

# Regularized Transformation Models: The tramnet Package

by Lucas Kook and Torsten Hothorn

**Abstract** The **tramnet** package implements regularized linear transformation models by combining the flexible class of transformation models from **tram** with constrained convex optimization implemented in **CVXR**. Regularized transformation models unify many existing and novel regularized regression models under one theoretical and computational framework. Regularization strategies implemented for transformation models in **tramnet** include the LASSO, ridge regression and the elastic net and follow the parametrization in **glmnet**. Several functionalities for optimizing the hyperparameters, including model-based optimization based on the **mlrMBO** package, are implemented. A multitude of S3 methods are deployed for visualization, handling and simulation purposes. This work aims at illustrating all facets of **tramnet** in realistic settings and comparing regularized transformation models with existing implementations of similar models.

## Introduction

A plethora of R packages exist to estimate generalized linear regression models via penalized maximum likelihood, such as **penalized** (Goeman, 2010) and **glmnet** (Friedman et al., 2010). Both packages come with an extension to fit a penalized form of the Cox proportional hazard model. The **tramnet** package aims at unifying the above-mentioned and several novel models using the theoretical and computational framework of transformation models. Novel models in this class include Continuous Outcome Logistic Regression (COLR), as introduced by Lohse et al. (2017) and Box-Cox type regression models with a transformed conditionally normal response (Box and Cox, 1964; Hothorn, 2020c).

The disciplined convex optimization package **CVXR** (Fu et al., 2020) is applied to solve the constrained convex optimization problems that arise when fitting regularized transformation models. Transformation models are introduced in Section 2.1.1. For a more theoretical treatise, we refer to Hothorn et al. (2014, 2018); Hothorn (2020b). Convex optimization and domain-specific languages are briefly discussed in Section 2.1.3, followed by a treatment of model-based optimization for hyperparameter tuning (2.1.4).

## Transformation models

In stark contrast to penalized generalized linear models, regularized transformation models aim at estimating the response's whole conditional distribution instead of focusing on a single moment, e.g., the conditional mean. This conditional distribution function of a response  $Y$  is decomposed into an *a priori* chosen absolute continuous and log-concave error distribution  $F$  and a conditional transformation function  $h(y|x, s)$  that depends on the measured covariates  $x$  and stratum variables  $s$  and is monotone increasing in  $y$ . Although the model class is more flexible, packages **tram** and **tramnet** focus on stratified linear transformation models of the form

$$\mathbb{P}(Y \leq y | X = x, S = s) = F(h(y|x, s)) = F(h(y|s) - x^\top \beta). \quad (1)$$

Here, the baseline transformation is allowed to vary with stratum variables  $s$ , while covariate effects  $\beta$  are restricted to be shifts common to all baseline transformations  $h(y|s)$ .

In order for the model to represent a valid cumulative distribution function,  $F(h(y|x, s))$  has to be monotone increasing in  $y$  and thus in  $h$  for all possible strata  $s$  and all possible configurations of the covariates  $x$ . To ensure monotonicity,  $h$  is parametrized in terms of a basis expansion using Bernstein polynomials as implemented in the **basefun** package (Hothorn, 2020b). Hence,  $h$  is of the form

$$h(y) = a_{Bs,p}(y)^\top \boldsymbol{\vartheta},$$

where  $a_{Bs,p}(y)$  denotes the vector of basis functions in  $y$  of order  $p$  and  $\boldsymbol{\vartheta}$  are the coefficients for each basis function. Conveniently,  $a_{Bs,p}(y)^\top \boldsymbol{\vartheta}$  is monotone increasing in  $y$  as long as

$$\vartheta_i \leq \vartheta_{i+1} \quad \forall i = 0, \dots, p-1 \quad (2)$$

holds. For the concrete parameterization of stratified linear transformation models, the reader is referred to Hothorn (2020c).

Many contemporary models can be understood as linear transformation models, such as the normal linear regression model, logistic regression for binary, ordered and continuous responses, as well as exponential, Weibull and Rayleigh regression, and the Cox model in survival analysis. Thus, by appropriately choosing and parametrizing  $F$  and  $h$ , one can understand all those models in the same maximum likelihood-based framework. One can formulate the corresponding likelihood contributions not only for exact observations, but under any form of random censoring and truncation for continuous and discrete or ordered categorical responses.

Given a univariate response  $Y$  and a set of covariates  $\mathbf{X}$ , one can specify the following cumulative distribution function and density valid for any linear transformation model,

$$F_{Y|X=x}(y|\mathbf{s}, \mathbf{x}) = F\left(h(y|\mathbf{s}) - \mathbf{x}^\top \boldsymbol{\beta}\right),$$

$$f_{Y|X=x}(y|\mathbf{s}, \mathbf{x}) = F'\left(h(y|\mathbf{s}) - \mathbf{x}^\top \boldsymbol{\beta}\right) \cdot h'(y|\mathbf{s}).$$

From here, the log-likelihood contributions for exact, right, left, and interval-censored responses can be derived as

$$\ell_i(\boldsymbol{\theta}, \boldsymbol{\beta}; y_i, \mathbf{s}_i, \mathbf{x}_i) = \begin{cases} \log\left(F'\left(h(y_i|\mathbf{s}_i) - \mathbf{x}_i^\top \boldsymbol{\beta}\right)\right) + \log(h'(y_i|\mathbf{s}_i)) & y_i & \text{exact} \\ \log\left(F\left(h(\underline{y}|\mathbf{s}_i) - \mathbf{x}_i^\top \boldsymbol{\beta}\right)\right) & y_i \in (-\infty, \underline{y}] & \text{left} \\ \log\left(1 - F\left(h(\underline{y}|\mathbf{s}_i) - \mathbf{x}_i^\top \boldsymbol{\beta}\right)\right) & y_i \in (\underline{y}, \infty) & \text{right} \\ \log\left(F\left(h(\underline{y}|\mathbf{s}_i) - \mathbf{x}_i^\top \boldsymbol{\beta}\right) - F\left(h(\underline{y}|\mathbf{s}_i) - \mathbf{x}_i^\top \boldsymbol{\beta}\right)\right) & y_i \in (\underline{y}, \underline{y}] & \text{interval.} \end{cases}$$

The joint log-likelihood of several observations  $y_1, \dots, y_n$  is obtained by summing over the individual log-likelihood contributions  $\ell_i$  under the assumption that the individual samples are independent and identically distributed, the case exclusively dealt with by **tramnet**.

## Regularization

The aim of **tramnet** is to enable the estimation of regularized stratified linear transformation models. This is achieved by optimizing a penalized form of the log-likelihood introduced in the last section. The penalized log-likelihood,

$$\tilde{\ell}(\boldsymbol{\theta}, \boldsymbol{\beta}, \lambda, \alpha; y, \mathbf{s}, \mathbf{x}) = \ell(\boldsymbol{\theta}, \boldsymbol{\beta}; y, \mathbf{s}, \mathbf{x}) - \lambda \left( \alpha \|\boldsymbol{\beta}\|_1 + \frac{1}{2}(1 - \alpha) \|\boldsymbol{\beta}\|_2^2 \right),$$

consists of the unpenalized log-likelihood and an additional penalty term. Note that only the shift parameters  $\boldsymbol{\beta}$  are penalized, whereas the coefficients for the baseline transformation  $\boldsymbol{\theta}$  remain unpenalized. The parameterization of the penalty is chosen to be the same as in **glmnet**, consisting of a global penalization parameter  $\lambda$  and a mixing parameter  $\alpha$  controlling the amount of  $L_1$  compared to  $L_2$  penalization.

The two penalties and any combination thereof have unique properties and may be useful under different circumstances. A pure  $L_1$  penalty was first introduced by Tibshirani (1996) in an OLS framework and was dubbed the LASSO (Least Absolute Shrinkage and Selection Operator) due to its property of shrinking regression coefficients exactly to 0 for large enough  $\lambda$ . A pure LASSO penalty can be obtained in a regularized transformation model by specifying  $\alpha = 1$ . Applying an  $L_2$  penalty in an OLS problem was introduced more than five decades earlier by Tikhonov (1943) and later termed ridge regression (Hoerl and Kennard, 1970). In contrast to LASSO, ridge regression leads to shrunken regression coefficients, but does not perform automatic variable selection. Zou and Hastie (2005) picked up on both approaches, discussed their advantages, disadvantages, and overall characteristics and combined them into the elastic net penalty, a convex combination of an  $L_1$  and  $L_2$  penalty, controlled by the mixing parameter  $\alpha$ . Some of these properties will be illustrated for different models and a real-world data set in sections 2.1.6 and 2.2.2.

## Constrained convex optimization

Special algorithms were developed to optimize regularized objective functions, most prominently the LARS and LARS-EN algorithm (Efron et al., 2004) and variants thereof for the penalized Cox model (Goeman, 2010). However, the aim of **tramnet** is to solve the objective functions arising in regularized transformation models in a single computational framework. Due to the log-concavity of all choices for  $F$  in this package and  $h(y)$  being monotone increasing in  $y$ , the resulting log-likelihood contributions for any form of censoring and truncation are concave and can thus be solved by constrained convex optimization.

The fairly recent development of **CVXR** allows the specification of constrained convex optimization problems in terms of a domain-specific language, yielding an intuitive and highly flexible framework for constrained optimization. Because checking convexity of an arbitrarily complex expression is extremely hard, **CVXR** makes use of a library of smaller expressions, called atoms, with known monotonicity and curvature and tries to decompose the objective at hand using a set of rules from disciplined convex optimization (DCP) (Grant et al., 2006). Thus a complex expression's curvature can be more easily determined.

More formally, convex optimization aims at solving a problem of the form

$$\begin{aligned} & \underset{\boldsymbol{\theta}}{\text{minimize}} && g(\boldsymbol{\theta}) \\ & \text{subject to} && g_i(\boldsymbol{\theta}) \leq 0, \quad i = 1, \dots, K \\ & && A\boldsymbol{\theta} = \mathbf{b}, \end{aligned}$$

where  $\boldsymbol{\theta} \in \mathbb{R}^p$  is the parameter vector,  $g(\boldsymbol{\theta})$  is the objective function to be optimized,  $g_i(\boldsymbol{\theta})$  specify the inequality constraints and  $A \in \mathbb{R}^{n \times p}$  and  $\mathbf{b} \in \mathbb{R}^p$  parametrize any equality constraints on  $\boldsymbol{\theta}$ . Importantly, the objective function and all inequality constraint functions are convex (Boyd and Vandenberghe, 2004).

The likelihood  $\sum_i \ell_i(\boldsymbol{\theta}, \boldsymbol{\beta}; y_i, s_i, x_i)$  for transformation models of the form (1) are convex for error distributions with log-concave density, because log-convexity of  $-F'$  ensures the existence and uniqueness of the most likely transformation  $\hat{h}(y)$  and the convexity of  $-\ell(h; y, x)$ . Because the penalty term

$$\lambda \left( \alpha \|\boldsymbol{\beta}\|_1 + \frac{1}{2}(1 - \alpha) \|\boldsymbol{\beta}\|_2^2 \right)$$

is convex in  $\boldsymbol{\beta}$ , it can be added to the negative log-likelihood while conserving convexity. However, monotonicity of  $h$  imposes inequality constraints on the parameters of the baseline transformation as illustrated in equation (2). The elegance of domain-specific language-based optimizers comes to play when adding these and potential other inequality or equality constraints to the objective function, which will be showcased in Section 2.2.3. Thus, the optimisation routines implemented in package **CVXR** can be applied for computing maximum likelihood estimates of the parameters of model (1).

## Model-based optimization

The predictive capabilities of regularized regression models heavily depend on the hyperparameters  $\alpha$  and  $\lambda$ . Hyperparameter tuning can be addressed by a multitude of methods with varying computational complexity, advantages and disadvantages. Naïve or random grid search for more than one tuning parameter are computationally demanding, especially if the objective function is expensive to evaluate. Model-based optimization circumvents this issue by fitting a surrogate model, usually a Gaussian process, to the objective function. The objective function is evaluated at an initial, e.g., a random latin hypercube, design, to which the Gaussian process is subsequently fit. The surrogate model then proposes the next set of hyperparameters to evaluate the objective function at by some infill criterion (Horn and Bischl, 2016). Bischl et al. (2017) implement model-based optimization for multi-objective blackbox functions in the **mlrMBO** package. The objective function can, in theory, be vector-valued and the tuning parameter spaces may be categorical. In **tramnet**, the objective function is the cross-validated log-likelihood optimized using a Kriging surrogate model with expected improvement as the infill criterion. Model-based optimization for hyperparameter tuning is illustrated in section 2.2.

## Basic usage

The initial step is fitting a potentially stratified, transformation model of the form

```
R> m1 <- tram(y | s ~ 1, ...)
```

omitting all explanatory variables. This sets up the basis expansion for the transformation function, whose regression coefficients will not be penalized, as mentioned in section 2.1.2. Additionally, **tramnet()** needs a model matrix including the predictors, whose regression coefficients ought to be penalized. For numerical reasons, it is useful to provide a scaled model matrix instead of the original data, such that every parameter is equally affected by the regularization. Lastly, **tramnet()** will need the tuning parameters  $\alpha \in [0, 1]$  and  $\lambda \in \mathbb{R}^+$ , with  $\alpha$  representing a mixing parameter and  $\lambda$  controlling the extent of regularization. Setting  $\lambda = 0$  will result in an unpenalized model, regardless of the value of  $\alpha$ .

```
R> x <- model.matrix(~ 0 + x, ...)
```

Table caption is below the table

**Table 1:** Combinations of model classes and censoring types that are possible in the `tramnet` package. Due to missing disciplined geometric optimization rules in `CVXR`, not every combination of error distribution and censoring type yield solvable objective functions. This will change with coming updates in the `CVXR` package.

Model Class	Censoring Type			
	Exact	Left	Right	Interval
BoxCox	✓	✗	✗	✗
Colr	✓	✓	✓	✗
Coxph	✓	✗	✓	✗
Lehmann	✓	✓	✗	✗

```
R> x_scaled <- scale(x)
R> mt <- tramnet(model = m1, x = x_scaled, lambda, alpha, ...)

S3 methods accompanying the "tramnet" class will be discussed in section 2.3.
```

Censoring and likelihood forms

Specific combinations of  $F$  and the form of censoring yield log-log-concave log-likelihoods. Under these circumstances, `tramnet` is not yet able to solve the resulting optimization problem. Table 1 indicates which model class can be fitted under what type of censoring in the current version of `tramnet`.

Prostate cancer data analysis

The regularized normal linear and extensions to transformed normal regression models will be illustrated using the Prostate data set (Stamey et al., 1989), which was used by Zou and Hastie (2005) to highlight properties of the elastic net.

```
R> data("Prostate", package = "lasso2")
R> Prostate$psa <- exp(Prostate$lpsa)
R> Prostate[, 1:8] <- scale(Prostate[, 1:8])
```

The data set contains 97 observations and 9 covariates. In the original paper, the authors chose the log-transformed prostate specific antigen concentration (`lpsa`) as the response and used the eight remaining predictors: log cancer volume (`lcavol`), log prostate weight (`lweight`), age of the patient (`age`), log benign prostatic hyperplasia amount (`lbph`), seminal vesicle invasion (`svi` coded as 1 for yes, 0 for no), log capsular penetration (`lcp`), Gleason score (`gleason`), and percentage Gleason score 4 or 5 (`pgg45`) as covariates.

Linear and Box-Cox type regression models

Zou and Hastie (2005) imposed an assumption on the conditional distribution of the response by log-transforming and fitting a linear model. In the following, it is shown that the impact of this assumption may be assessed by estimating the baseline transformation from the data, followed by a comparison with the log-transformation applied by Zou and Hastie (2005). The linear models in `lpsa` and `log(psa)` are compared to transformation models with basis expansions in both `log(psa)` and `psa`, while specifying conditional normality of the transformed response. Additionally, the models are compared to an alternative implementation of regularized normal linear regression in **penalized**. Five different models will be used to illustrate important facets of transformation models, including parametrization and interpretation. The models are summarized in Table 2 and will be elaborated on throughout this section. The comparison is based on unpenalized models first. Later, the section highlights the penalized models together with hyperparameter tuning.

```
R> fm_Pr <- psa ~ lcavol + lweight + age + lbph + svi + lcp + gleason + pgg45
R> fm_Pr1 <- update(fm_Pr, ~ 0 + .)
R> x <- model.matrix(fm_Pr1, data = Prostate)
```

The normal linear regression model is implemented in `tram`'s `Lm()` function. `Lm()`'s parametrization differs from the usual linear model, hence caution has to be taken when interpreting the resulting

Table caption should be under the table.

**Table 2:** Summary of the five models illustrated in section 2.2, including their name throughout the manuscript, the R code to fit them, and the mathematical formulation of their conditional cumulative distribution function. For comparison, `mp` is included as an ordinary linear model, which is equivalent to model `mt` in terms of log-likelihood, but differs in the parametrization of the transformation function  $h$  and thus yields scaled coefficient estimates (cf. Table 3). Model `mtp` is a linear model parametrized in terms of a Bernstein basis of maximum order 1. This will yield the same coefficient estimates as `mt`, but a log-likelihood that is comparable to models `mt1` and `mt2`, whose transformation functions are parametrized in terms of higher-order Bernstein bases. The `log_first` argument specifies whether the basis expansion is calculated on the log-transformed or untransformed response.

Name	Code	Model for $F_{Y X=x}(y x)$
<code>mp</code>	<code>penalized(response = lpsa, penalized = x)</code>	$\Phi(\vartheta_1 + \vartheta_2 \log(y) - \mathbf{x}^\top \boldsymbol{\beta})$
<code>mt</code>	<code>Lm(lpsa ~ .)</code>	$\Phi(\vartheta_1 + \vartheta_2 \log(y) - \mathbf{x}^\top \boldsymbol{\beta})$
<code>mtp</code>	<code>BoxCox(psa ~ ., log_first = TRUE, order = 1)</code>	$\Phi(\mathbf{a}_{Bs,1}(\log(y))^\top \boldsymbol{\vartheta} - \mathbf{x}^\top \boldsymbol{\beta})$
<code>mt1</code>	<code>BoxCox(psa ~ ., log_first = TRUE, order = 7)</code>	$\Phi(\mathbf{a}_{Bs,7}(\log(y))^\top \boldsymbol{\vartheta} - \mathbf{x}^\top \boldsymbol{\beta})$
<code>mt2</code>	<code>BoxCox(psa ~ ., log_first = FALSE, order = 11)</code>	$\Phi(\mathbf{a}_{Bs,11}(y)^\top \boldsymbol{\vartheta} - \mathbf{x}^\top \boldsymbol{\beta})$

regression coefficients  $\boldsymbol{\beta}$ . In order to compare the results to an equivalent, already existing implementation, the same model is fitted using `penalized`.

```
R> m0 <- Lm(lpsa ~ 1, data = Prostate)
R> mt <- tramnet(m0, x = x, alpha = 0, lambda = 0)
R> mp <- penalized(response = Prostate$lpsa, penalized = x,
+                 lambda1 = 0, lambda2 = 0)
```

A linear model of the form

$$Y = \tilde{\alpha} + \mathbf{x}^\top \tilde{\boldsymbol{\beta}} + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

can be understood as a transformation model through reparametrization as

$$\mathbb{P}(Y \leq y | \mathbf{X} = \mathbf{x}) = \Phi(\vartheta_1 + \vartheta_2 y - \mathbf{x}^\top \boldsymbol{\beta}).$$

Here,  $\vartheta_1 = -\tilde{\alpha}/\sigma$  is a reparametrized intercept term,  $\vartheta_2 = 1/\sigma$  is the slope of the baseline transformation, and the regression coefficients  $\boldsymbol{\beta} = \tilde{\boldsymbol{\beta}}/\sigma$  represent scaled shift terms, influencing only the intercept. To recover the usual parametrization, `tramnet::coef.Lm()` offers the `as.lm = TRUE` argument.

```
R> cfx_tramnet <- coef(mt, as.lm = TRUE)
```

The transformation function for the linear model is depicted in Figure 1 (pink line). Because a linear baseline transformation imposes restrictive assumptions on the response's conditional distribution, it is advantageous to replace the linear baseline transformation by a more flexible one. In the case of the Box-Cox type regression model, the linear baseline transformation  $h(y) = \vartheta_1 + \vartheta_2 \log y$  is replaced by the basis expansion  $h(y) = \mathbf{a}_{Bs,7}(\log y)^\top \boldsymbol{\vartheta}$ .

```
R> ord <- 7 # flexible baseline transformation
R> m01 <- BoxCox(psa ~ 1, data = Prostate, order = ord,
+              extrapolate = TRUE, log_first = TRUE)
R> mt1 <- tramnet(m01, x = x, alpha = 0, lambda = 0)
```

The Box-Cox type regression model is then estimated with the `BoxCox()` function, while specifying the appropriate maximum order of the Bernstein polynomial. Because the more flexible transformation slightly deviates from being linear, the normal linear model yields a smaller log-likelihood (cf. Table 3). To make sure that this improvement is not due to the increased number of parameters and hence overfitting, the models' predictive capacities could be compared via cross-validation.

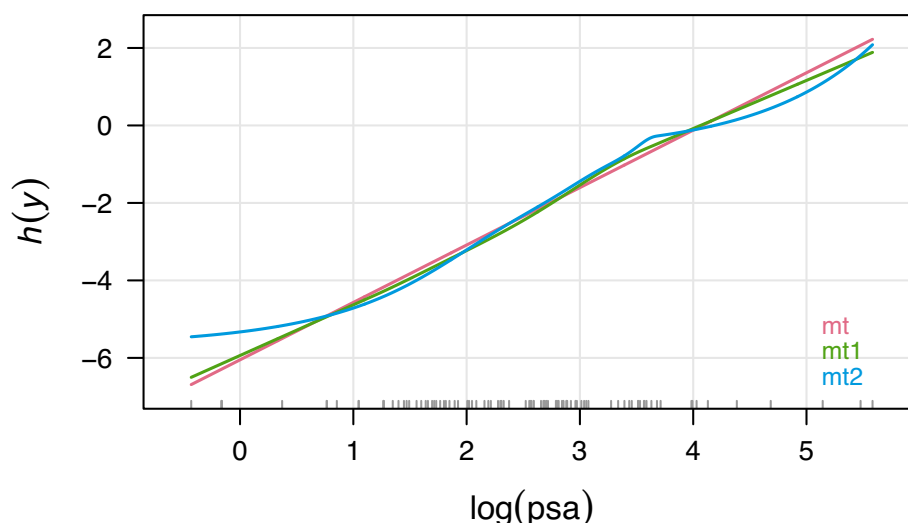
These results hold for the pre-specified log transformation of the response and a basis expansion thereof. Instead of prespecifying the log-transformation, its 'logarithmic nature' can be estimated from the data. Afterwards, one can compare the deviation from a log-linear baseline transformation graphically and by inspecting the predictive performance of the model in terms of the out-of-sample log-likelihood.

```
R> m02 <- BoxCox(psa ~ 1, order = 11, data = Prostate, extrapolate = TRUE)
R> mt2 <- tramnet(m02, x = x, lambda = 0, alpha = 0)
```

Indeed, the baseline transformation in Figure 1 is similar to the basis expansion in the log-transformed response upon visual inspection. Because  $mt$  is estimated using the log-transformed response and  $mt1$  and  $mt2$  are based on the original scale of the response, the resulting model log-likelihoods are not comparable. To overcome this issue, one can fit a Box-Cox type model with maximum order 1, as this results in a linear, but alternatively parametrized baseline transformation.

```
R> m0p <- BoxCox(psa ~ 1, order = 1, data = Prostate, log_first = TRUE)
R> mtp <- tramnet(m0p, x = x, lambda = 0, alpha = 0)
```

Figure 1 plots the three distinct baseline transformations resulting from models  $mt$ ,  $mt1$ , and  $mt2$ . The initial assumption to model the prostate-specific antigen concentration linearly on the log-scale seems to be valid when comparing the three transformation functions. The linear transformation in  $lpsa$  used in  $mt$ , and the basis expansion in  $\log(psa)$  ( $mt1$ ) are almost indistinguishable and yield very similar coefficient estimates, as well as log-likelihoods (*cf.* Table 3,  $mtp$  vs.  $mt1$ ). The basis expansion in  $psa$  ( $mt2$ ) is expected to be less stable due to the highly skewed untransformed response. This is reflected in Figure 1, where the baseline transformation deviates from being linear towards the bounds of the response's support. However, the log-linear behaviour of  $h$  was clearly captured by this model and further supports the initial assumption of conditional log-normality of the response. For the same reasons, the resulting log-likelihood of  $mt2$  is smaller than for  $mt1$  (Table 3). Taken together, this exemplary analysis highlights the flexibility and usefulness of transformation models to judge crucial modelling assumptions.



**Figure 1:** Comparison of different choices for the baseline transformation of the response (prostate-specific antigen concentration) in the Prostate data. The original analysis prespecified a log-transformation of the response and then assumed conditional normality on this scale. Hence the baseline transformation of  $mt$  is of the form:  $h(lpsa) = \theta_1 + \theta_2 \cdot lpsa$ . Now one can allow a more flexible transformation function in  $\log(psa)$  to judge any deviations of  $h(\log(psa))$  from linearity, leading to a baseline transformation in terms of basis functions:  $a_{Bs,7}(\log(psa))^T \theta$  in  $mt1$ . Lastly, instead of presuming a log-transformation one could estimate the baseline transformation from the raw response ( $psa$ ), i.e.  $h(psa) = a_{Bs,11}(psa)^T \theta$  in  $mt2$ . In this case, a higher order basis expansion was chosen to account for the skewness of the raw response. Notably, all three baseline transformations are fairly comparable. The basis expansion in  $psa$  deviates from being log-linear towards the boundaries of the response's support, as there are only a few observations.

## Hyperparameter tuning

This section features cross-validation, model-based optimization, and profiling functions for hyperparameter tuning, whose appropriate values are highly problem-dependent and hard to know in advance. **tramnet** implements naive grid search and model-based optimization in the functions `cv1_tramnet()` and `tramnet_mbo()`, respectively. In the framework of regularized transformation models, it is very natural to choose the out-of-sample log-likelihood as the objective function, because the notion of a mean square loss does not make sense for survival, let alone censored outcomes. The out-of-sample log-likelihood is, in fact, the log score, which is a proper scoring rule (Gneiting and Raftery, 2007).



Table Caption should be below the table.

**Table 3:** Comparison of the three transformation models on the Prostate data. Coefficient estimates are shown for each model, together with the in-sample log-likelihood in the last column. The first three models, mp, mt, and mtp use a linear baseline transformation in lpsa and log(psa), respectively. The mp model was fit using **penalized** and <sup>gave</sup> gives the scaled version of the regression coefficients in mt, but the same log-likelihood. At the same time, mt and mtp differ only in their response variable and its subsequent log-transformation in mtp, yielding the same coefficient estimates but a different log-likelihood. Models mt1 and mt2 allow a flexible basis expansion in log(psa) and psa, respectively. Model mt1, allowing for a flexible basis expansion in lpsa, fits the data the best. <sup>H</sup> However, the resulting coefficient estimates are similar for all models.

Model	Coefficient estimates								logLik
	lcavol	lweight	age	lbph	svi	lcp	gleason	pgg45	
mp	0.69	0.23	-0.15	0.16	0.32	-0.15	0.03	0.13	-99.5
mt	1.03	0.33	-0.22	0.23	0.47	-0.22	0.05	0.19	-99.5
mtp	1.03	0.33	-0.22	0.23	0.47	-0.22	0.05	0.19	-339.9
mt1	1.03	0.34	-0.21	0.22	0.48	-0.23	0.04	0.22	-338.0
mt2	0.97	0.32	-0.19	0.22	0.48	-0.21	0.07	0.21	-343.5

```
R> m0 <- BoxCox(lpsa ~ 1, data = Prostate, order = 7, extrapolate = TRUE)
R> mt <- tramnet(m01, x = x, alpha = 1, lambda = 0)
```

**tramnet** offers cross-validation in `cvl_tramnet()`, comparable to the `optL1()` and `optL2()` functions in **penalized**, which takes a sequence of values for  $\lambda$  and  $\alpha$  and performs a simple – and arguably slow – grid search. Per default, it computes 2-fold cross-validation, the user is encouraged, however, to judge the resulting bias-variance trade-off accordingly.

```
R> lambdas <- c(0, 10^seq(-4, log10(15), length.out = 4))
R> cvlt <- cvl_tramnet(object = mt, fold = 2, lambda = lambdas, alpha = 1)
```

In order to compare cross-validation across multiple packages and functions, it is also possible to supply the folds for each row in the design matrix as a numeric vector, as for example returned by `penalized::optL1()`.

```
R> pen_cvl <- optL1(response = lpsa, penalized = x, lambda2 = 0, data = Prostate,
+                 fold = 2)
R> cvlt <- cvl_tramnet(object = mt, lambda = lambdas, alpha = 1,
+                 folds = pen_cvl$fold)
```

The resulting object is of class “`cvl_tramnet`” and contains a table for the cross-validated log-likelihoods for each fold and the sum thereof, the ‘optimal’ tuning parameter constellation, which resulted in the largest cross-validated log-likelihood, tables for the cross-validated regularization paths, the folds and lastly the full fit based on the ‘optimal’ tuning parameters. Additionally, the resulting object can be used to visualize the log-likelihood and coefficient trajectories. These trajectories highly depend on the chosen folds, and the user is referred to the full profiling functions discussed in section 2.2.2.

Model-based optimization

In contrast to naive grid search, model-based optimization comprises more elegant methods for hyperparameter tuning. **tramnet** offers the `mbo_tramnet()` and `mbo_recommended()` functions. The former implements Kriging-based hyperparameter tuning for the elastic net, the LASSO and ridge regression. `mbo_tramnet()` takes a “`tramnet`” object as input and computes the cross-validated log-likelihood based on the provided fold or folds argument. The initial design is a random latin hypercube design with `n_design` rows per parameter. The number of sequential fits of the surrogate models is specified through `n_iter` and the range of the tuning parameters can be specified by `max/min` arguments. The default infill criterion is expected improvement. However, [Bischl et al. \(2017\)](#) encourage the use of the lower confidence bound over expected improvement, which can be achieved in `mbo_tramnet()` by specifying `opt_crit = makeMBOInfillCritCB()`. 10-fold cross-validation is used to compute the objective function for the initial design and each iteration. The recommended model is then extracted using `mbo_recommended()`.

```
R> tmbo <- mbo_tramnet(mt, obj_type = "elnet", fold = 10)
R> mtmbo <- mbo_recommended(tmbo, m0, x)
```

Unlike in the previous section, one can directly optimize the tuning parameters in an elastic net problem instead of optimizing over one hyperparameter at a time or having to specify LASSO or ridge regression *a priori*. The output of `mbo_tramnet()` is quite verbose and can be shortened by using the helper function `print_mbo()`.

```
R> print_mbo(tmbo)
```

```
## Recommended parameters:
## lmb=1.04e-05; alp=0.751
## Objective: y = 710
```

Interpreting the output, model-based optimization suggests an unpenalized model with  $\alpha = 0.75$  and  $\lambda = 0$ . This result stresses the advantages of model-based optimization over naive or random grid search in terms of complexity and computational efficiency. In the end, the proposed model is unpenalized and thus does not introduce sparsity in the regression coefficients.

```
R> coef(mtmb)
```

```
## lcavol lweight age lbph svi lcp gleason pgg45
## 1.0312 0.3380 -0.2068 0.2198 0.4801 -0.2329 0.0437 0.2157
```

```
R> summary(mtmb)$sparsity
## [1] "8 regression coefficients, 8 of which are non-zero"
```

## Regularization paths

As discussed before, it may be useful to inspect the full regularization paths over the tuning parameters  $\lambda$  and  $\alpha$ . Akin to the functions `profL1()` and `profL2()` in package **penalized**, **tramnet** offers `prof_lambda()` and `prof_alpha()`. Since these functions take a fitted model of class "tramnet" as input, which is internally updated, it is crucial to correctly specify the other tuning parameter in the model fitting step. In the example to come, `mt` was fit using  $\alpha = 1$  and  $\lambda = 0$ , resulting in a LASSO penalty only when profiling over  $\lambda$ . The resulting profile is depicted in Figure 2.

```
R> pfl <- prof_lambda(mt)
```

`prof_lambda()` takes `min_lambda`, `max_lambda` and `nprof` as arguments and internally generates an equi-spaced sequence from `min_lambda` to `max_lambda` on the log scale of length `nprof`. By default, this sequence ranges from 0 to 15 and is of length 5.

```
R> plot_path(pfl, plot_logLik = FALSE, las = 1, col = coll)
```

## Additional constraints

In some applications, the specification of additional constraints on the shift parameters  $\beta$  are of interest. Most commonly, either positivity or negativity for some or all regression coefficients is aimed at. In **tramnet**, additional inequality constraints can be specified via the `constraints` argument, which are internally handled as `constraints[[1]] %>= constraints[[2]]`. Hence, to specify the constraint of strictly positive regression coefficients, one would supply an identity matrix of dimension  $p$  for the left-hand side and the zero  $p$ -vector for the left-hand side, as done in the following example.

```
R> m0 <- BoxCox(lpsa ~ 1, data = Prostate, extrapolate = TRUE)
R> mt <- tramnet(m0, x, alpha = 0, lambda = 0, constraints = list(diag(8),
+                                                                rep(0, 8)))
R> coef(mt)
```

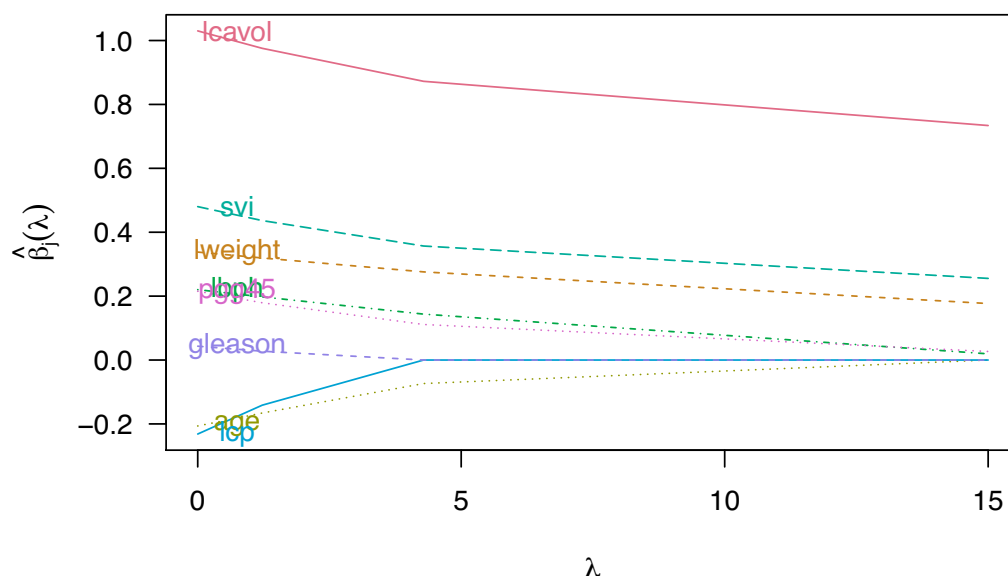
```
## lcavol lweight lbph svi gleason pgg45
## 0.9111 0.2996 0.1684 0.3969 0.0133 0.1125
```

The coefficients with a negative sign in the model without additional positivity constraints now shrink to zero and the other coefficient estimates change as well.

```
R> summary(mt)$sparsity
```

```
## [1] "8 regression coefficients, 6 of which are non-zero"
```





**Figure 2:** Full regularization paths for the tuning parameter  $\lambda$  using the default values of `plot_path()`.

One can compare this model to the implementation in **tram**, where it is also possible to specify linear inequality constraints on the regression coefficients  $\beta$ . Here, it is sufficient to specify `constraints = c("age >= 0", "lcp >= 0")` for the two non-positive coefficient estimates, due to the convexity of the underlying problem.

```
R> m <- BoxCox(lpsa ~ . - psa, data = Prostate, extrapolate = TRUE,
+             constraints = c("age >= 0", "lcp >= 0"))
R> max(abs(coef(m) - coef(mt, tol = 0)))
```

```
## [1] 1.28e-05
```

Indeed, both optimizers arrive at virtually the same coefficient estimates.

### S3 Methods

Building on the S3 infrastructure of the packages **mlt** and **tram**, this package provides corresponding methods for the following generics: `coef()`, `logLik()`, `plot()`, `predict()`, `simulate()` and `residuals()`. The methods' additional "tramnet"-specific arguments will be briefly discussed in this section.

`coef.tramnet()` suppresses the baseline transformation's coefficient estimates  $\hat{\theta}$  by default and considers shift parameter estimates  $\hat{\beta}$  below  $10^{-6}$  as 0 to stress the selected variables only. This threshold can be controlled by the `tol` argument. Hence, `coef(mt, with_baseline = TRUE, tol = 0)` returns all coefficients.

```
R> coef(mtmbo, with_baseline = TRUE, tol = 0)
```

```
## Bs1(lpsa) Bs2(lpsa) Bs3(lpsa) Bs4(lpsa) Bs5(lpsa) Bs6(lpsa) Bs7(lpsa)
## -1.9775 -1.5055 -1.0335 -0.2778 -0.2778 1.0723 1.5150
## Bs8(lpsa) lcavol lweight age lbph svi lcp
## 1.9576 1.0312 0.3380 -0.2068 0.2198 0.4801 -0.2329
## gleason pgg45
## 0.0437 0.2157
```

The `logLik.tramnet()` method allows the log-likelihoods re-computation under new data (i.e. out-of-sample) and different coefficients (`parm`) and weights (`w`), as illustrated below.

```
R> logLik(mtmbo)
```

```
## 'log Lik.' -97.7 (df=NA)
```

```
R> cfx <- coef(mtmbo, with_baseline = TRUE, tol = 0)
R> cfx[5:8] <- 0.5
R> logLik(mtmbo, parm = cfx)
```

```
## 'log Lik.' -561 (df=NA)
```

```
R> logLik(mtmbo, newdata = Prostate[1:10,])
```

```
## 'log Lik.' -14.3 (df=NA)
```

```
R> logLik(mtmbo, w = runif(n = nrow(mtmbo$x)))
```

```
## 'log Lik.' -41.8 (df=NA)
```

In the spirit of `mlt`'s plotting methods for classes `"mlt"` and `"ctm"`, `plot.tramnet()` offers diverse plotting options for objects of class `"tramnet"`. The specification of new data and the type of plot is illustrated in the following code chunk and Figure 3.

```
R> par(mfrow = c(3, 2)); K <- 1e3
R> plot(mtmbo, type = "distribution", K = K, main = "A") # A, default
R> plot(mtmbo, type = "survivor", K = K, main = "B") # B
R> plot(mtmbo, type = "trafo", K = K, main = "C") # C
R> plot(mtmbo, type = "density", K = K, main = "D") # D
R> plot(mtmbo, type = "hazard", K = K, main = "E") # E
R> plot(mtmbo, type = "trafo", newdata = Prostate[1, ], col = 1, K = K, main = "F") # F
```

The `predict.tramnet()` method works in the same way as `predict.mlt()` and as such supports the types `trafo`, `distribution`, `survivor`, `density`, `logdensity`, `hazard`, `loghazard`, `cumhazard` and `quantile`. For type = `"quantile"`, the corresponding probabilities (`prob`) have to be supplied as an argument to evaluate the quantile function at.

```
R> predict(mtmbo, type = "quantile", prob = 0.2, newdata = Prostate[1:5,])
```

```
## prob      [,1] [,2] [,3] [,4] [,5]
##      0.2   3.4  3.55  3.74  3.72  2.68
```

Another method offered by this package implements parametric bootstrap-based sampling. In particular, `simulate.tramnet()` calls the `simulate.ctm()` function after converting the `"tramnet"` object to a `"ctm"` object.

```
R> simulate(mtmbo, nsim = 1, newdata = Prostate[1:5,], seed = 1)
```

```
## [1] 3.56 3.97 4.57 5.48 2.69
```

Lastly, `residuals.tramnet()` computes the generalized residual  $r$  defined as the score contribution for sample  $i$  with respect to a newly introduced intercept parameter  $\gamma$ , which is restricted to be zero. In particular,

$$r = \left. \frac{\partial \ell(\boldsymbol{\vartheta}, \boldsymbol{\beta}, \gamma; y, \mathbf{s}, \mathbf{x})}{\partial \gamma} \right|_{\gamma=0}$$

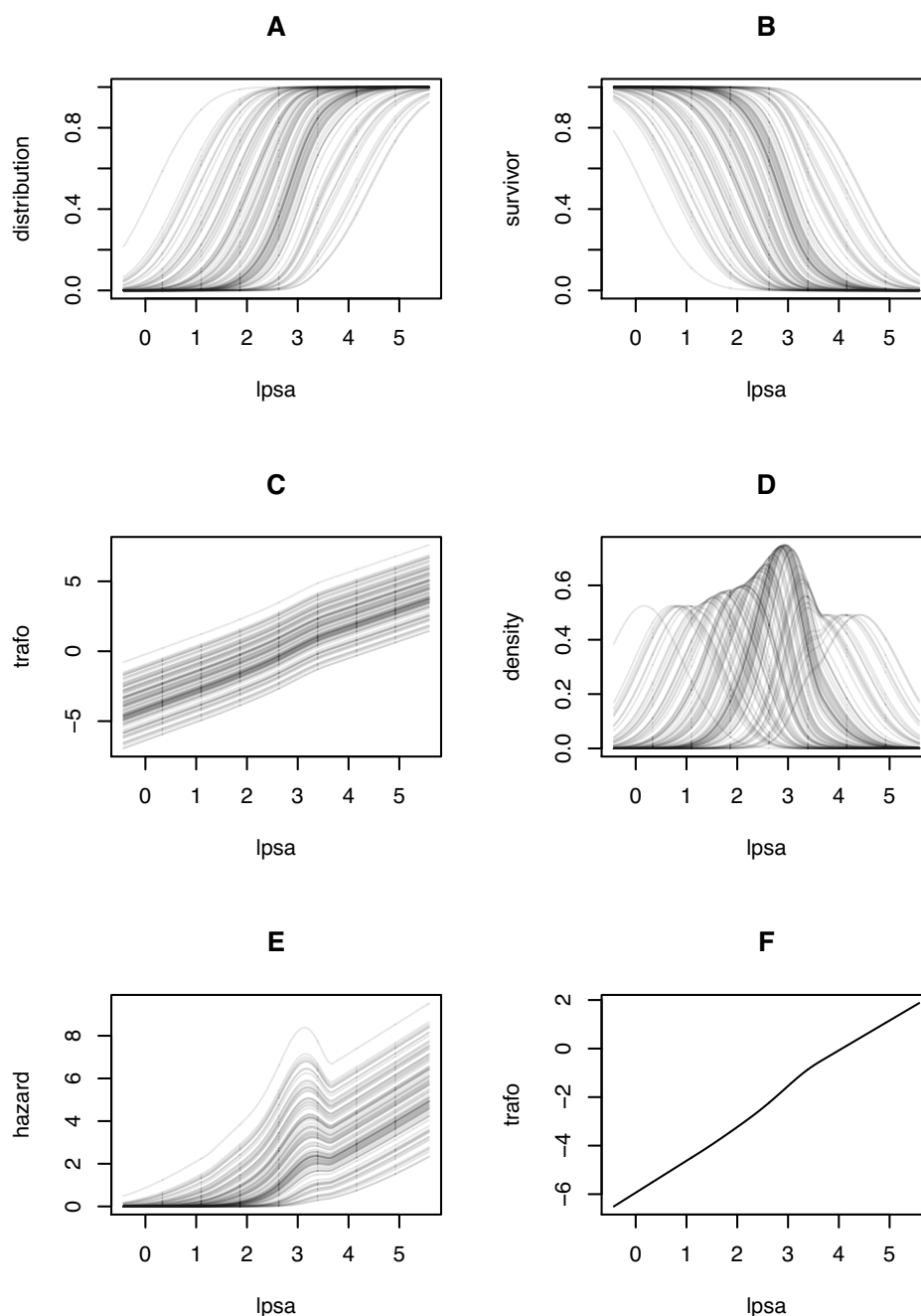
yields the generalized residual with respect to  $\gamma$  for the model

$$F_Y(y|\mathbf{s}, \mathbf{x}) = F\left(h(y|\mathbf{s}) - \mathbf{x}^\top \boldsymbol{\beta} - \gamma\right).$$

```
R> residuals(mtmbo)[1:5]
```

```
## [1] -6.50 -6.36 -6.60 -6.57 -4.17
```

In residual analysis and boosting it is common practice to check for associations between residuals and covariates that are not included in the model. In the prostate cancer example one could investigate whether the variables `age` and `lcp` should be included in the model. To illustrate this particular case, a non-parametric `independence_test()` from package `coin` can be used (Hothorn et al., 2008). First, the unconditional transformation model `m0` is fit. Afterwards, the `tramnet` models excluding `age` and `lcp` are estimated, and their residuals extracted using the `residuals.tramnet()` method. Lastly, an independence test using a maximum statistic (`teststat = "max"`) and a Monte Carlo-based approximation



**Figure 3:** Illustration of `plot.tramnet()`'s versatility in visualizing the response's estimated conditional distribution on various scales including cdf, survivor, transformation scale, and pdf. Note that by default the plot is produced for each row in the design matrix. In unstratified linear transformation models, this leads to shifted versions of the same curve on the transformation function's scale. A: Estimated conditional distribution function for every observation. B: Estimated conditional survivor function for every observation. The conditional survivor function is defined as  $S(y|x) = 1 - F_Y(y|x)$ . C: Conditional most likely transformation for every observation. Note that every conditional transformation function is a shifted version of the same curve. D: The conditional density for every observation can be calculated using  $f_Y(y|x) = F'(a(y)^\top \boldsymbol{\theta} - x^\top \boldsymbol{\beta}) a'(y)^\top \boldsymbol{\theta}$ . E: A distribution function is fully characterized by its hazard function  $\lambda(y|x) = f_Y(y|x)/S(y|x)$ , which is depicted in this panel. F: The `newdata` argument can be used to plot the predicted most likely transformation for the provided data, in this case the first row of the Prostate data.

of the null distribution based on resampling  $10^6$  times (`distribution = approximate(1e6)`), yields the results printed below.

```
R> library("coin")
R> m0 <- BoxCox(lpsa ~ 1, data = Prostate, extrapolate = TRUE)
R> x_no_age_lcp <- x[, !colnames(x) %in% c("age", "lcp")]
R> mt_no_age_lcp <- tramnet(m0, x_no_age_lcp, alpha = 0, lambda = 0)
R> r <- residuals(mt_no_age_lcp)
R> it <- independence_test(r ~ age + lcp, data = Prostate,
+                           teststat = "max", distribution = approximate(1e6))
R> pvalue(it, "single-step")

## age  0.023748
## lcp <0.000001
```

Because there is substantial evidence against <sup>the</sup>independence of the models' residuals to either lcp or age, we can conclude that it is worthwhile to include age and lcp in the model. Packages **trtf** (Hothorn, 2019b) and **tbm** (Hothorn, 2020a, 2019a) make use of this definition of a residual for estimating and boosting transformation models, trees, and random forests. For more theoretical insight, the reader is referred to the above mentioned publications.

## Bibliography

- B. Bischl, J. Richter, J. Bossek, D. Horn, J. Thomas, and M. Lang. *mlrMBO: A Modular Framework for Model-Based Optimization of Expensive Black-Box Functions*, 2017. URL <http://arxiv.org/abs/1703.03373>. [p]
- G. E. Box and D. R. Cox. An analysis of transformations. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 26(2):211–243, 1964. doi: 10.1111/j.2517-6161.1964.tb00553.x. [p]
- S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004. doi: 10.1017/CBO9780511804441. [p]
- B. Efron, T. Hastie, I. Johnstone, R. Tibshirani, et al. Least angle regression. *The Annals of Statistics*, 32(2):407–499, 2004. doi: 10.1214/009053604000000067. [p]
- J. Friedman, T. Hastie, and R. Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33(1):1–22, 2010. doi: 10.18637/jss.v033.i01. [p]
- A. Fu, B. Narasimhan, and S. Boyd. CVXR: An R package for disciplined convex optimization. *Journal of Statistical Software*, 2020. doi: <https://arxiv.org/abs/1711.07582>. Accepted for publication. [p]
- T. Gneiting and A. E. Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007. doi: 10.1198/016214506000001437. [p]
- J. J. Goeman. L1 penalized estimation in the Cox proportional hazards model. *Biometrical Journal*, 52(1):–14, 2010. doi: 10.1002/bimj.200900028. [p]
- M. Grant, S. Boyd, and Y. Ye. Disciplined convex programming. In *Global Optimization*, pages 155–210. Springer, 2006. doi: 10.1007/s11590-019-01422-z. [p]
- A. E. Hoerl and R. W. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970. doi: 10.1080/00401706.1970.10488634. [p]
- D. Horn and B. Bischl. Multi-objective parameter configuration of machine learning algorithms using model-based optimization. In *2016 IEEE Symposium Series on Computational Intelligence (SSCI)*, pages 1–8. IEEE, 2016. doi: 10.1109/SSCI.2016.7850221. [p]
- T. Hothorn. *tbm: Transformation Boosting Machines*, 2019a. URL <https://CRAN.R-project.org/package=tbm>. R package version 0.3-0. [p]
- T. Hothorn. *trtf: Transformation Trees and Forests*, 2019b. URL <https://CRAN.R-project.org/package=trtf>. R package version 0.3-5. [p]
- T. Hothorn. Transformation boosting machines. *Statistics and Computing*, 30:141–152, 2020a. doi: 10.1007/s11222-019-09870-4. [p]
- T. Hothorn. Most likely transformations: The mlt package. *Journal of Statistical Software*, 92(1):1–68, 2020b. doi: 10.18637/jss.v092.i01. [p]
- T. Hothorn. *tram: Transformation Models*, 2020c. URL <http://ctm.R-forge.R-project.org>. R package version 0.3-2. [p]
- T. Hothorn, K. Hornik, M. A. van de Wiel, and A. Zeileis. Implementing a class of permutation tests: The coin package. *Journal of Statistical Software*, 28(8):1–23, 2008. doi: 10.18637/jss.v028.i08. [p]
- T. Hothorn, T. Kneib, and P. Bühlmann. Conditional transformation models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):3–27, 2014. doi: 10.1111/rssb.12017. [p]
- T. Hothorn, L. Möst, and P. Bühlmann. Most likely transformations. *Scandinavian Journal of Statistics*, 45(1):110–134, 2018. doi: 10.1111/sjos.12291. [p]
- T. Lohse, S. Rohrmann, D. Faeh, and T. Hothorn. Continuous outcome logistic regression for analyzing body mass index distributions. *F1000Research*, 6(1933), 2017. doi: 10.12688/f1000research.12934.1. [p]
- T. A. Stamey, J. N. Kabalin, J. E. McNeal, I. M. Johnstone, F. Freiha, E. A. Redwine, and N. Yang. Prostate specific antigen in the diagnosis and treatment of adenocarcinoma of the prostate. ii. radical prostatectomy treated patients. *The Journal of Urology*, 141(5):1076–1083, 1989. doi: 10.1016/S0022-5347(17)41176-1. [p]

- R. Tibshirani. Regression shrinkage and selection via the LASSO. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 58(1):267–288, 1996. doi: 10.1111/j.2517-6161.1996.tb02080.x. [p]
- A. N. Tikhonov. On the stability of inverse problems. In *Dokl. Akad. Nauk SSSR*, volume 39, pages 195–198, 1943. doi: 10.1155/2011/450269. [p]
- H. Zou and T. 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. [p]

Lucas Kook, Torsten Hothorn  
Institut für Epidemiologie, Biostatistik und Prävention  
Universität Zürich  
Hirschengraben 84, CH-8001 Zürich  
Switzerland  
[lucasheinrich.kook@uzh.ch](mailto:lucasheinrich.kook@uzh.ch), [Torsten.Hothorn@R-project.org](mailto:Torsten.Hothorn@R-project.org)