

# Best subset selection for transformation models

Lucas Kook

March 10, 2023

## Abstract

The **tramvs** package implements best subset selection for various kinds of transformation models via the **abess** algorithm. The optimal subset is elicited based on greedily updating the active set of covariates via changes in log-likelihood when in- or excluding a variable. This vignette illustrates the package’s functionalities, **S3** classes and -methods using simulated and real datasets.

## 1 Introduction

After introducing notation, the **abess** algorithm is described for linear transformation models. Extensions to more general transformation models are presented thereafter. Applications and simulation studies involving **tramvs** for linear location-scale transformation models can be found in (Siegfried et al., 2022).

### 1.1 Notation

Let  $Y$  denote a univariate response with at least ordered sample space  $\mathcal{Y}$ ,  $\mathbf{x} \in \mathbb{R}^p$  the observed covariates,  $F_Z$  an inverse link function, and  $h$  the transformation function. We model the conditional cumulative distribution function of  $Y|\mathbf{X} = \mathbf{x}$  using parametric linear transformation models (Hothorn et al., 2014, 2018)

$$F_{Y|\mathbf{X}=\mathbf{x}}(y) = F_Z(\mathbf{a}(y)^\top \boldsymbol{\vartheta} - \mathbf{x}^\top \boldsymbol{\beta}). \quad (1)$$

The parameters of the basis expansion  $h(y) := \mathbf{a}(y)^\top \boldsymbol{\vartheta}$  remain unpenalized and we consider only  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$  for penalization. We take a shorthand in writing  $\ell(\boldsymbol{\vartheta}, \boldsymbol{\beta}) := -\sum_{i=1}^n \ell_i(\boldsymbol{\vartheta}, \boldsymbol{\beta}; y_i, \mathbf{x}_i)$  for the negative log-likelihood of a transformation model assuming conditionally independent observations.

Let  $\mathcal{S} = [p]$  denote the set of all integers up to  $p$ , *i.e.*  $\{1, \dots, p\}$ . For any  $\mathcal{A} \subset \mathcal{S}$ ,  $\mathcal{A}^c = \mathcal{S} \setminus \mathcal{A}$  denotes the complement of  $\mathcal{A}$  and  $|\mathcal{A}|$  its cardinality. The support (or active set) of  $\boldsymbol{\beta}$  is denoted by  $\text{supp } \boldsymbol{\beta} = \{j : \beta_j \neq 0\}$ . By  $\boldsymbol{\beta}^{\mathcal{A}}$  we denote the restriction of  $\boldsymbol{\beta}$  to the support set  $\mathcal{A}$ , *i.e.*  $\beta_j^{\mathcal{A}} = 0$  if  $j \notin \mathcal{A}$ . The  $\ell_0$  norm can be written as  $\|\boldsymbol{\beta}\|_0 = |\text{supp } \boldsymbol{\beta}|$ , *i.e.* the number of non-zero entries in  $\boldsymbol{\beta}$ .

### 1.2 The abess algorithm

The **abess** algorithm (Zhu et al., 2020) performs best subset selection for a fixed support size  $s \in [p]$ ,

$$\min_{\boldsymbol{\vartheta}, \boldsymbol{\beta}} \ell(\boldsymbol{\vartheta}, \boldsymbol{\beta}), \quad \text{s.t. } \|\boldsymbol{\beta}\|_0 \leq s, \quad (2)$$

for a general class of models. It requires the computation of two “sacrifices”, namely a backward sacrifice  $\xi_j$  and a forward sacrifice  $\zeta_j$ . The backward sacrifice measures the drop in goodness of fit (as measured by the negative log-likelihood, *i.e.* smaller is better) when discarding variable  $j$  via

$$\xi_j := \ell(\hat{\boldsymbol{\beta}}^{\mathcal{A}}) - \ell(\hat{\boldsymbol{\beta}}^{\mathcal{A} \setminus \{j\}}). \quad (3)$$

The forward sacrifice measures the benefit of adding variable  $j$  via

$$\zeta_j := \ell(\hat{\boldsymbol{\beta}}^{\{j\}})|_{\hat{\boldsymbol{\beta}}^{\mathcal{A}}} - \ell(\hat{\boldsymbol{\beta}}^{\mathcal{A}}), \quad (4)$$

where  $\ell(\hat{\beta}^{\{j\}})|_{\hat{\beta}^{\mathcal{A}}}$  denotes the maximum likelihood when estimating  $\beta^{\{j\}}$  while keeping  $\hat{\beta}^{\mathcal{A}}$  fixed.

Based on both sacrifices, the **abess** algorithm looks for improvements of the active (and inactive) set for any splicing size  $k \leq s$  via

$$\mathcal{A}_k := \left\{ j \in \mathcal{A} : \sum_{i \in \mathcal{A}} \mathbf{1}(\xi_j \geq \xi_i) \leq k \right\}, \quad (5)$$

and

$$\mathcal{I}_k := \left\{ j \in \mathcal{I} : \sum_{i \in \mathcal{I}} \mathbf{1}(\zeta_j \leq \zeta_i) \leq k \right\}, \quad (6)$$

where  $\mathcal{I} := \mathcal{S} \setminus \mathcal{A}$  denotes the inactive set. Then, the active set is updated via

$$\tilde{\mathcal{A}} := (\mathcal{A} \setminus \mathcal{A}_k) \cup \mathcal{I}_k, \quad (7)$$

if there is an improvement in the negative log-likelihood (controlled via a tuning parameter  $\tau_s$ ). We choose  $\tau_s = 0.01s \log(p) \log(\log(n))/n$  per default. Further detail can be found in Zhu et al. (2020).

When the support size  $s$  is unknown, Zhu et al. (2020) recommend tuning  $s$  via a high-dimensional Bayesian information criterion (SIC), given by

$$\text{SIC}(\mathcal{A}) := \ell(\beta^{\mathcal{A}}) + \|\beta^{\mathcal{A}}\|_0 \log(p) \log \log n. \quad (8)$$

For varying support sizes  $s$ , the model with minimal SIC is selected. We illustrate this tuning in Section 2.4 and plot regularization and tuning paths.

**Choosing the initial support.** For choosing the initial  $\mathcal{A}$ , Zhu et al. (2020) recommend choosing those  $k$  covariates most correlated with the response  $Y$ . For transformation models this is problematic because empirical correlations are not well-defined for censored responses. Instead, the default in the **tramvs** package is to choose those  $k$  covariates most correlated with the score residuals of transformation model containing mandatory or no (*i.e.* an unconditional model) covariates. For a single observation  $(y, \mathbf{x})$ , the score residual is computed as,

$$s(\hat{\boldsymbol{\vartheta}}, \hat{\boldsymbol{\beta}}; y, \mathbf{x}) := \partial_{\alpha} \ell(\alpha; \boldsymbol{\vartheta}, \boldsymbol{\beta}, y, \mathbf{x}) \Big|_{\alpha=0, \boldsymbol{\vartheta}=\hat{\boldsymbol{\vartheta}}, \boldsymbol{\beta}=\hat{\boldsymbol{\beta}}}, \quad (9)$$

where  $\ell(\alpha; \boldsymbol{\vartheta}, \boldsymbol{\beta}, y, \mathbf{x})$  is the likelihood contribution of the observation in the model  $F_{Y|\mathbf{X}=\mathbf{x}}(y) = F_Z(\mathbf{a}(y)^{\top} \boldsymbol{\vartheta} - \mathbf{x}^{\top} \boldsymbol{\beta} - \alpha)$  (for more detail on score residuals, see Kook et al., 2021).

## 2 Illustrations

First, the basic usage of **tramvs** is explained. Then, we illustrate the package using various simulated and real datasets.

### 2.1 Basic usage

The function `abess_tram()` implements the core algorithm for best subset selection for a fixed support size  $s$ .

```
args(abess_tram)
```

```
## function (formula, data, modFUN, supp, mandatory = NULL, k_max = supp,
##   thresh = NULL, init = TRUE, m_max = 10, m0 = NULL, ...)
## NULL
```

`abess_tram()` is called by the main function `tramvs()`, which loops over the possible range of supports supplied to the function. It computes a high-dimensional information criterion (SIC) for model selection.

```
args(tramvs)

## function (formula, data, modFUN, mandatory = NULL, supp_max = NULL,
##      k_max = NULL, thresh = NULL, init = TRUE, m_max = 10, m0 = NULL,
##      ...)
## NULL
```

However, instead of supplying `modFUN` (a transformation model function), one can call the respective aliases, `<tram>VS()`, *e.g.* `CoxphVS()`, to skip this step. Further arguments to `modFUN` can be supplied via the `...` argument.

## 2.2 Interfacing tram

We generate a toy example with three out of ten non-zero coefficients in a normal linear regression model. We benchmark directly against the OLS alternative in the **abess** package (Zhu et al., 2020).

```
N <- 1e2; P <- 10; nz <- 3
beta <- rep(c(3, 0), c(nz, P - nz))
X <- matrix(rnorm(N * P), nrow = N, ncol = P)
Y <- 1 + X %*% beta + rnorm(N)

dat <- data.frame(y = Y, x = X)
cont_res <- tramvs(y ~ ., data = dat, modFUN = Lm)
```

```
res_abess <- abess(y ~ ., data = dat, family = "gaussian")
```

The two methods agree on the optimal subset, which can be easily extracted using `support()`.

```
support(cont_res)

## [1] "x.1" "x.2" "x.3"
```

```
extract(res_abess, support.size = res_abess$best.size)$support.vars

## [1] "x.1" "x.2" "x.3"
```

However, one can leave the Gaussian world behind by using a more flexible transformation function and a simple alias, like `BoxCoxVS()`.

```
BoxCoxVS(y ~ ., data = dat)
```

More low-level arguments to `BoxCox()` such as `order` or `extrapolate` can be supplied via the `...` argument, as shown below.

```
BoxCoxVS(y ~ ., data = dat, order = 3, extrapolate = TRUE)
```

## 2.3 Handling mandatory covariates

Mandatory covariates are covariates which should remain in the active set at all times. However, their coefficient estimates do not stay constant when other covariates are in- or excluded. In `tramvs`, mandatory covariates can be specified via a formula supplied to the `mandatory` argument.

```
BoxCoxVS(y ~ ., data = dat, mandatory = y ~ x.1)
```

Note that supplying mandatory covariates also alters the initialization of the active set for the `abess` algorithm. Now, instead of the residuals of the unconditional model, the residuals of `modFUN(mandatory, ...)` will be used.

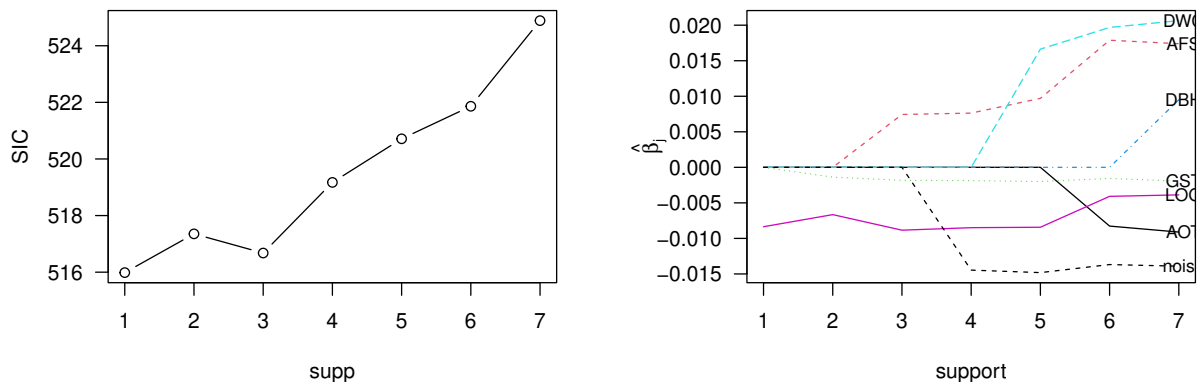


Figure 1: Demonstrating the plotting methods in a **cotram** example.

## 2.4 Interfacing cotram

Other transformation model add-on packages can be easily included in **tramvs**. The only requirements are a `logLik` method with arguments `newdata` and `parm`, and a `fixed` argument in `modFUN()`. For instance, best subset selection for models in the **cotram** add-on package (Siegfried and Hothorn, 2020) can be done as shown below.

```
library(cotram)

data("birds", package = "TH.data")
birds$noise <- rnorm(nrow(birds), sd = 10)

# Estimate support size via HBIC
count_res <- tramvs(SG5 ~ AOT + AFS + GST + DBH + DWC + LOG + noise, data = birds,
                    modFUN = cotram)
```

## 2.5 Location-scale transformation models

For linear location-scale transformation models (Siegfried et al., 2022),

$$F_Y(y|\mathbf{x}) = F_Z\left(\sqrt{\exp(\mathbf{x}^\top \boldsymbol{\gamma})} \mathbf{a}(y)^\top \boldsymbol{\vartheta} - \mathbf{x}^\top \boldsymbol{\beta}\right), \quad (10)$$

the initialization of the active set involves the model matrix for the location- and scale-terms and the correlation with the location- and scale-residuals, respectively. Thus, the residuals for the scale term are computed w.r.t. an a scale parameter constrained to one,

$$s(\hat{\boldsymbol{\vartheta}}, \hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\beta}}; y, \mathbf{x}) = \partial_\sigma \ell(\sigma; y, \mathbf{x}, \boldsymbol{\vartheta}, \boldsymbol{\gamma}, \boldsymbol{\beta}) \Big|_{\sigma=1, \boldsymbol{\vartheta}=\hat{\boldsymbol{\vartheta}}, \boldsymbol{\gamma}=\hat{\boldsymbol{\gamma}}, \boldsymbol{\beta}=\hat{\boldsymbol{\beta}}}, \quad (11)$$

for the transformation model (with  $\sigma > 0$ )

$$F_{Y|\mathbf{X}=\mathbf{x}}(y) = F_Z\left(\sigma \sqrt{\exp(\mathbf{x}^\top \boldsymbol{\gamma})} \mathbf{a}(y)^\top \boldsymbol{\beta} - \mathbf{x}^\top \boldsymbol{\beta}\right). \quad (12)$$

Location-scale transformation models can be specified via the formula interface of the usual **tram** functions by using a pipe on the right-hand side of the formula.

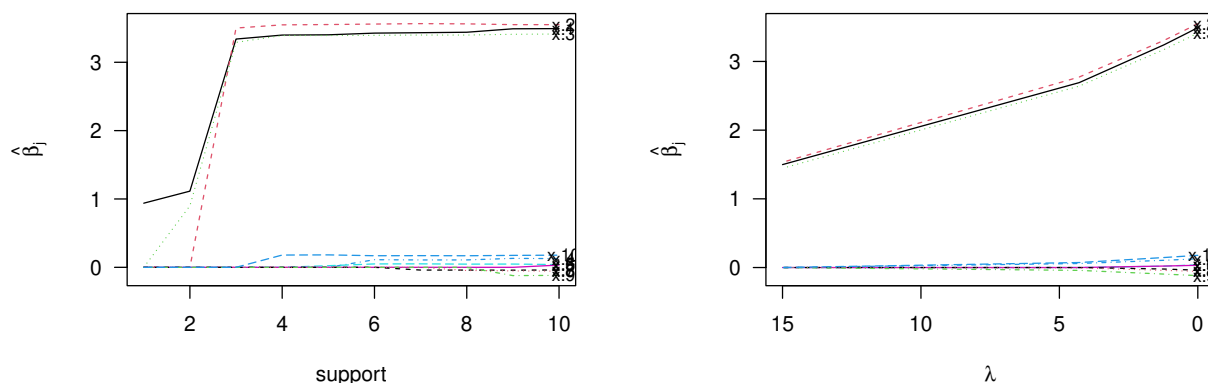


Figure 2: Comparing  $\ell_0$ - and  $\ell_1$ -regularized transformation models from **tramvs** and **tramnet**, respectively.

## 2.6 Comparison with tramnet

Transformation models with  $\ell_1$ - and  $\ell_2$ -penalties have been described previously and implemented in the **tramnet** package (Kook and Hothorn, 2021). The LASSO and elastic net penalties can be used for variable selection, but also shrink the regression coefficients. Especially the LASSO has trouble dealing with highly correlated covariates and tends to select a single random covariate from a group of highly correlated ones (Zou and Hastie, 2005). Figure 2 juxtaposes  $\ell_0$ - and  $\ell_1$ -regularization paths in the example with continuous response from Section 2.2.

```
m0 <- Lm(y ~ 1, data = dat)
X <- model.matrix(y ~ 0 + ., data = dat)
mt <- tramnet(m0, X, lambda = 0, alpha = 1)
pfl <- prof_lambda(mt, nprof = 5)
```

## 2.7 S3 methods

Below, all S3 methods for "**tramvs**" are showcased. Further arguments to the tram-specific methods can again be supplied via the ellipses. The plotting methods are illustrated in Fig. 1.

The **summary** method prints the full regularization path alongside a standard **summary.tram** of the best model.

```
# More elaborate summary
summary(cont_res)

##
## L0-penalized tram:
##
## Normal Linear Regression Model
##
## Call:
## modFUN(formula = formula, data = data, fixed = fix0, theta = theta_init[!names(theta_init) %in%
## I0])
##
## Coefficients:
## x.1 x.2 x.3 x.4 x.5 x.6 x.7 x.8 x.9 x.10
## 3.34 3.50 3.29 0.00 0.00 0.00 0.00 0.00 0.00 0.00
##
## Log-Likelihood:
```

```
## -134 (df = 5)
##
##
## SIC:
##   supp SIC
## 1     1 280
## 2     2 258
## 3     3 145
## 4     4 147
## 5     5 150
## 6     6 153
## 7     7 157
## 8     8 160
## 9     9 163
## 10    10 166
##
##
## Active set: x.1 x.2 x.3
```

The log-likelihood can be computed in and out-of-sample for the best model (`best_only=TRUE`) or all models. Additional arguments to `modFUN` can be supplied, *e.g.* `with_baseline=TRUE`.

```
# logLik of best model
logLik(cont_res)

## 'log Lik.' -134 (df=5)

nparm <- coef(cont_res, with_baseline = TRUE, best_only = TRUE)
logLik(cont_res, newdata = dat[1:5, ], parm = nparm)

## 'log Lik.' -5.05 (df=NULL)
```

The SIC that is printed in the summary can be obtained via `SIC()`.

```
# High-dimensional information criterion
SIC(cont_res)

##   supp SIC
## 1     1 280
## 2     2 258
## 3     3 145
## 4     4 147
## 5     5 150
## 6     6 153
## 7     7 157
## 8     8 160
## 9     9 163
## 10    10 166

SIC(cont_res, best_only = TRUE)

## [1] 145
```

Coefficients are returned as a sparse matrix for all model. In case of `best_only=TRUE`, a numeric vector is returned. Additional arguments to `coef.tram` can be supplied, as shown below.

```
coef(cont_res)

## 10 x 10 sparse Matrix of class "dgCMatrix"
##
```

```
## x.1 0.938 1.114 3.34 3.398 3.4017 3.4266 3.4326 3.4392 3.4898 3.4925
## x.2 . . 3.50 3.547 3.5519 3.5592 3.5653 3.5619 3.5505 3.5503
## x.3 . 0.899 3.29 3.391 3.3910 3.3948 3.3982 3.3981 3.4096 3.4134
## x.4 . . . . . 0.1112 0.1077 0.1080 0.1308 0.1272
## x.5 . . . . 0.0232 0.0496 0.0507 0.0464 0.0476 0.0401
## x.6 . . . . . . . . . 0.0327
## x.7 . . . . . . -0.0404 -0.0407 -0.0418 -0.0387
## x.8 . . . . . . . -0.0467 -0.0529 -0.0559
## x.9 . . . . . . . -0.1214 -0.1176
## x.10 . . . 0.179 0.1819 0.1694 0.1706 0.1691 0.1746 0.1776
```

```
coef(cont_res, best_only = TRUE)
```

```
## x.1 x.2 x.3 x.4 x.5 x.6 x.7 x.8 x.9 x.10
## 3.34 3.50 3.29 0.00 0.00 0.00 0.00 0.00 0.00 0.00
```

```
coef(cont_res, as.lm = TRUE)
```

```
## 11 x 10 sparse Matrix of class "dgCMatrix"
##
## (Intercept) 0.26 0.947 0.899 0.901 0.8981 0.892 0.8966 0.8893 0.8936 0.8889
## x.1 3.60 3.322 3.089 3.099 3.1020 3.108 3.1108 3.1136 3.1390 3.1398
## x.2 . . 3.238 3.236 3.2389 3.228 3.2311 3.2246 3.1935 3.1918
## x.3 . 2.680 3.044 3.093 3.0922 3.079 3.0797 3.0763 3.0668 3.0687
## x.4 . . . . . 0.101 0.0976 0.0978 0.1176 0.1143
## x.5 . . . . 0.0212 0.045 0.0459 0.0420 0.0428 0.0360
## x.6 . . . . . . . . . 0.0294
## x.7 . . . . . . -0.0366 -0.0368 -0.0376 -0.0348
## x.8 . . . . . . . -0.0423 -0.0476 -0.0502
## x.9 . . . . . . . . -0.1092 -0.1057
## x.10 . . . 0.164 0.1659 0.154 0.1546 0.1531 0.1571 0.1597
```

Several `tram` methods are applicable for the best model in an object of class `"tramvs"`, such as `predict`, `simulate`, and `residuals`.

```
head(predict(cont_res, which = "distribution", type = "trafo"))
simulate(cont_res)[1:5]
head(residuals(cont_res))
```

## References

- Torsten Hothorn, Thomas Kneib, and Peter Bühlmann. Conditional transformation models. *Journal of the Royal Statistical Society. Series B: Statistical Methodology*, 76(1):3–27, 2014. doi: 10.1111/rssb.12017.
- Torsten Hothorn, Lisa Möst, and Peter Bühlmann. Most Likely Transformations. *Scandinavian Journal of Statistics*, 45(1):110–134, 2018. doi: 10.1111/sjos.12291.
- Lucas Kook and Torsten Hothorn. Regularized Transformation Models: The tramnet Package. *The R Journal*, 13(1):581–594, 2021. doi: 10.32614/RJ-2021-054.
- Lucas Kook, Beate Sick, and Peter Bühlmann. Distributional Anchor Regression. *preprint arXiv:2101.08224*, 2021. URL <http://arxiv.org/abs/2101.08224>.
- Sandra Siegfried and Torsten Hothorn. Count transformation models. *Methods in Ecology and Evolution*, 11(7):818–827, 2020. doi: 10.1111/2041-210X.13383.
- Sandra Siegfried, Lucas Kook, and Torsten Hothorn. Distribution-free location-scale regression. *arXiv preprint arXiv:2208.05302*, 2022. doi: 10.48550/arXiv.2208.05302.

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

Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005. doi: 10.1111/j.1467-9868.2005.00503.x.

### 3 Session info

```
## R version 4.1.2 (2021-11-01)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Ubuntu 20.04.5 LTS
##
## Matrix products: default
## BLAS: /usr/lib/x86_64-linux-gnu/blas/libblas.so.3.9.0
## LAPACK: /usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.9.0
##
## locale:
## [1] LC_CTYPE=en_US.UTF-8 LC_NUMERIC=C
## [3] LC_TIME=de_CH.UTF-8 LC_COLLATE=en_US.UTF-8
## [5] LC_MONETARY=de_CH.UTF-8 LC_MESSAGES=en_US.UTF-8
## [7] LC_PAPER=de_CH.UTF-8 LC_NAME=C
## [9] LC_ADDRESS=C LC_TELEPHONE=C
## [11] LC_MEASUREMENT=de_CH.UTF-8 LC_IDENTIFICATION=C
##
## attached base packages:
## [1] stats graphics grDevices utils datasets methods base
##
## other attached packages:
## [1] cotram_0.3-1 colorspace_2.0-3 tramnet_0.0-6 mlrMBO_1.1.5
## [5] smooof_1.6.0.2 checkmate_2.0.0 mlr_2.19.0 ParamHelpers_1.14
## [9] CVXR_1.0-9 abess_0.3.0 tramvs_0.0-3 tram_0.7-2
## [13] mlt_1.4-3 basefun_1.1-3 variables_1.1-1 knitr_1.41
##
## loaded via a namespace (and not attached):
## [1] httr_1.4.4 tidyr_1.2.1 viridisLite_0.4.1
## [4] jsonlite_1.8.4 bit64_4.0.5 splines_4.1.2
## [7] ECOSolveR_0.5.4 Formula_1.2-5 highr_0.10
## [10] numDeriv_2016.8-1.1 pillar_1.8.1 backports_1.4.1
## [13] lattice_0.20-45 glue_1.6.2 quadprog_1.5-8
## [16] alabama_2022.4-1 digest_0.6.31 RColorBrewer_1.1-3
## [19] sandwich_3.0-2 htmltools_0.5.4 Matrix_1.5-3
## [22] pkgconfig_2.0.3 lhs_1.1.3 misc3d_0.9-1
## [25] purrr_1.0.1 mvtnorm_1.1-3 scales_1.2.1
## [28] parallelMap_1.5.1 mco_1.15.6 tibble_3.1.8
## [31] gmp_0.6-2 generics_0.1.3 ggplot2_3.4.0
## [34] TH.data_1.1-1 lazyeval_0.2.2 Rmpfr_0.8-6
## [37] cli_3.6.0 survival_3.2-13 RJSONIO_1.3-1.6
## [40] magrittr_2.0.3 evaluate_0.20 fansi_1.0.4
## [43] MASS_7.3-54 tools_4.1.2 data.table_1.14.6
## [46] lifecycle_1.0.3 BBmisc_1.11 multcomp_1.4-20
## [49] stringr_1.5.0 plotly_4.10.0 munsell_0.5.0
## [52] orthopolynom_1.0-6.1 compiler_4.1.2 rlang_1.0.6
## [55] plot3D_1.4 grid_4.1.2 coneproj_1.16
## [58] htmlwidgets_1.5.4 tcltk_4.1.2 gtable_0.3.1
## [61] codetools_0.2-18 BB_2019.10-1 polynom_1.4-1
```



```
## [64] R6_2.5.1          zoo_1.8-11          dplyr_1.1.0
## [67] fastmap_1.1.0       bit_4.0.5           utf8_1.2.3
## [70] fastmatch_1.1-3     stringi_1.7.12     parallel_4.1.2
## [73] Rcpp_1.0.10        vctrs_0.5.2        tidyselect_1.2.0
## [76] xfun_0.36
```