> .
<code>warm.start</code>	Whether to use the last solution as a warm start. Default is <code>warm.start = TRUE</code> .
<code>nfolds</code>	The number of folds in cross-validation. Default is <code>nfolds = 5</code> .
<code>foldid</code>	an optional integer vector of values between 1, ..., <code>nfolds</code> identifying what fold each observation is in. The default <code>foldid = NULL</code> would generate a random <code>foldid</code> .
<code>cov.update</code>	A logical value only used for <code>family = "gaussian"</code> . If <code>cov.update = TRUE</code> , use a covariance-based implementation; otherwise, a naive implementation. The naive method is more computational efficient than covariance-based method when $p \gg n$ and <code>important.search</code> is much large than its default value. Default: <code>cov.update = FALSE</code> .
<code>newton</code>	A character specify the Newton's method for fitting generalized linear models, it should be either <code>newton = "exact"</code> or <code>newton = "approx"</code> . If <code>newton = "exact"</code> , then the exact hessian is used, while <code>newton = "approx"</code> uses diagonal entry of the hessian, and can be faster (especially when <code>family = "cox"</code>).
<code>newton.thresh</code>	a numeric value for controlling positive convergence tolerance. The Newton's iterations converge when $ dev - dev_{old} /(dev + 0.1) < \text{newton.thresh}$.

<code>max.newton.iter</code>	a integer giving the maximal number of Newton's iteration iterations. Default is <code>max.newton.iter = 10</code> if <code>newton = "exact"</code> , and <code>max.newton.iter = 60</code> if <code>newton = "approx"</code> .
<code>early.stop</code>	A boolean value decide whether early stopping. If <code>early.stop = TRUE</code> , algorithm will stop if the last tuning value less than the existing one. Default: <code>early.stop = FALSE</code> .
<code>ic.scale</code>	A non-negative value used for multiplying the penalty term in information criterion. Default: <code>ic.scale = 1</code> .
<code>num.threads</code>	An integer decide the number of threads to be concurrently used for cross-validation (i.e., <code>tune.type = "cv"</code>). If <code>num.threads = 0</code> , then all of available cores will be used. Default: <code>num.threads = 0</code> .
<code>seed</code>	Seed to be used to divide the sample into cross-validation folds. Default is <code>seed = 1</code> .
<code>...</code>	further arguments to be passed to or from methods.
<code>formula</code>	an object of class "formula": a symbolic description of the model to be fitted. The details of model specification are given in the "Details" section of " formula ".
<code>data</code>	a data frame containing the variables in the formula.
<code>subset</code>	an optional vector specifying a subset of observations to be used.
<code>na.action</code>	a function which indicates what should happen when the data contain NAs. Defaults to <code>getOption("na.action")</code> .

Details

Best-subset selection aims to find a small subset of predictors, so that the resulting model is expected to have the most desirable prediction accuracy. Best-subset selection problem under the support size s is

$$\min_{\beta} -2 \log L(\beta) \quad \text{s.t.} \quad \|\beta\|_0 \leq s,$$

where $L(\beta)$ is arbitrary convex functions. In the GLM case, $\log L(\beta)$ is the log-likelihood function; in the Cox model, $\log L(\beta)$ is the log partial-likelihood function. The best subset selection problem is solved by the splicing algorithm in this package, see Zhu (2020) for details. Under mild conditions, the algorithm exactly solve this problem in polynomial time. This algorithm exploits the idea of sequencing and splicing to reach a stable solution in finite steps when s is fixed. The parameters `c.max`, `splicing.type` and `max.splicing.iter` allow user control the splicing technique flexibly. On the basis of our numerical experiment results, we assign properly parameters to the these parameters as the default such that the precision and runtime are well balanced, we suggest users keep the default values unchanged. Please see [this online page](#) for more details about the splicing algorithm.

To find the optimal support size s , we provide various criterion like GIC, AIC, BIC and cross-validation error to determine it. More specifically, the sequence of models implied by `support.size` are fit by the splicing algorithm. And the solved model with least information criterion or cross-validation error is the optimal model. The sequential searching for the optimal model is somehow time-wasting. A faster strategy is golden section (GS), which only need to specify `gs.range`. More details about GS is referred to Zhang et al (2021).

It is worthy to note that the parameters `newton`, `max.newton.iter` and `newton.thresh` allows user control the parameter estimation in non-gaussian models. The parameter estimation procedure use Newton method or approximated Newton method (only consider the diagonal elements in the Hessian matrix). Again, we suggest to use the default values unchanged because the same reason for the parameter `c.max`.

abess support some well-known advanced statistical methods to analyze data, including

- sure independent screening: helpful for ultra-high dimensional predictors (i.e., $p \gg n$). Use the parameter `screening.num` to retain the marginally most important predictors. See Fan et al (2008) for more details.
- best subset of group selection: helpful when predictors have group structure. Use the parameter `group.index` to specify the group structure of predictors. See Zhang et al (2021) for more details.
- l_2 regularization best subset selection: helpful when signal-to-ratio is relatively small. Use the parameter `lambda` to control the magnitude of the regularization term.
- nuisance selection: helpful when the prior knowledge of important predictors is available. Use the parameter `always.include` to retain the important predictors.

The arbitrary combination of the four methods are definitely support. Please see [online vignettes](#) for more details about the advanced features support by abess.

Value

A S3 abess class object, which is a list with the following components:

<code>beta</code>	A p -by- <code>length(support.size)</code> matrix of coefficients for univariate family, stored in column format; while a list of <code>length(support.size)</code> coefficients matrix (with size p -by- <code>ncol(y)</code>) for multivariate family.
<code>intercept</code>	An intercept vector of length <code>length(support.size)</code> for univariate family; while a list of <code>length(support.size)</code> intercept vector (with size <code>ncol(y)</code>) for multivariate family.
<code>dev</code>	the deviance of length <code>length(support.size)</code> .
<code>tune.value</code>	A value of tuning criterion of length <code>length(support.size)</code> .
<code>nobs</code>	The number of sample used for training.
<code>nvars</code>	The number of variables used for training.
<code>family</code>	Type of the model.
<code>tune.path</code>	The path type for tuning parameters.
<code>support.size</code>	The actual <code>support.size</code> values used. Note that it is not necessary the same as the input if the later have non-integer values or duplicated values.
<code>edf</code>	The effective degree of freedom. It is the same as <code>support.size</code> when <code>lambda = 0</code> .
<code>best.size</code>	The best support size selected by the tuning value.
<code>tune.type</code>	The criterion type for tuning parameters.
<code>tune.path</code>	The strategy for tuning parameters.
<code>screening.vars</code>	The character vector specify the feature selected by feature screening. It would be an empty character vector if <code>screening.num = 0</code> .
<code>call</code>	The original call to abess.

Author(s)

Jin Zhu, Junxian Zhu, Canhong Wen, Heping Zhang, Xueqin Wang

References

A polynomial algorithm for best-subset selection problem. Junxian Zhu, Canhong Wen, Jin Zhu, Heping Zhang, Xueqin Wang. Proceedings of the National Academy of Sciences Dec 2020, 117 (52) 33117-33123; doi:10.1073/pnas.2014241117

A Splicing Approach to Best Subset of Groups Selection. Zhang, Yanhang, Junxian Zhu, Jin Zhu, and Xueqin Wang (2023). INFORMS Journal on Computing, 35:1, 104-119. doi:10.1287/ijoc.2022.1241

abess: A Fast Best-Subset Selection Library in Python and R. Zhu Jin, Xueqin Wang, Liyuan Hu, Junhao Huang, Kangkang Jiang, Yanhang Zhang, Shiyun Lin, and Junxian Zhu. Journal of Machine Learning Research 23, no. 202 (2022): 1-7.

Sure independence screening for ultrahigh dimensional feature space. Fan, J. and Lv, J. (2008), Journal of the Royal Statistical Society: Series B (Statistical Methodology), 70: 849-911. doi:10.1111/j.14679868.2008.00674.x

Targeted Inference Involving High-Dimensional Data Using Nuisance Penalized Regression. Qiang Sun & Heping Zhang (2020). Journal of the American Statistical Association, doi:10.1080/01621459.2020.1737079

See Also

[print.abess](#), [predict.abess](#), [coef.abess](#), [extract.abess](#), [plot.abess](#), [deviance.abess](#).

Examples

```
library(abess)
Sys.setenv("OMP_THREAD_LIMIT" = 2)
n <- 100
p <- 20
support.size <- 3

##### linear model #####
dataset <- generate.data(n, p, support.size)
abess_fit <- abess(dataset[["x"]], dataset[["y"]])
## helpful generic functions:
print(abess_fit)
coef(abess_fit, support.size = 3)
predict(abess_fit,
  newx = dataset[["x"]][1:10, ],
  support.size = c(3, 4)
)
str(extract(abess_fit, 3))
deviance(abess_fit)
plot(abess_fit)
plot(abess_fit, type = "tune")

##### logistic model #####
dataset <- generate.data(n, p, support.size, family = "binomial")
```

```

## allow cross-validation to tuning
abess_fit <- abess(dataset[["x"]], dataset[["y"]],
  family = "binomial", tune.type = "cv"
)
abess_fit

##### poisson model #####
dataset <- generate.data(n, p, support.size, family = "poisson")
abess_fit <- abess(dataset[["x"]], dataset[["y"]],
  family = "poisson", tune.type = "cv"
)
abess_fit

##### Cox model #####
dataset <- generate.data(n, p, support.size, family = "cox")
abess_fit <- abess(dataset[["x"]], dataset[["y"]],
  family = "cox", tune.type = "cv"
)

##### Multivariate gaussian model #####
dataset <- generate.data(n, p, support.size, family = "mgaussian")
abess_fit <- abess(dataset[["x"]], dataset[["y"]],
  family = "mgaussian", tune.type = "cv"
)
plot(abess_fit, type = "l2norm")

##### Multinomial model (multi-classification) #####
dataset <- generate.data(n, p, support.size, family = "multinomial")
abess_fit <- abess(dataset[["x"]], dataset[["y"]],
  family = "multinomial", tune.type = "cv"
)
predict(abess_fit,
  newx = dataset[["x"]][1:10, ],
  support.size = c(3, 4), type = "response"
)

##### Ordinal regression #####
dataset <- generate.data(n, p, support.size, family = "ordinal", class.num = 4)
abess_fit <- abess(dataset[["x"]], dataset[["y"]],
  family = "ordinal", tune.type = "cv"
)
coef <- coef(abess_fit, support.size = abess_fit[["best.size"]][1])
predict(abess_fit,
  newx = dataset[["x"]][1:10, ],
  support.size = c(3, 4), type = "response"
)

##### Best group subset selection #####
dataset <- generate.data(n, p, support.size)
group_index <- rep(1:10, each = 2)
abess_fit <- abess(dataset[["x"]], dataset[["y"]], group.index = group_index)
str(extract(abess_fit))

```

```
##### Golden section searching #####
dataset <- generate.data(n, p, support.size)
abess_fit <- abess(dataset[["x"]], dataset[["y"]], tune.path = "gsection")
abess_fit

##### Feature screening #####
p <- 1000
dataset <- generate.data(n, p, support.size)
abess_fit <- abess(dataset[["x"]], dataset[["y"]],
  screening.num = 100
)
str(extract(abess_fit))

##### Sparse predictor #####
require(Matrix)
p <- 1000
dataset <- generate.data(n, p, support.size)
dataset[["x"]][abs(dataset[["x"]]) < 1] <- 0
dataset[["x"]] <- Matrix(dataset[["x"]])
abess_fit <- abess(dataset[["x"]], dataset[["y"]])
str(extract(abess_fit))

##### Formula interface #####
data("trim32")
abess_fit <- abess(y ~ ., data = trim32)
abess_fit
```

abesspca

Adaptive best subset selection for principal component analysis

Description

Adaptive best subset selection for principal component analysis

Usage

```
abesspca(
  x,
  type = c("predictor", "gram"),
  sparse.type = c("fpc", "kpc"),
  cor = FALSE,
  kpc.num = NULL,
  support.size = NULL,
  gs.range = NULL,
  tune.path = c("sequence", "gsection"),
  tune.type = c("gic", "aic", "bic", "ebic", "cv"),
  nfolds = 5,
```

```

    foldid = NULL,
    ic.scale = 1,
    c.max = NULL,
    always.include = NULL,
    group.index = NULL,
    screening.num = NULL,
    splicing.type = 1,
    max.splicing.iter = 20,
    warm.start = TRUE,
    num.threads = 0,
    ...
)

```

Arguments

<code>x</code>	A matrix object. It can be either a predictor matrix where each row is an observation and each column is a predictor or a sample covariance/correlation matrix. If <code>x</code> is a predictor matrix, it can be in sparse matrix format (inherit from class "dgCMatrix" in package Matrix).
<code>type</code>	If <code>type = "predictor"</code> , <code>x</code> is considered as the predictor matrix. If <code>type = "gram"</code> , <code>x</code> is considered as a sample covariance or correlation matrix.
<code>sparse.type</code>	If <code>sparse.type = "fpc"</code> , then best subset selection performs on the first principal component; If <code>sparse.type = "kpc"</code> , then best subset selection would be sequentially performed on the first <code>kpc.num</code> number of principal components. If <code>kpc.num</code> is supplied, the default is <code>sparse.type = "kpc"</code> ; otherwise, is <code>sparse.type = "fpc"</code> .
<code>cor</code>	A logical value. If <code>cor = TRUE</code> , perform PCA on the correlation matrix; otherwise, the covariance matrix. This option is available only if <code>type = "predictor"</code> . Default: <code>cor = FALSE</code> .
<code>kpc.num</code>	A integer decide the number of principal components to be sequentially considered.
<code>support.size</code>	It is a flexible input. If it is an integer vector. It represents the support sizes to be considered for each principal component. If it is a list object containing <code>kpc.num</code> number of integer vectors, the <i>i</i> -th principal component consider the support size specified in the <i>i</i> -th element in the list. Only used for <code>tune.path = "sequence"</code> . The default is <code>support.size = NULL</code> , and some rules in details section are used to specify <code>support.size</code> .
<code>gs.range</code>	A integer vector with two elements. The first element is the minimum model size considered by golden-section, the later one is the maximum one. Default is <code>gs.range = c(1, min(n, round(n/(log(log(n))log(p))))))</code> .
<code>tune.path</code>	The method to be used to select the optimal support size. For <code>tune.path = "sequence"</code> , we solve the best subset selection problem for each size in <code>support.size</code> . For <code>tune.path = "gsection"</code> , we solve the best subset selection problem with support size ranged in <code>gs.range</code> , where the specific support size to be considered is determined by golden section.

<code>tune.type</code>	The type of criterion for choosing the support size. Available options are "gic", "ebic", "bic", "aic" and "cv". Default is "gic". <code>tune.type = "cv"</code> is available only when <code>type = "predictor"</code> .
<code>nfolds</code>	The number of folds in cross-validation. Default is <code>nfolds = 5</code> .
<code>foldid</code>	an optional integer vector of values between 1, ..., <code>nfolds</code> identifying what fold each observation is in. The default <code>foldid = NULL</code> would generate a random <code>foldid</code> .
<code>ic.scale</code>	A non-negative value used for multiplying the penalty term in information criterion. Default: <code>ic.scale = 1</code> .
<code>c.max</code>	an integer splicing size. The default of <code>c.max</code> is the maximum of 2 and <code>max(support.size) / 2</code> .
<code>always.include</code>	An integer vector containing the indexes of variables that should always be included in the model.
<code>group.index</code>	A vector of integers indicating the which group each variable is in. For variables in the same group, they should be located in adjacent columns of <code>x</code> and their corresponding index in <code>group.index</code> should be the same. Denote the first group as 1, the second 2, etc. If you do not fit a model with a group structure, please set <code>group.index = NULL</code> (the default).
<code>screening.num</code>	An integer number. Preserve <code>screening.num</code> number of predictors with the largest marginal maximum likelihood estimator before running algorithm.
<code>splicing.type</code>	Optional type for splicing. If <code>splicing.type = 1</code> , the number of variables to be spliced is <code>c.max, ..., 1</code> ; if <code>splicing.type = 2</code> , the number of variables to be spliced is <code>c.max, c.max/2, ..., 1</code> . Default: <code>splicing.type = 1</code> .
<code>max.splicing.iter</code>	The maximum number of performing splicing algorithm. In most of the case, only a few times of splicing iteration can guarantee the convergence. Default is <code>max.splicing.iter = 20</code> .
<code>warm.start</code>	Whether to use the last solution as a warm start. Default is <code>warm.start = TRUE</code> .
<code>num.threads</code>	An integer decide the number of threads to be concurrently used for cross-validation (i.e., <code>tune.type = "cv"</code>). If <code>num.threads = 0</code> , then all of available cores will be used. Default: <code>num.threads = 0</code> .
<code>...</code>	further arguments to be passed to or from methods.

Details

Adaptive best subset selection for principal component analysis (abessPCA) aim to solve the non-convex optimization problem:

$$- \arg \min_v v^\top \Sigma v, s.t. \quad v^\top v = 1, \|v\|_0 \leq s,$$

where s is support size. Here, Σ is covariance matrix, i.e.,

$$\Sigma = \frac{1}{n} X^\top X.$$

A generic splicing technique is implemented to solve this problem. By exploiting the warm-start initialization, the non-convex optimization problem at different support size (specified by `support.size`) can be efficiently solved.

The abessPCA can be conducted sequentially for each component. Please see the multiple principal components Section on the [website](#) for more details about this function. For abesspca function, the arguments `kpc.num` control the number of components to be considered.

When `sparse.type = "fpc"` but `support.size` is not supplied, it is set as `support.size = 1:min(ncol(x), 100)` if `group.index = NULL`; otherwise, `support.size = 1:min(length(unique(group.index)), 100)`. When `sparse.type = "kpc"` but `support.size` is not supplied, then for `20` it is set as `min(ncol(x), 100)` if `group.index = NULL`; otherwise, `min(length(unique(group.index)), 100)`.

Value

A S3 abesspca class object, which is a list with the following components:

<code>coef</code>	A p -by- <code>length(support.size)</code> loading matrix of sparse principal components (PC), where each row is a variable and each column is a support size;
<code>nvars</code>	The number of variables.
<code>sparse.type</code>	The same as input.
<code>support.size</code>	The actual support.size values used. Note that it is not necessary the same as the input if the later have non-integer values or duplicated values.
<code>ev</code>	A vector with size <code>length(support.size)</code> . It records the cumulative sums of explained variance at each support size.
<code>tune.value</code>	A value of tuning criterion of length <code>length(support.size)</code> .
<code>kpc.num</code>	The number of principal component being considered.
<code>var.pc</code>	The variance of principal components obtained by performing standard PCA.
<code>cum.var.pc</code>	Cumulative sums of <code>var.pc</code> .
<code>var.all</code>	If <code>sparse.type = "fpc"</code> , it is the total standard deviations of all principal components.
<code>pev</code>	A vector with the same length as <code>ev</code> . It records the percent of explained variance (compared to <code>var.all</code>) at each support size.
<code>pev.pc</code>	It records the percent of explained variance (compared to <code>var.pc</code>) at each support size.
<code>tune.type</code>	The criterion type for tuning parameters.
<code>tune.path</code>	The strategy for tuning parameters.
<code>call</code>	The original call to abess.

It is worthy to note that, if `sparse.type == "kpc"`, the `coef`, `support.size`, `ev`, `tune.value`, `pev` and `pev.pc` in list are list objects.

Note

Some parameters not described in the Details Section is explained in the document for [abess](#) because the meaning of these parameters are very similar.

Author(s)

Jin Zhu, Junxian Zhu, Ruihuang Liu, Junhao Huang, Xueqin Wang

References

A polynomial algorithm for best-subset selection problem. Junxian Zhu, Canhong Wen, Jin Zhu, Heping Zhang, Xueqin Wang. Proceedings of the National Academy of Sciences Dec 2020, 117 (52) 33117-33123; doi:10.1073/pnas.2014241117

Sparse principal component analysis. Hui Zou, Hastie Trevor, and Tibshirani Robert. Journal of computational and graphical statistics 15.2 (2006): 265-286. doi:10.1198/106186006X113430

See Also

[print.abesspca](#), [coef.abesspca](#), [plot.abesspca](#).

Examples

```
library(abess)
Sys.setenv("OMP_THREAD_LIMIT" = 2)

## predictor matrix input:
head(USArrests)
pca_fit <- abesspca(USArrests)
pca_fit
plot(pca_fit)

## covariance matrix input:
cov_mat <- stats::cov(USArrests) * (nrow(USArrests) - 1) / nrow(USArrests)
pca_fit <- abesspca(cov_mat, type = "gram")
pca_fit

## robust covariance matrix input:
rob_cov <- MASS::cov.rob(USArrests)[["cov"]]
rob_cov <- (rob_cov + t(rob_cov)) / 2
pca_fit <- abesspca(rob_cov, type = "gram")
pca_fit

## K-component principal component analysis
pca_fit <- abesspca(USArrests,
  sparse.type = "kpc",
  support.size = 1:4
)
coef(pca_fit)
plot(pca_fit)
plot(pca_fit, "coef")

## select support size via cross-validation ##
n <- 500
p <- 50
support_size <- 3
dataset <- generate.spc.matrix(n, p, support_size, snr = 20)
spca_fit <- abesspca(dataset[["x"]], tune.type = "cv", nfolds = 5)
plot(spca_fit, type = "tune")
```

abessrpca

Adaptive best subset selection for robust principal component analysis

Description

Decompose a matrix into the summation of low-rank matrix and sparse matrix via the best subset selection approach

Usage

```
abessrpca(
  x,
  rank,
  support.size = NULL,
  tune.path = c("sequence", "gsection"),
  gs.range = NULL,
  tune.type = c("gic", "aic", "bic", "ebic"),
  ic.scale = 1,
  lambda = 0,
  always.include = NULL,
  group.index = NULL,
  c.max = NULL,
  splicing.type = 2,
  max.splicing.iter = 1,
  warm.start = TRUE,
  important.search = NULL,
  max.newton.iter = 1,
  newton.thresh = 0.001,
  num.threads = 0,
  seed = 1,
  ...
)
```

Arguments

<code>x</code>	A matrix object.
<code>rank</code>	A positive integer value specify the rank of the low-rank matrix.
<code>support.size</code>	An integer vector representing the alternative support sizes. Only used for <code>tune.path = "sequence"</code> . Strongly suggest its minimum value larger than <code>min(dim(x))</code> .
<code>tune.path</code>	The method to be used to select the optimal support size. For <code>tune.path = "sequence"</code> , we solve the best subset selection problem for each size in <code>support.size</code> . For <code>tune.path = "gsection"</code> , we solve the best subset selection problem with support size ranged in <code>gs.range</code> , where the specific support size to be considered is determined by golden section.

<code>gs.range</code>	A integer vector with two elements. The first element is the minimum model size considered by golden-section, the later one is the maximum one. Default is <code>gs.range = c(1, min(n, round(n/(log(log(n))log(p)))))</code> .
<code>tune.type</code>	The type of criterion for choosing the support size. Available options are "gic", "ebic", "bic" and "aic". Default is "gic".
<code>ic.scale</code>	A non-negative value used for multiplying the penalty term in information criterion. Default: <code>ic.scale = 1</code> .
<code>lambda</code>	A single lambda value for regularized best subset selection. Default is 0.
<code>always.include</code>	An integer vector containing the indexes of variables that should always be included in the model.
<code>group.index</code>	A vector of integers indicating the which group each variable is in. For variables in the same group, they should be located in adjacent columns of x and their corresponding index in <code>group.index</code> should be the same. Denote the first group as 1, the second 2, etc. If you do not fit a model with a group structure, please set <code>group.index = NULL</code> (the default).
<code>c.max</code>	an integer splicing size. Default is: <code>c.max = 2</code> .
<code>splicing.type</code>	Optional type for splicing. If <code>splicing.type = 1</code> , the number of variables to be spliced is <code>c.max, ..., 1</code> ; if <code>splicing.type = 2</code> , the number of variables to be spliced is <code>c.max, c.max/2, ..., 1</code> . (Default: <code>splicing.type = 2</code> .)
<code>max.splicing.iter</code>	The maximum number of performing splicing algorithm. In most of the case, only a few times of splicing iteration can guarantee the convergence. Default is <code>max.splicing.iter = 20</code> .
<code>warm.start</code>	Whether to use the last solution as a warm start. Default is <code>warm.start = TRUE</code> .
<code>important.search</code>	An integer number indicating the number of important variables to be splicing. When <code>important.search</code> \ll p variables, it would greatly reduce runtimes. Default: <code>important.search = 128</code> .
<code>max.newton.iter</code>	a integer giving the maximal number of Newton's iteration iterations. Default is <code>max.newton.iter = 10</code> if <code>newton = "exact"</code> , and <code>max.newton.iter = 60</code> if <code>newton = "approx"</code> .
<code>newton.thresh</code>	a numeric value for controlling positive convergence tolerance. The Newton's iterations converge when $ dev - dev_{old} /(dev + 0.1) < newton.thresh$.
<code>num.threads</code>	An integer decide the number of threads to be concurrently used for cross-validation (i.e., <code>tune.type = "cv"</code>). If <code>num.threads = 0</code> , then all of available cores will be used. Default: <code>num.threads = 0</code> .
<code>seed</code>	Seed to be used to divide the sample into cross-validation folds. Default is <code>seed = 1</code> .
<code>...</code>	further arguments to be passed to or from methods.

Details

Adaptive best subset selection for robust principal component analysis aim to find two latent matrices L and S such that the original matrix X can be appropriately approximated:

$$x = L + S + N,$$

where L is a low-rank matrix, S is a sparse matrix, N is a dense noise matrix. Generic splicing technique can be employed to solve this problem by iteratively improve the quality of the estimation of S .

For a given support set Ω , the optimization problem:

$$\min_S \|x - L - S\|_F^2 \text{ s.t. } S_{ij} = 0 \text{ for } (i, j) \in \Omega^c,$$

still a non-convex optimization problem. We use the hard-impute algorithm proposed in one of the reference to solve this problem. The hard-impute algorithm is an iterative algorithm, people can set `max.newton.iter` and `newton.thresh` to control the solution precision of the optimization problem. (Here, the name of the two parameters are somehow abused to make the parameters cross functions have an unified name.) According to our experiments, we assign properly parameters to the two parameter as the default such that the precision and runtime are well balanced, we suggest users keep the default values unchanged.

Value

A S3 `abessrpca` class object, which is a list with the following components:

<code>S</code>	A list with <code>length(support.size)</code> elements, each of which is a sparse matrix estimation;
<code>L</code>	The low rank matrix estimation.
<code>nobs</code>	The number of sample used for training.
<code>nvars</code>	The number of variables used for training.
<code>rank</code>	The rank of matrix <code>L</code> .
<code>loss</code>	The loss of objective function.
<code>tune.value</code>	A value of tuning criterion of length <code>length(support.size)</code> .
<code>support.size</code>	The actual <code>support.size</code> values used. Note that it is not necessary the same as the input if the later have non-integer values or duplicated values.
<code>tune.type</code>	The criterion type for tuning parameters.
<code>call</code>	The original call to <code>abessrpca</code> .

Note

Some parameters not described in the Details Section is explained in the document for `abess` because the meaning of these parameters are very similar.

At present, l_2 regularization and group selection are not support, and thus, set `lambda` and `group.index` have no influence on the output. This feature will coming soon.

References

A polynomial algorithm for best-subset selection problem. Junxian Zhu, Canhong Wen, Jin Zhu, Heping Zhang, Xueqin Wang. Proceedings of the National Academy of Sciences Dec 2020, 117 (52) 33117-33123; doi:10.1073/pnas.2014241117

Emmanuel J. Candès, Xiaodong Li, Yi Ma, and John Wright. 2011. Robust principal component analysis? Journal of the ACM. 58, 3, Article 11 (May 2011), 37 pages. doi:10.1145/1970392.1970395

Mazumder, Rahul, Trevor Hastie, and Robert Tibshirani. Spectral regularization algorithms for learning large incomplete matrices. *The Journal of Machine Learning Research* 11 (2010): 2287-2322.

Examples

```
library(abess)
Sys.setenv("OMP_THREAD_LIMIT" = 2)
n <- 30
p <- 30
true_S_size <- 60
true_L_rank <- 2
dataset <- generate.matrix(n, p, support.size = true_S_size, rank = true_L_rank)
res <- abessrpca(dataset[["x"]], rank = true_L_rank, support.size = 50:70)
print(res)
coef(res)
plot(res, type = "tune")
plot(res, type = "loss")
plot(res, type = "S")
```

coef.abess

Extract Model Coefficients from a fitted "abess" object.

Description

This function provides estimated coefficients from a fitted "abess" object.

Usage

```
## S3 method for class 'abess'
coef(object, support.size = NULL, sparse = TRUE, ...)
```

Arguments

object	An "abess" project.
support.size	An integer vector specifies the coefficient fitted at given support.size. If support.size = NULL, then all coefficients would be returned. Default: support.size = NULL.
sparse	A logical value, specifying whether the coefficients should be presented as sparse matrix or not. Default: sparse = TRUE.
...	Other arguments.

Value

A coefficient matrix when fitting an univariate model including gaussian, binomial, poisson, and cox; otherwise, a list containing coefficient matrices. For a coefficient matrix, each row is a variable, and each column is a support size.

See Also

[print.abess](#), [predict.abess](#), [coef.abess](#), [extract.abess](#), [plot.abess](#), [deviance.abess](#).

coef.abesspca	<i>Extract Sparse Loadings from a fitted "abesspca" object.</i>
---------------	---

Description

This function provides estimated coefficients from a fitted "abesspca" object.

Usage

```
## S3 method for class 'abesspca'
coef(object, support.size = NULL, kpc = NULL, sparse = TRUE, ...)
```

Arguments

object	An "abesspca" project.
support.size	An integer vector specifies the coefficient fitted at given support.size. If support.size = NULL, then all coefficients would be returned. Default: support.size = NULL. This parameter is omitted if sparse.type = "kpc".
kpc	An integer vector specifies the coefficient fitted at given principal component. If kpc = NULL, then all coefficients would be returned. Default: kpc = NULL. This parameter is omitted if sparse.type = "fpc".
sparse	A logical value, specifying whether the coefficients should be presented as sparse matrix or not. Default: sparse = TRUE.
...	Other arguments.

Value

A matrix with length(support.size) columns. Each column corresponds to a sparse loading for the first principal component, where the number of non-zeros entries depends on the support.size.

See Also

[print.abesspca](#), [coef.abesspca](#), [plot.abesspca](#).

coef.abessrpca	<i>Extract sparse component from a fitted "abessrpca" object.</i>
----------------	---

Description

This function provides estimated coefficients from a fitted "abessrpca" object.

Usage

```
## S3 method for class 'abessrpca'
coef(object, support.size = NULL, sparse = TRUE, ...)
```

Arguments

object	An "abessrpca" project.
support.size	An integer vector specifies the sparse matrix fitted at given support.size to be returned. If support.size = NULL, then the sparse matrix with the least tuning value would be returned. Default: support.size = NULL.
sparse	A logical value, specifying whether the coefficients should be presented as sparse matrix or not. Default: sparse = TRUE.
...	Other arguments.

Value

A list with length(support.size) number of dgCMatrix, each of which is the estimation the sparse component.

deviance.abess	<i>Extract the deviance from a fitted "abess" object.</i>
----------------	---

Description

Similar to other deviance methods, which returns deviance from a fitted "abess" object.

Usage

```
## S3 method for class 'abess'
deviance(object, type = c("standard", "gic", "ebic", "bic", "aic"), ...)
```

Arguments

object	A "abess" object.
type	The type of deviance. One of the following: "standard", "gic", "ebic", "bic" and "aic". Default is "standard".
...	additional arguments

Value

A numeric vector.

See Also

[print.abess](#), [predict.abess](#), [coef.abess](#), [extract.abess](#), [plot.abess](#), [deviance.abess](#).

extract

Extract one model from a fitted "abess" object.

Description

Extract the fixed-support-size model's information such as the selected predictors, coefficient estimation, and so on.

Usage

```
extract(object, support.size = NULL, ...)
```

```
## S3 method for class 'abess'
extract(object, support.size = NULL, ...)
```

Arguments

object	An "abess" project.
support.size	An integer value specifies the model size fitted at given support.size. If support.size = NULL, then the model with the best tuning value would be returned. Default: support.size = NULL.
...	Other arguments.

Value

A list object including the following components:

beta	A p -by-1 matrix of sparse matrix, stored in column format.
intercept	The fitted intercept value.
support.size	The support.size used in the function.
support.beta	The support.size-length vector of fitted coefficients on the support set.
support.vars	The character vector gives variables in the support set.
tune.value	The tuning value of the model.
dev	The deviance of the model.

See Also

[print.abess](#), [predict.abess](#), [coef.abess](#), [extract.abess](#), [plot.abess](#), [deviance.abess](#).

generate.bmn.data *Generate a dataset from binary markov network*

Description

Generate a dataset with binary variables upon the interaction matrix. If the interaction matrix is not provided, and then, an interaction matrix would be generated.

Usage

```
generate.bmn.data(
  n,
  p,
  type = 1,
  seed = NULL,
  graph.seed = NULL,
  theta = NULL,
  beta = 0.7,
  degree = 3,
  alpha = 0.4,
  method = "freq"
)
```

Arguments

n	The number of observations.
p	The number of predictors of interest.
type	a numeric value that indicates the type of an interaction matrix
seed	Seed to be used to reproduce dataset
graph.seed	Seed used to reproduce the interaction matrix
theta	An user specified interaction matrix. It should be a symmetric matrix.
beta	a numeric value that specifies the remaining signal of the interaction matrix
degree	an integer value that specifies the maximum node degree in the interaction matrix
alpha	a numeric value that specifies the minimum signal of the interaction matrix
method	the method used for generating samples. One is the "freq", another is it "gibbs". The former is more suitable when $p < 18$, while the later is more suitable when $p > 18$.

Value

a list that includes three components: data, weight, theta. data is an n-by-p matrix that record samples; weight: an n-dimensional vector that records the frequency of each row of data; theta is the interaction matrix used for generating samples.

Examples

```
library(abess)
p <- 16
n <- 1e3
train <- generate.bmn.data(n, p, type = 7, graph.seed = 1, seed = 1, beta = 1.3, alpha = 0.2)
res <- slide(train[["data"]], train[["weight"]], tune.type = "gic")
```

generate.data

*Generate simulated data***Description**

Generate simulated data under the generalized linear model and Cox proportional hazard model.

Usage

```
generate.data(
  n,
  p,
  support.size = NULL,
  rho = 0,
  family = c("gaussian", "binomial", "poisson", "cox", "mgaussian", "multinomial",
    "gamma", "ordinal"),
  beta = NULL,
  cortype = 1,
  snr = 10,
  sigma = NULL,
  weibull.shape = 1,
  uniform.max = 1,
  y.dim = 3,
  class.num = 3,
  seed = 1
)
```

Arguments

n	The number of observations.
p	The number of predictors of interest.
support.size	The number of nonzero coefficients in the underlying regression model. Can be omitted if beta is supplied.
rho	A parameter used to characterize the pairwise correlation in predictors. Default is 0.
family	The distribution of the simulated response. "gaussian" for univariate quantitative response, "binomial" for binary classification response, "poisson" for counting response, "cox" for left-censored response, "mgaussian" for multivariate quantitative response, "mgaussian" for multi-classification response, "ordinal" for ordinal response.

beta	The coefficient values in the underlying regression model. If it is supplied, support.size would be omitted.
cortype	The correlation structure. cortype = 1 denotes the independence structure, where the covariance matrix has (i, j) entry equals $I(i \neq j)$. cortype = 2 denotes the exponential structure, where the covariance matrix has (i, j) entry equals $rho^{ i-j }$. cortype = 3 denotes the constant structure, where the non-diagonal entries of covariance matrix are rho and diagonal entries are 1.
snr	A numerical value controlling the signal-to-noise ratio (SNR). The SNR is defined as as the variance of $x\beta$ divided by the variance of a gaussian noise: $\frac{Var(x\beta)}{\sigma^2}$. The gaussian noise ϵ is set with mean 0 and variance. The noise is added to the linear predictor $\eta = x\beta$. Default is snr = 10. Note that this arguments's effect is overridden if sigma is supplied with a non-null value.
sigma	The variance of the gaussian noise. Default sigma = NULL implies it is determined by snr.
weibull.shape	The shape parameter of the Weibull distribution. It works only when family = "cox". Default: weibull.shape = 1.
uniform.max	A parameter controlling censored rate. A large value implies a small censored rate; otherwise, a large censored rate. It works only when family = "cox". Default is uniform.max = 1.
y.dim	Response's Dimension. It works only when family = "mgaussian". Default: y.dim = 3.
class.num	The number of class. It works only when family = "multinomial". Default: class.num = 3.
seed	random seed. Default: seed = 1.

Details

For family = "gaussian", the data model is

$$Y = X\beta + \epsilon.$$

The underlying regression coefficient β has uniform distribution $[m, 100m]$ and $m = 5\sqrt{2\log(p)/n}$.

For family= "binomial", the data model is

$$Prob(Y = 1) = \exp(X\beta + \epsilon)/(1 + \exp(X\beta + \epsilon)).$$

The underlying regression coefficient β has uniform distribution $[2m, 10m]$ and $m = 5\sqrt{2\log(p)/n}$.

For family = "poisson", the data is modeled to have an exponential distribution:

$$Y = Exp(\exp(X\beta + \epsilon)).$$

The underlying regression coefficient β has uniform distribution $[2m, 10m]$ and $m = \sqrt{2\log(p)/n}/3$.

For family = "gamma", the data is modeled to have a gamma distribution:

$$Y = Gamma(X\beta + \epsilon + 10, shape),$$

where shape is shape parameter in a gamma distribution. The underlying regression coefficient β has uniform distribution $[2m, 100m]$ and $m = \sqrt{2\log(p)/n}$.

For family = "ordinal", the data is modeled to have an ordinal distribution.

For family = "cox", the model for failure time T is

$$T = (-\log(U/\exp(X\beta)))^{1/weibull.shape},$$

where U is a uniform random variable with range $[0, 1]$. The centering time C is generated from uniform distribution $[0, uniform.max]$, then we define the censor status as $\delta = I(T \leq C)$ and observed time as $R = \min\{T, C\}$. The underlying regression coefficient β has uniform distribution $[2m, 10m]$, where $m = 5\sqrt{2\log(p)/n}$.

For family = "mgaussian", the data model is

$$Y = X\beta + E.$$

The non-zero values of regression matrix β are sampled from uniform distribution $[m, 100m]$ and $m = 5\sqrt{2\log(p)/n}$.

For family= "multinomial", the data model is

$$Prob(Y = 1) = \exp(X\beta + E)/(1 + \exp(X\beta + E)).$$

The non-zero values of regression coefficient β has uniform distribution $[2m, 10m]$ and $m = 5\sqrt{2\log(p)/n}$.

In the above models, $\epsilon \sim N(0, \sigma^2)$ and $E \sim MVN(0, \sigma^2 \times I_{q \times q})$, where σ^2 is determined by the snr and q is y . dim.

Value

A list object comprising:

x	Design matrix of predictors.
y	Response variable.
beta	The coefficients used in the underlying regression model.

Author(s)

Jin Zhu

Examples

```
# Generate simulated data
n <- 200
p <- 20
support.size <- 5
dataset <- generate.data(n, p, support.size)
str(dataset)
```

generate.matrix *Generate matrix composed of a sparse matrix and low-rank matrix*

Description

Generate simulated matrix that is the superposition of a low-rank component and a sparse component.

Usage

```
generate.matrix(  
  n,  
  p,  
  rank = NULL,  
  support.size = NULL,  
  beta = NULL,  
  snr = Inf,  
  sigma = NULL,  
  seed = 1  
)
```

Arguments

n	The number of observations.
p	The number of predictors of interest.
rank	The rank of low-rank matrix.
support.size	The number of nonzero coefficients in the underlying regression model. Can be omitted if beta is supplied.
beta	The coefficient values in the underlying regression model. If it is supplied, support.size would be omitted.
snr	A positive value controlling the signal-to-noise ratio (SNR). A larger SNR implies the identification of sparse matrix is much easier. Default snr = Inf enforces no noise exists.
sigma	A numerical value supplied the variance of the gaussian noise. Default sigma = NULL implies it is determined by snr.
seed	random seed. Default: seed = 1.

Details

The low rank matrix L is generated by $L = UV$, where U is an n -by- $rank$ matrix and V is a $rank$ -by- p matrix. Each element in U (or V) are i.i.d. drawn from $N(0, 1/n)$.

The sparse matrix S is an n -by- $rank$ matrix. It is generated by choosing a support set of size support.size uniformly at random. The non-zero entries in S are independent Bernoulli (-1, +1) entries.

The noise matrix N is an n -by- $rank$ matrix, the elements in N are i.i.d. gaussian random variable with standard deviation σ .

The SNR is defined as the variance of vectorized matrix $L + S$ divided by σ^2 .

The matrix x is the superposition of L, S, N :

$$x = L + S + N.$$

Value

A list object comprising:

x	An n -by- p matrix.
L	The latent low rank matrix.
S	The latent sparse matrix.

Author(s)

Jin Zhu

Examples

```
# Generate simulated data
n <- 30
p <- 20
dataset <- generate.matrix(n, p)

stats::heatmap(as.matrix(dataset[["S"]]),
  Rowv = NA,
  Colv = NA,
  scale = "none",
  col = grDevices::cm.colors(256),
  frame.plot = TRUE,
  margins = c(2.4, 2.4)
)
```

generate.spc.matrix *Generate matrix with sparse principal component*

Description

Generate simulated matrix that its principal component are sparse linear combination of its columns.

Usage

```
generate.spc.matrix(
  n,
  p,
  support.size = 3,
  snr = 20,
  sigma = NULL,
  sparse.loading = NULL,
  seed = 1
)
```

Arguments

<code>n</code>	The number of observations.
<code>p</code>	The number of predictors of interest.
<code>support.size</code>	A integer specify the number of non-zero entries in the first column of loading matrix.
<code>snr</code>	A positive value controlling the signal-to-noise ratio (SNR). A larger SNR implies the identification of sparse matrix is much easier. Default <code>snr = Inf</code> enforces no noise exists.
<code>sigma</code>	A numerical vector with length <code>p</code> specify the standard deviation of each columns. Default <code>sigma = NULL</code> implies it is determined by <code>snr</code> . If it is supplied, <code>support.size</code> would be omit.
<code>sparse.loading</code>	A <code>p</code> -by- <code>p</code> sparse orthogonal matrix. If it is supplied, <code>support.size</code> would be omit.
<code>seed</code>	random seed. Default: <code>seed = 1</code> .

Details

The methods for generating the matrix is detailedly described in the APPENDIX A: Data generation Section in Schipper et al (2021).

Value

A list object comprising:

<code>x</code>	An n -by- p matrix.
<code>coef</code>	The sparse loading matrix used to generate <code>x</code> .
<code>support.size</code>	A vector recording the number of non-zero entries in each .

References

Model selection techniques for sparse weight-based principal component analysis. de Schipper, Niek C and Van Deun, Katrijn. Journal of Chemometrics. 2021. doi:10.1002/cem.3289.

plot.abess	<i>Creat plot from a fitted "abess" object</i>
------------	--

Description

Produces a coefficient/deviance/tuning-value plot for a fitted "abess" object.

Usage

```
## S3 method for class 'abess'
plot(
  x,
  type = c("coef", "l2norm", "dev", "dev.ratio", "tune"),
  label = FALSE,
  ...
)
```

Arguments

x	A "abess" object.
type	The type of terms to be plot in the y-axis. One of the following: "coef" (i.e., coefficients), "l2norm" (i.e., L2-norm of coefficients), "dev" (i.e., deviance), and "tune" (i.e., tuning value). Default is "coef".
label	A logical value. If label = TRUE (the default), label the curves with variable sequence numbers.
...	Other graphical parameters to plot

Value

No return value, called for side effects.

Note

If family = "mgaussian", family = "ordinal" or family = "multinomial", a coefficient plot is produced for each dimension of multivariate response.

See Also

[print.abess](#), [predict.abess](#), [coef.abess](#), [extract.abess](#), [plot.abess](#), [deviance.abess](#).

Examples

```
dataset <- generate.data(100, 20, 3)
abess_fit <- abess(dataset[["x"]], dataset[["y"]])
plot(abess_fit)
plot(abess_fit, type = "l2norm")
plot(abess_fit, type = "dev")
plot(abess_fit, type = "tune")
```

plot.abesspca *Creat plot from a fitted "abess" object*

Description

Produces a coefficient/deviance/tuning-value plot for a fitted "abess" object.

Usage

```
## S3 method for class 'abesspca'  
plot(x, type = c("pev", "coef", "tune"), label = FALSE, ...)
```

Arguments

x	A "abess" object.
type	The type of terms to be plot in the y-axis. One of the following: "pev" (i.e., percent of explained variance), "coef" (i.e., coefficients), and "tune" (i.e., tuning value). Default is "coef".
label	A logical value. If label = TRUE (the default), label the curves with variable sequence numbers.
...	Other graphical parameters to plot

Value

No return value, called for side effects.

Note

If family = "mgaussian" or family = "multinomial", a coefficient plot is produced for each dimension of multivariate response.

See Also

[print.abesspca](#), [coef.abesspca](#), [plot.abesspca](#).

Examples

```
abess_fit <- abesspca(USArrests, support.size = 1:4, sparse.type = "kpc")  
plot(abess_fit)  
plot(abess_fit, type = "coef")  
plot(abess_fit, type = "tune")
```

plot.abessrpca *Creat plot from a fitted "abessrpca" object*

Description

Produces a sparse-matrix/loss/tuning-value plot for a fitted "abessrpca" object.

Usage

```
## S3 method for class 'abessrpca'
plot(x, type = c("S", "loss", "tune"), support.size = NULL, label = TRUE, ...)
```

Arguments

x	A "abessrpca" object.
type	The plot type. One of the following: "S" (i.e., a heatmap for the sparse matrix estimation), "loss" (i.e., a support.size versus loss plot), and "tune" (i.e., a support.size versus tuning value plot). Default is "coef".
support.size	An integer vector specifies the sparse matrix fitted at given support.size to be returned. If support.size = NULL, then the sparse matrix with the least tuning value would be returned. Default: support.size = NULL.
label	A logical value. If label = TRUE (the default), label the curves with variable sequence numbers.
...	Other graphical parameters to plot or stats::heatmap function

Value

No return value, called for side effects.

predict.abess *Make predictions from a fitted "abess" object.*

Description

Make predictions from a fitted "abess" object.

Usage

```
## S3 method for class 'abess'
predict(object, newx, type = c("link", "response"), support.size = NULL, ...)
```

Arguments

object	An "abess" project.
newx	New data used for prediction. If omitted, the fitted linear predictors are used.
type	type = "link" gives the linear predictors for "binomial", "poisson" or "cox" models; for "gaussian" models it gives the fitted values. type = "response" gives the fitted probabilities for "binomial" and "ordinal", fitted mean for "poisson" and the fitted relative-risk for "cox"; for "gaussian", type = "response" is equivalent to type = "link".
support.size	An integer value specifies the model size fitted at given support.size. If support.size = NULL, then the model with the best tuning value would be returned. Default: support.size = NULL.
...	Additional arguments affecting the predictions produced.

Value

The object returned depends on the types of family.

See Also

[print.abess](#), [predict.abess](#), [coef.abess](#), [extract.abess](#), [plot.abess](#), [deviance.abess](#).

print.abess	<i>Print method for a fitted "abess" object</i>
-------------	---

Description

Prints the fitted model and returns it invisibly.

Usage

```
## S3 method for class 'abess'
print(x, digits = max(5, getOption("digits") - 5), ...)
```

Arguments

x	A "abess" object.
digits	Minimum number of significant digits to be used.
...	additional print arguments

Details

Print a data.frame with three columns: the first column is support size of model; the second column is deviance of model; the last column is the tuning value of the certain tuning type.

Value

No return value, called for side effects

See Also

[print.abess](#), [predict.abess](#), [coef.abess](#), [extract.abess](#), [plot.abess](#), [deviance.abess](#).

print.abesspca	<i>Print method for a fitted "abesspca" object</i>
----------------	--

Description

Prints the fitted model and returns it invisibly.

Usage

```
## S3 method for class 'abesspca'
print(x, digits = max(5, getOption("digits") - 5), ...)
```

Arguments

x	A "abesspca" object.
digits	Minimum number of significant digits to be used.
...	additional print arguments

Details

Print a `data.frame` with three columns: the first column is support size of model; the second column is the explained variance of model; the last column is the percent of explained variance of model.

Value

No return value, called for side effects

See Also

[print.abesspca](#), [coef.abesspca](#), [plot.abesspca](#).

`print.abessrpca` *Print method for a fitted "abessrpca" object*

Description

Prints the fitted model and returns it invisibly.

Usage

```
## S3 method for class 'abessrpca'  
print(x, digits = max(5, getOption("digits") - 5), ...)
```

Arguments

<code>x</code>	A "abessrpca" object.
<code>digits</code>	Minimum number of significant digits to be used.
<code>...</code>	additional print arguments

Details

Print a `data.frame` with three columns: the first column is support size of model; the second column is the explained variance of model; the last column is the percent of explained variance of model.

Value

No return value, called for side effects

`slide` *Sparsity Learning for Ising moDel rEconstruction (SLIDE)*

Description

Sparsity Learning for Ising moDel rEconstruction (SLIDE)

Usage

```
slide(  
  x,  
  weight = NULL,  
  c.max = 8,  
  max.support.size = NULL,  
  tune.type = "cv",  
  foldid = NULL,  
  support.size = NULL,
```

```

    ic.scale = 1,
    graph.threshold = 0,
    newton = "approx"
  )

```

Arguments

<code>x</code>	Input matrix, of dimension $n \times p$; each row is an observation vector and each column is a predictor/feature/variable. Can be in sparse matrix format (inherit from class "dgCMatrix" in package Matrix).
<code>weight</code>	Observation weights. When <code>weight = NULL</code> , we set <code>weight = 1</code> for each observation as default.
<code>c.max</code>	an integer splicing size. Default is: <code>c.max = 2</code> .
<code>max.support.size</code>	The maximum node degree in the estimated graph. If prior information is available, we recommend setting this value accordingly. Otherwise, it is internally set to $n/(\log p \log \log n)$ by default.
<code>tune.type</code>	The type of criterion for choosing the support size. Available options are "gic", "ebic", "bic", "aic" and "cv". Default is "gic".
<code>foldid</code>	an optional integer vector of values between 1, ..., <code>n</code> folds identifying what fold each observation is in. The default <code>foldid = NULL</code> would generate a random <code>foldid</code> .
<code>support.size</code>	An integer vector representing the alternative support sizes. Only used for <code>tune.path = "sequence"</code> . Default is <code>0:min(n, round(n/(log(log(n))log(p))))</code> .
<code>ic.scale</code>	A non-negative value used for multiplying the penalty term in information criterion. Default: <code>ic.scale = 1</code> .
<code>graph.threshold</code>	A numeric value specifying the post-thresholding level for graph estimation. If prior knowledge about the minimum signal strength is available, this can be set to approximately half of that value. The default is <code>0.0</code> , which means no thresholding is applied.
<code>newton</code>	A character specify the Newton's method for fitting generalized linear models, it should be either <code>newton = "exact"</code> or <code>newton = "approx"</code> . If <code>newton = "exact"</code> , then the exact hessian is used, while <code>newton = "approx"</code> uses diagonal entry of the hessian, and can be faster (especially when <code>family = "cox"</code>).

Value

a sparse interaction matrix estimation

Examples

```

p <- 16
n <- 1e3
library(abess)
train <- generate.bmn.data(n, p, type = 3, graph.seed = 1, seed = 1, beta = 0.4)
res <- slide(train[["data"]], train[["weight"]], tune.type = "gic",

```

```
max.support.size = rep(4, p), support.size = rep(4, p))
all((res[[1]] != 0) == (train[["theta"]] != 0))
```

trim32

The Bardet-Biedl syndrome Gene expression data

Description

Gene expression data (500 gene probes for 120 samples) from the microarray experiments of mammalian eye tissue samples of Scheetz et al. (2006).

Format

A data frame with 120 rows and 501 variables, where the first variable is the expression level of TRIM32 gene, and the remaining 500 variables are 500 gene probes.

Details

In this study, laboratory rats (*Rattus norvegicus*) were studied to learn about gene expression and regulation in the mammalian eye. Inbred rat strains were crossed and tissue extracted from the eyes of 120 animals from the F2 generation. Microarrays were used to measure levels of RNA expression in the isolated eye tissues of each subject. Of the 31,000 different probes, 18,976 were detected at a sufficient level to be considered expressed in the mammalian eye. For the purposes of this analysis, we treat one of those genes, Trim32, as the outcome. Trim32 is known to be linked with a genetic disorder called Bardet-Biedl Syndrome (BBS): the mutation (P130S) in Trim32 gives rise to BBS.

Note

This data set contains 120 samples with 500 predictors. The 500 predictors are features with maximum marginal correlation to Trim32 gene.

References

T. Scheetz, k. Kim, R. Swiderski, A. Philp, T. Braun, K. Knudtson, A. Dorrance, G. DiBona, J. Huang, T. Casavant, V. Sheffield, E. Stone. Regulation of gene expression in the mammalian eye and its relevance to eye disease. Proceedings of the National Academy of Sciences of the United States of America, 2006.

Index

abess, [13](#), [17](#)
abess (abess.default), [3](#)
abess.default, [3](#)
abesspca, [10](#)
abessrpca, [15](#)

coef.abess, [8](#), [18](#), [19](#), [21](#), [29](#), [32](#), [33](#)
coef.abesspca, [14](#), [19](#), [19](#), [30](#), [33](#)
coef.abessrpca, [20](#)

deviance.abess, [8](#), [19](#), [20](#), [21](#), [29](#), [32](#), [33](#)

extract, [21](#)
extract.abess, [8](#), [19](#), [21](#), [29](#), [32](#), [33](#)

formula, [6](#)

generate.bmn.data, [22](#)
generate.data, [23](#)
generate.matrix, [26](#)
generate.spc.matrix, [27](#)

plot.abess, [8](#), [19](#), [21](#), [29](#), [29](#), [32](#), [33](#)
plot.abesspca, [14](#), [19](#), [30](#), [30](#), [33](#)
plot.abessrpca, [31](#)
predict.abess, [8](#), [19](#), [21](#), [29](#), [31](#), [32](#), [33](#)
print.abess, [8](#), [19](#), [21](#), [29](#), [32](#), [32](#), [33](#)
print.abesspca, [14](#), [19](#), [30](#), [33](#), [33](#)
print.abessrpca, [34](#)

slide, [34](#)

trim32, [36](#)