



Transformation boosting machines

Torsten Hothorn¹

Received: 26 November 2018 / Accepted: 27 March 2019
© The Author(s) 2019

Abstract

The broad class of conditional transformation models includes interpretable and simple as well as potentially very complex models for conditional distributions. This makes conditional transformation models attractive for predictive distribution modelling, especially because models featuring interpretable parameters and black-box machines can be understood as extremes in a whole cascade of models. So far, algorithms and corresponding theory was developed for special forms of conditional transformation models only: maximum likelihood inference is available for rather simple models, there exists a tailored boosting algorithm for the estimation of additive conditional transformation models, and a special form of random forests targets the estimation of interaction models. Here, I propose boosting algorithms capable of estimating conditional transformation models of arbitrary complexity, starting from simple shift transformation models featuring linear predictors to essentially unstructured conditional transformation models allowing complex nonlinear interaction functions. A generic form of the likelihood is maximized. Thus, the novel boosting algorithms for conditional transformation models are applicable to all types of univariate response variables, including randomly censored or truncated observations.

Keywords Transformation model · Distribution regression · Conditional distribution function · Conditional quantile function · Conditional hazard function · Probabilistic forecasting

1 Introduction

The future remains unknown, yet we have witnessed considerably improved predictions owing to advances in statistical and machine learning over the last two decades. Numerous procedures, such as support vector machines, random forests, and tree boosting, deliver accurate point predictions of conditional means. However, in many applications, a mean prediction is not good enough. Full predictive distributions, also known as probabilistic forecasts, are required in applications where an assessment of the associated uncertainty is essential, for example in models of future disease progression (Küffner et al. 2015), electricity demand (Cabrera and Schulz 2017), stock asset returns (Mitrodima and Griffin 2017), and counterfactual distributions (Chernozhukov et al. 2013). In these applications, the prediction “takes the form of a predictive probability distribution over future quantities

or events of interest” (Gneiting and Katzfuss 2014). Here, I present a novel generic boosting approach to the estimation of full predictive distributions under mild assumptions.

Apart from completely model-free procedures (such as kernel smoothing, Li and Racine 2008), four main approaches of obtaining predictive distributions exist. (1) Flexible parametric models for conditional density functions rely on a strict parametric model of the response distribution those parameters might be linked to predictor variables in complex ways, for example, in generalized additive models for location, scale, and shape (GAMLSS, Rigby and Stasinopoulos 2005) and in heteroscedastic Bayesian additive regression tree ensembles (Pratola et al. 2017). (2) Quantile regression models for conditional quantiles of interest can be modelled in a linear or nonlinear additive form (Koenker 2005); more complex relationships can be estimated by quantile regression forests (Meinshausen 2006; Athey et al. 2019). (3) Distribution regression and transformation models potentially allow response-varying (or time-varying) effects (Foresi and Peracchi 1995; Rothe and Wied 2013; Chernozhukov et al. 2013; Wu and Tian 2013; Leorato and Peracchi 2015) in models for conditional distribution functions on the probit, logit, or complementary

✉ Torsten Hothorn
Torsten.Hothorn@uzh.ch

¹ Institut für Epidemiologie, Biostatistik und Prävention,
Universität Zürich, Hirschengraben 84, 8001 Zürich,
Switzerland

log–log scale. (4) Hazard regression (Koopberg et al. 1995) aims at estimating conditional nonproportional hazard functions directly.

Boosting, and especially the statistical view on boosting (Friedman et al. 2000; Bühlmann and Hothorn 2007), have already proved helpful in these four different approaches. Mayr et al. (2012) developed boosting for GAMLSS, conditional quantile boosting was introduced by Fenske et al. (2011), and nonproportional hazard boosting was recently introduced by Lee and Chen (2018). Distribution regression is a special case of conditional transformation models (Hothorn et al. 2014). What is interesting about conditional transformation models is that very simple models, such as the linear proportional odds and hazards models, and essentially unstructured models for conditional distribution functions can be understood in a unified theoretical framework (Hothorn et al. 2018). The same level of generality is, however, lacking from an algorithmic perspective. The boosting algorithm introduced by Hothorn et al. (2014) is limited to additive models and explicitly excludes tree-based interaction models. Furthermore, the target function is approximate and applicable to responses observed without censoring or truncation only. The aim of this work is to establish a general computational framework that allows specification, estimation, evaluation, and comparison in a cascade of models starting with very simple linear models and ending with essentially unstructured models for conditional distribution functions for arbitrary response variables.

Section 2 gives a dense introduction to transformation models. An elaborate description and connections to well-established models can be found in Hothorn et al. (2014) and Hothorn et al. (2018). Sections 3 and 4 develop two boosting algorithms for complex and simple transformation models based on a generic form of the likelihood (technical details regarding the definition of the likelihood for all types of response variables, including random censoring and truncation, are discussed by Hothorn et al. 2018). Empirical evaluations are presented in Sect. 5.

2 Transformation models

Let Y denote a univariate and at least ordered response variable on a measurable space $(\mathcal{E}, \mathcal{C})$ and $\mathbf{X} \in \mathcal{X}$ a set of predictor variables with joint distribution $(Y, \mathbf{X}) \sim \mathbb{P}_{Y, \mathbf{X}}$. Based on random samples from $\mathbb{P}_{Y, \mathbf{X}}$, the goal is to estimate the conditional distribution $\mathbb{P}_{Y|\mathbf{X}=\mathbf{x}}$ of a response given predictors. For each conditional cumulative distribution function $F_{Y|\mathbf{X}=\mathbf{x}}(y) = \mathbb{P}_{Y|\mathbf{X}=\mathbf{x}}(\{v \in \mathcal{E} \mid v < y\})$, a unique conditional transformation function $h : \mathcal{E} \times \mathcal{X} \rightarrow \mathbb{R}$ exists such that $F_{Y|\mathbf{X}=\mathbf{x}}(y) = F_Z(h(y | \mathbf{x}))$, assuming $F_Z : \mathbb{R} \rightarrow [0, 1]$ is an a priori given cumulative distribution function of an

absolutely continuous random variable Z with log-concave density f_Z (Hothorn et al. 2018). The conditional transformation function h is monotonic in y

$$h(\underline{y} | \mathbf{x}) \leq h(\bar{y} | \mathbf{x}) \quad \text{for all } \underline{y} < \bar{y} \in \mathcal{E}, \mathbf{x} \in \mathcal{X}. \quad (1)$$

Starting with Box and Cox (1964), shift transformation functions based on the decomposition $h(y | \mathbf{x}) = h_Y(y) - \beta(\mathbf{x})$ featuring a baseline transformation function $h_Y : \mathcal{E} \rightarrow \mathbb{R}$ and a shift term $\beta : \mathcal{X} \rightarrow \mathbb{R}$ have been studied intensively. The proportional hazards (with $F_Z(z) = 1 - \exp(-\exp(z))$) and proportional odds (with $F_Z(z) = \text{expit}(z)$) models are the most well-known representatives of this class of shift transformation models (STM, often also referred to as linear or nonlinear transformation models, depending on the functional form of $\beta(\mathbf{x})$). Boosting procedures that allow flexible estimation of $\beta(\mathbf{x})$ have been studied for proportional hazards models under right censoring (Ridgeway 1999; Schmid and Hothorn 2008; Lu and Li 2008; Yue et al. 2017) and proportional odds models have been studied for ordered responses (Schmid et al. 2011). A comparison of prominent and less prominent members of this model class is given in Hothorn et al. (2018).

Structured additive transformation functions that allow interactions between the two arguments y and \mathbf{x} of the form $h(y | \mathbf{x}) = \sum_{j=1}^J h_j(y | \mathbf{x})$ lead to conditional transformation models (CTM, Hothorn et al. 2014). The J partial transformation functions h_j allow formulation of problem-specific effects of the predictors \mathbf{x} , such as linear, nonlinear, spatio-temporal, or other model terms. Distribution regression models featuring response-varying effects are an important special case of this model class. When $\mathbf{x} = (x_1, \dots, x_J) \in \mathbb{R}^J$, a distribution regression model is characterized by partial transformation functions $h_j(y | x_j) = \beta_j(y)x_j$ and corresponding interpretable response-varying effects $\beta_j : \mathcal{E} \rightarrow \mathbb{R}$. The analogon of an additive model features partial transformation functions $h_j(y | x_j)$, i.e. bivariate smooth functions of both y and x_j . These bivariate terms are more complex than the one-dimensional coefficient functions $\beta_j(y)$ but can still be visualized and interpreted. If x_j is more complex, for example, if it describes a spatial location, $h_j(y | x_j)$ might be a spatially smooth term that captures unexplained spatial heterogeneity (Hothorn et al. 2014).

Models with transformation function $h(y | \mathbf{x}) = \sum_{j=1}^J h_j(y | x_j)$ and potential applications are discussed in Hothorn et al. (2014) and Hothorn et al. (2018). The standard estimation of maximizing the continuously ranked probability score over a discrete grid covering \mathcal{E} (Foresi and Peracchi 1995; Chernozhukov et al. 2013; Hothorn et al. 2014), potentially with inverse probability of censoring weight adjustment for right censoring (Möst and Hothorn 2015; Garcia et al. 2019), does not allow essentially unstructured transforma-

tion functions $h(y | \mathbf{x})$, including higher-order interactions, and thus relaxation of the additivity assumption on the scale of the transformation function h . Furthermore, it is computationally inefficient (because the data have to be expanded) and is unable to handle censoring or truncation directly.

This paper addresses these issues by introducing computationally efficient boosted likelihood estimation for unstructured or structured additive conditional transformation functions (Sect. 3) and shift transformation functions (Sect. 4) under all forms of random censoring and truncation for at least ordered responses based on potentially correlated observations.

3 Boosting the likelihood of conditional transformation models

In the following, the conditional transformation function $h(y | \mathbf{x}) = \mathbf{a}(y)^\top \boldsymbol{\vartheta}(\mathbf{x})$ is parameterized in terms of basis functions $\mathbf{a} : \mathcal{E} \rightarrow \mathbb{R}^P$ of the response and a conditional parameter function $\boldsymbol{\vartheta} : \mathcal{X} \rightarrow \mathbb{R}^P$; the latter function will be estimated.

3.1 Definition of the likelihood

The parameterisation of h implies a conditional cumulative distribution function $\mathbb{P}(Y \leq y | \mathbf{X} = \mathbf{x}) = F_Z(\mathbf{a}(y)^\top \boldsymbol{\vartheta}(\mathbf{x}))$ and thus a conditional density

$$\begin{aligned} f_Y(y | \boldsymbol{\vartheta}(\mathbf{x})) &= \frac{\partial F_Z(\mathbf{a}(y)^\top \boldsymbol{\vartheta}(\mathbf{x}))}{\partial y} \\ &= f_Z(\mathbf{a}(y)^\top \boldsymbol{\vartheta}(\mathbf{x})) \mathbf{a}'(y)^\top \boldsymbol{\vartheta}(\mathbf{x}) \end{aligned}$$

when $y \in \mathbb{R}$ comes from an absolutely continuous distribution (\mathbf{a}' is the derivative of \mathbf{a}). For discrete $y \in \mathcal{E} = \{y_1, y_2, \dots\}$, the density function is

$$\begin{aligned} f_Y(y_k | \boldsymbol{\vartheta}(\mathbf{x})) &= \begin{cases} F_Z(\mathbf{a}(y_k)^\top \boldsymbol{\vartheta}(\mathbf{x})) & k = 1 \\ F_Z(\mathbf{a}(y_k)^\top \boldsymbol{\vartheta}(\mathbf{x})) - F_Z(\mathbf{a}(y_{k-1})^\top \boldsymbol{\vartheta}(\mathbf{x})) & k > 1. \end{cases} \end{aligned}$$

There are also other forms of the density, for example, in mixed discrete-continuous distributions. The population optimizer for the conditional parameter function $\boldsymbol{\vartheta}$ is

$$\begin{aligned} \boldsymbol{\vartheta} &:= \arg \max_{\boldsymbol{\vartheta} : \mathcal{X} \rightarrow \mathbb{R}^P} \int \log\{f_Y[y | \boldsymbol{\vartheta}(\mathbf{x})]\} d\mathbb{P}_{Y, \mathbf{X}}(y, \mathbf{x}) \\ \text{st. (1).} \end{aligned}$$

Based on N independent observations (y_i, \mathbf{x}_i) , $i = 1, \dots, N$ from $\mathbb{P}_{Y, \mathbf{X}}$, empirical risk minimization with negative log-likelihood loss

$$\hat{\boldsymbol{\vartheta}}_N = \arg \max_{\boldsymbol{\vartheta} : \mathcal{X} \rightarrow \mathbb{R}^P} \sum_{i=1}^N \ell_i(\boldsymbol{\vartheta}(\mathbf{x}_i)) \quad \text{st. (1)}$$

can be applied to estimate the conditional parameter function $\boldsymbol{\vartheta}$. The log-likelihood contribution $\ell_i : \mathbb{R}^P \rightarrow \mathbb{R}$ for the i th observation is given by

$$\begin{aligned} \ell_i(\boldsymbol{\vartheta}(\mathbf{x}_i)) &= \begin{cases} \log\{f_Y[y_i | \boldsymbol{\vartheta}(\mathbf{x}_i)]\} & y_i \in \mathcal{E} \\ \log\left\{\int_{y_i} f_Y[y | \boldsymbol{\vartheta}(\mathbf{x}_i)] d\mu(y)\right\} & y_i \in \mathcal{C} \setminus \mathcal{E}, \end{cases} \end{aligned} \quad (2)$$

where the first case corresponds to an observation y_i from an absolutely continuous or ordered response and the second case corresponds to the situation where a set or interval was observed (for example, for a left-, right-, or interval-censored observation y_i). Integration is with respect to the measure μ dominating $\mathbb{P}_{Y|X=\mathbf{x}}$. Details on the likelihood function for models parameterized by $\mathbf{a}(y)^\top \boldsymbol{\vartheta}(\mathbf{x})$, including gradients (denoted \mathbf{u}_i in Algorithm 1) and Hessians under censoring and truncation, are given in Hothorn et al. (2018).

3.2 Boosting the likelihood

The proposed boosting Algorithm 1 outputs a model of the form

$$\boldsymbol{\vartheta}(\mathbf{x}) = \boldsymbol{\vartheta}^{[0]}(\mathbf{x}) + \sum_{b=1}^B \boldsymbol{\Gamma}^{[b]} \mathbf{b}_{j(b)}(\mathbf{x}), \quad (3)$$

with offset term $\boldsymbol{\vartheta}^{[0]}(\mathbf{x})$ and $j = 1, \dots, J$ a priori defined basis functions $\mathbf{b}_j : \mathcal{X} \rightarrow \mathbb{R}^{P_j}$ of the predictor variables. The function $j(b)$ returns the index of the basis function \mathbf{b}_j which was selected in the b th iteration of the algorithm. Each basis may be equipped with an explicit penalty function Pen_j . The corresponding penalty parameter λ_j is chosen such that the degrees of freedom are the same for all J basis functions to facilitate unbiased model selection (Hofner et al. 2011). The number of terms B , selected basis functions $j(b)$, and corresponding coefficient matrices $\boldsymbol{\Gamma}^{[b]} \in \mathbb{R}^{P \times P_{j(b)}}$ are unknowns and are estimated from data. The basis functions \mathbf{b}_j may feature unknown parameters. With relatively deep regression trees \mathbf{b}_j (where the tree structure is estimated from the data in every boosting iteration and $\boldsymbol{\Gamma}$ are the parameters in each terminal node), model (3) is the sum of B trees and as such is potentially highly unstructured. Similar to GAMLSS-boosting (Mayr et al. 2012), a parameter vector $\boldsymbol{\vartheta}$ is modelled instead of a scalar predictor function. The main difference is that all dimensions of the parameter vector $\boldsymbol{\vartheta}$ are updated simultaneously whereas each dimension is assigned its own predictor function in GAMLSS-boosting.

Algorithm 1 is essentially a multivariate version of L_2 boosting (Bühlmann and Yu 2003) using the negative trans-

Algorithm 1 Boosting CTM Likelihoods

- Start with $b = 0$ and offset $\boldsymbol{\vartheta}^{[0]}(\mathbf{x}_i)$. Initialize base learners $j = 1, \dots, J$ via basis functions \mathbf{b}_j , stepsize $v \in (0, 1)$, and $b_{\text{stop}} > 0$
- Iterate

1. $b \rightarrow b + 1$; stop if $b > b_{\text{stop}}$
2. Compute multivariate negative gradient $\mathbf{u}_i^{[b]} \in \mathbb{R}^P$ of $-\ell_i$ evaluated at $\boldsymbol{\vartheta}^{[b-1]}(\mathbf{x}_i)$

$$\mathbf{u}_i^{[b]} = \left. \frac{\partial \ell_i(\boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}} \right|_{\boldsymbol{\vartheta}=\boldsymbol{\vartheta}^{[b-1]}(\mathbf{x}_i)}$$

with components $\mathbf{u}_i^{[b]} = (u_{i,1}^{[b]}, \dots, u_{i,P}^{[b]})$ for $i = 1, \dots, N$

3. Fit base learners $j = 1, \dots, J$ by (penalized) multivariate least squares

$$\hat{\boldsymbol{\Gamma}}_j = \arg \min_{\boldsymbol{\Gamma}^\top = (\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_P)} \sum_{p=1}^P \left\{ \left[\sum_{i=1}^N \left(u_{ip}^{[b]} - \mathbf{b}_j(\mathbf{x}_i)^\top \boldsymbol{\gamma}_p \right)^2 \right] + \lambda_j \text{Pen}_j(\boldsymbol{\gamma}_p) \right\}$$

4. Select base learner with minimal quadratic error

$$j(b) = \arg \min_{j=1, \dots, J} \sum_{i=1}^N \left(u_i^{[b]} - \hat{\boldsymbol{\Gamma}}_j \mathbf{b}_j(\mathbf{x}_i) \right)^\top \left(u_i^{[b]} - \hat{\boldsymbol{\Gamma}}_j \mathbf{b}_j(\mathbf{x}_i) \right)$$

and $\boldsymbol{\Gamma}^{[b]} = v \hat{\boldsymbol{\Gamma}}_{j(b)}$

5. Update $\boldsymbol{\vartheta}^{[b]}(\mathbf{x}_i) = \boldsymbol{\vartheta}^{[b-1]}(\mathbf{x}_i) + \boldsymbol{\Gamma}^{[b]} \mathbf{b}_{j(b)}(\mathbf{x}_i)$ for $i = 1, \dots, N$

- Output $\hat{\boldsymbol{\vartheta}}_N(\mathbf{x}_i) = \boldsymbol{\vartheta}^{[b_{\text{stop}}]}(\mathbf{x}_i)$, $i = 1, \dots, N$

formation log-likelihood (2) as loss function. This choice makes the algorithm agnostic with respect to the scale of the response variable and potential censoring or truncation. The default offset is the unconditional maximum-likelihood estimator $\boldsymbol{\vartheta}^{[0]}(\mathbf{x}_i) \equiv \hat{\boldsymbol{\vartheta}}_{\text{ML}}$ for $i = 1, \dots, N$ that maximizes $\sum_{i=1}^N \ell_i(\boldsymbol{\vartheta})$. The algorithm is also applicable to the high-dimensional setting where the number of predictor variables exceeds the number of observations N . The number of boosting iterations b_{stop} is a tuning parameter that has to be chosen by the out-of-sample log-likelihood for a validation sample $i = N + 1, \dots, N + \tilde{N}$

$$\hat{B} = \hat{b}_{\text{stop}} = \arg \max_{b=0, 1, \dots} \sum_{i=N+1}^{N+\tilde{N}} \ell_i(\boldsymbol{\vartheta}^{[b]}(\mathbf{x}_i)).$$

Model choice, for example using cross-validation, subsampling, or the bootstrap, can also be implemented conveniently by comparing this out-of-sample log-likelihood of different candidate models.

An additional advantage of this algorithm over boosted continuously ranked probability scores (“CTM-CRPS-boosting”, Hothorn et al. 2014) is that computations of tensor products in $\mathbf{a}(y)^\top \otimes \mathbf{b}_j(\mathbf{x})^\top \text{vec}(\boldsymbol{\Gamma}) = \mathbf{a}(y)^\top \boldsymbol{\Gamma} \mathbf{b}_j(\mathbf{x})$ are never explicitly required because the linear array model

formulation (*i.e.* the right-hand side of the equation, see Currie et al. 2006, formula 2.5) formula 2.5 is implemented by Algorithm 1. This allows estimation of potentially highly unstructured models by choosing relatively deep multivariate regression trees as basis functions \mathbf{b} . Moreover, the algorithm does not require expansion of the data set (to size sample size N^2 , in the worst case).

3.3 Model interpretation

The partial transformation functions h_j can be obtained from the boosted model (3)

$$\begin{aligned} h_j(y | \mathbf{x}) &= \mathbf{a}(y)^\top \left(\sum_{b:j(b)=j} \boldsymbol{\Gamma}^{(b)} \mathbf{b}_{j(b)}(\mathbf{x}) \right) \\ &= \mathbf{a}(y)^\top \boldsymbol{\vartheta}_j(\mathbf{x}). \end{aligned}$$

The choice $\mathbf{b}_j(\mathbf{x}) = x_j$ results in the distribution regression model $F_{Y|X=\mathbf{x}}(y) = F_Z(h_Y(y) - \mathbf{x}^\top \boldsymbol{\beta}(y))$ with partial transformation functions $h_j(y | x_j) = x_j \beta_j(y) = x_j \mathbf{a}(y)^\top \boldsymbol{\vartheta}_j$ and corresponding response-varying effects $\beta_j(y) = \mathbf{a}(y)^\top \boldsymbol{\vartheta}_j$. Thus, this boosting procedure can also be used to estimate Cox models with time-varying effects under all forms of random censoring and truncation. Nonlinear effects can be implemented by a B-spline basis $\mathbf{b}_j(\mathbf{x}) = \mathbf{b}_j(x_j)$, and more complex bases allow specification of terms that capture spatio-temporal correlations or other forms of unexplained heterogeneity (Kneib et al. 2009; Hofner et al. 2011). A collection of commonly used basis functions \mathbf{b} , along with corresponding penalty functions and interpretable model terms, is reviewed in Mayr and Hofner (2018). Specific choices of basis functions underlying the empirical results presented in Sect. 5 are discussed in detail in Hothorn (2019).

4 Boosting the likelihood of shift transformation models

A comparison of conditional transformation models that allow interactions of y and \mathbf{x} in the transformation function h with shift transformation models in which these terms are absent can help to identify situations where the simpler models perform as good as or even better than the more complex models. Likelihood boosting for shift transformation models of the form $h(y | \mathbf{x}) = \mathbf{a}(y)^\top \boldsymbol{\vartheta} - \beta(\mathbf{x})$ is presented as Algorithm (2).

The procedure outputs a model

$$\beta(\mathbf{x}) = \beta^{[0]}(\mathbf{x}) + \sum_{b=1}^B \boldsymbol{\Gamma}^{[b]} \mathbf{b}_{j(b)}(\mathbf{x}),$$

Algorithm 2 Boosting STM Likelihoods

- Start with offset $\beta^{[0]}(\mathbf{x}) = 0$; other choices as in Algorithm 1
- Iterate
 1. $b \rightarrow b + 1$; stop if $b > b_{\text{stop}}$
 2. Update $\boldsymbol{\vartheta}^{[b]} = \arg \max_{\boldsymbol{\vartheta} \in \mathbb{R}^P} \sum_{i=1}^N \ell_i(\boldsymbol{\vartheta}, \beta^{[b-1]}(\mathbf{x}_i))$ and compute univariate negative gradient $u_i^{[b]} \in \mathbb{R}$ evaluated at $\beta^{[b-1]}(\mathbf{x}_i)$

$$u_i^{[b]} = \left. \frac{\partial \ell_i(\boldsymbol{\vartheta}^{[b]}, \beta)}{\partial \beta} \right|_{\beta=\beta^{[b-1]}(\mathbf{x}_i)}$$
 3. Fit base learners $j = 1, \dots, J$ by (penalized) multivariate least squares

$$\hat{\boldsymbol{\gamma}}_j = \arg \min_{\boldsymbol{\gamma}} \left\{ \left[\sum_{i=1}^N \left(u_i^{[b]} - \mathbf{b}_j(\mathbf{x}_i)^\top \boldsymbol{\gamma} \right)^2 \right] + \lambda_j \text{Pen}_j(\boldsymbol{\gamma}) \right\}$$
 4. Select base learner with minimal quadratic error

$$j(b) = \arg \min_{j=1, \dots, J} \sum_{i=1}^N \left(u_i^{[b]} - \mathbf{b}_j(\mathbf{x}_i)^\top \hat{\boldsymbol{\gamma}}_j \right)^2$$
 and $\boldsymbol{\gamma}^{[b]} = v \hat{\boldsymbol{\gamma}}_{j(b)}$
 5. Update $\beta^{[b]}(\mathbf{x}_i) = \beta^{[b-1]}(\mathbf{x}_i) + \mathbf{b}_{j(b)}(\mathbf{x}_i)^\top \boldsymbol{\gamma}^{[b]}$ for $i = 1, \dots, N$
- Output $\hat{\boldsymbol{\vartheta}}_N = \arg \max_{\boldsymbol{\vartheta} \in \mathbb{R}^P} \sum_{i=1}^N \ell_i(\boldsymbol{\vartheta}, \hat{\beta}_N(\mathbf{x}_i))$ and $\hat{\beta}_N(\mathbf{x}_i) = \beta^{[b_{\text{stop}}]}(\mathbf{x}_i), i = 1, \dots, N$

with univariate shift function $\beta(\mathbf{x}) \in \mathbb{R}$, i.e. with $\boldsymbol{\gamma}^\top = \boldsymbol{\Gamma} \in \mathbb{R}^{1 \times P_j}$. In contrast to conditional transformation models, the model term $\beta(\mathbf{x})$ does not depend on y and thus shift transformation models are easier to interpret. L_2 boosting in this setup is performed based on log-likelihood contributions $\ell_i(\boldsymbol{\vartheta}, \beta(\mathbf{x}_i))$ for $\ell_i : \mathbb{R}^{P+1} \rightarrow \mathbb{R}$ from densities

$$\begin{aligned} f_Y(y \mid \boldsymbol{\vartheta}, \beta(\mathbf{x})) &= \frac{\partial F_Z(\mathbf{a}(y)^\top \boldsymbol{\vartheta} - \beta(\mathbf{x}))}{\partial y} \\ &= f_Z(\mathbf{a}(y)^\top \boldsymbol{\vartheta} - \beta(\mathbf{x})) \mathbf{a}'(y)^\top \boldsymbol{\vartheta} \end{aligned}$$

for the absolutely continuous case and

$$\begin{aligned} f_Y(y_k \mid \boldsymbol{\vartheta}, \beta(\mathbf{x})) &= \begin{cases} F_Z(\mathbf{a}(y_k)^\top \boldsymbol{\vartheta} - \beta(\mathbf{x})) & k = 1 \\ F_Z(\mathbf{a}(y_k)^\top \boldsymbol{\vartheta} - \beta(\mathbf{x})) - F_Z(\mathbf{a}(y_{k-1})^\top \boldsymbol{\vartheta} - \beta(\mathbf{x})) & k > 1 \end{cases} \end{aligned}$$

for the discrete case. The core idea is to update the nuisance parameter $\boldsymbol{\vartheta}$ before computing the gradients in every boosting iteration (following Schmid and Hothorn 2008). For discrete proportional odds models, Algorithm (2) is equivalent to the procedure proposed by Schmid et al. (2011). The ability to handle censoring and truncation in the likelihood allows proportional hazards models with potentially very flexible log-hazard ratios $\beta(\mathbf{x})$ to be fitted to randomly

censored (including left- and interval-censoring) or truncated responses.

5 Empirical evaluation

The rationale for the empirical evaluation of Algorithms 1 and 2 was to demonstrate that interpretable transformation models can be estimated for applications where information about predictive distribution matters (Sect. 5.1) and to investigate the robustness of boosted transformation models under model misspecification (Sect. 5.2).

5.1 Applications

Eight life science applications in which estimation of a predictive distribution is of special interest are listed in Table 1. Four applications are described by a continuous response, two feature an ordered categorical response, and two feature a right-censored response. Except for the Beetle Extinction Risk application, which requires a discrete basis \mathbf{a} , a Bernstein basis \mathbf{a} of order $M = 6$ (for technical details see Hothorn et al. 2018) was used to parameterize the transformation functions. Conditional transformation models (Algorithm 1) with nonlinear (N, using B-splines), linear (L), and tree-based (T, of depth two and thus allowing only two-way interactions) basis functions \mathbf{b} as well as shift transformation models (Algorithm 2) using the same bases were evaluated. The performance of these boosted transformation models was compared to the performance of transformation trees and transformation forests (Hothorn and Zeileis 2017). The latter two procedures estimate conditional transformation models of the form $F_Z(\mathbf{a}(y)^\top \boldsymbol{\vartheta}(\mathbf{x}))$, where $\boldsymbol{\vartheta}(\mathbf{x})$ is obtained either from a single tree (transformation trees) or from a nonlinear interaction function (transformation forest). I hypothesized a priori that transformation trees should perform worst across all applications because this method corresponds to the most simple (but easily interpretable) model. Also, I expected transformation forests to outperform transformation trees and to perform only slightly worse than the best performing boosting procedure because of the high adaptivity of the underlying random forest procedure. My motivation for this experiment was the hope that I would be able to find a simple and interpretable transformation model that outperforms the most complex transformation forests by means of either Algorithm 1 or 2.

Subsampling (with $n = 3/4N$ observations in the learning and $\tilde{n} = 1/4N$ observations in the validation sample) was performed 100 times. Performance of the competitors was assessed by the out-of-sample log-likelihood centered by the out-of-sample log-likelihood of the unconditional transformation model $F_Z(\mathbf{a}(y)^\top \hat{\boldsymbol{\vartheta}}_{\text{ML}})$. For a learning sample of size

Table 1 Applications. Eight prediction problems with continuous, ordered categorical (number of categories in parentheses), or right-censored response (per cent censored in parentheses)

Application	Source	N	Y	X	F_Z
Beetle Extinction Risk	Seibold et al. (2015)	1025	Ordered (6)	(7, 3)	expit
Birth Weight Prediction	Schild et al. (2008)	150	Continuous	(5, 0)	Φ
Body Fat Mass	Garcia et al. (2005)	71	Continuous	(9, 0)	Φ
CAO/ARO/AIO-04 DFS	Rödel et al. (2015)	1153	Survival (71%)	(3, 15)	MEV
Childhood Malnutrition	Fenske et al. (2011)	24166	Continuous	(6, 14)	Φ
Head Circumference	Fredriks et al. (2000)	7040	Continuous	(1, 0)	Φ
PRO-ACT ALSFRS	Küffner et al. (2015)	1013	Ordered (50)	(43, 16)	expit
PRO-ACT OS	Seibold et al. (2017)	2711	Survival (69%)	(3, 16)	MEV

Number of complete observations N and number of (numeric, categorical) predictor variables X . F_Z is the cumulative distribution function of the standard normal (Φ), standard logistic (expit), or standard minimum extreme value distribution (MEV, the inverse complementary log–log)

n and a validation sample $i = n + 1, \dots, \tilde{n}$, the centered out-of-sample log-likelihood is given by

$$\tilde{n}^{-1} \sum_{i=n+1}^{n+\tilde{n}} \ell_i \left(\hat{\boldsymbol{\theta}}^{[\hat{b}_{\text{stop}}]}(\mathbf{x}_i) \right) - \tilde{n}^{-1} \sum_{i=n+1}^{n+\tilde{n}} \ell_i(\hat{\boldsymbol{\theta}}_{\text{ML}}).$$

Values close to zero indicate that the conditional model did not outperform the unconditional model.

The results presented in Table 2 demonstrate that the best-performing method was always a boosted transformation model. Transformation forests performed only slightly worse than the top model for the Beetle Extinction Risk, Birth Weight Prediction, Body Fat Mass, and Childhood Malnutrition applications. In the remaining four applications, the best boosting procedure outperformed transformation forests substantially. Nonlinear conditional transformation models (N $\boldsymbol{\theta}(\mathbf{x})$) performed best twice, as did tree-based shift transformation models (T $\beta(\mathbf{x})$). Each of the remaining models ranked at the top once. Transformation trees outperformed transformation forests for two applications (Head Circumference and PRO-ACT ALSFRS) but never performed better than any of the boosted transformation models.

Graphical representations of the distributions of out-of-sample log-likelihoods along with the exact model and algorithm specification and corresponding software implementation are presented for all eight applications in Hothorn (2019).

5.2 Artificial data-generating processes

The response Y was generated conditionally on two groups and one numeric predictor variable $x \in [0, 1]$ following a transformation model of the form

$$\mathbb{P}(Y \leq y \mid \text{Group}, x) = \Phi(h(y \mid \text{Group}, x)),$$

where the conditional transformation function $h(y \mid \text{Group}, x) = \Phi^{-1}(\mathbb{P}(Y \leq y \mid \text{Group}, x))$ for four data-

generating processes (DGPs) is given in Table 3. The model labelled “Linear $\beta(\mathbf{x})$ ” is a shift transformation model with a main effect of group, a linear main effect of x , and a corresponding linear interaction effect. The linear main and interaction effects of x are replaced by nonlinear effects (a scaled sin function) of x in the shift transformation model “Nonlinear $\beta(\mathbf{x})$ ”. The extension to response-varying main and interaction effects defines the distribution regression model “Linear $\boldsymbol{\theta}(\mathbf{x})$ ” and the conditional transformation model “Nonlinear $\boldsymbol{\theta}(\mathbf{x})$ ”. The coefficients of the terms introduced in Table 3 are given in Table 4. Details of the implementation of these DGPs are explained in Hothorn (2019). The conditional densities associated with the four DGPs are shown in Fig. 1.

Models were evaluated by out-of-sample log-likelihoods centered by the out-of-sample log-likelihood of the true DGP (for test samples of size $\tilde{N} = 2000$)

$$\tilde{N}^{-1} \sum_{i=N+1}^{N+\tilde{N}} \ell_i \left(\hat{\boldsymbol{\theta}}^{[\hat{b}_{\text{stop}}]}(\mathbf{x}_i) \right) - \tilde{N}^{-1} \sum_{i=N+1}^{N+\tilde{N}} \ell_i(\boldsymbol{\theta}_{\text{True}}),$$

where $\boldsymbol{\theta}_{\text{True}}$ is as given in Tables 3 and 4.

In Part A of this simulation, nonlinear (N, using B-splines), linear (L), and tree-based (T, of depth six, which allows higher-order interactions) basis functions \mathbf{b} for shift transformation models, *i.e.* models for $\beta(\mathbf{x})$, and for conditional transformation models, *i.e.* models for $\boldsymbol{\theta}(\mathbf{x})$, were evaluated for sample sizes $N = 75, 150, 300$ under correctly specified models; this means that models were fitted using the correct distribution function $F_Z = \Phi$, the correct order $M = 6$ of \mathbf{a} , no uninformative predictor variables, and the correct basis functions. Both the linear and nonlinear models were fitted with basis functions representing a main effect of group, a main effect of x , and a corresponding interaction effect, whereas trees had to learn this structure from the data. In Part B, these models were evaluated under model misspecification, *i.e.* using the incorrect distribution function

Table 2 Applications. Median out-of-sample log-likelihoods (centered by out-of-sample log-likelihoods of the corresponding unconditional model) from 100 subsampled divisions into learning and test samples

	Parm. $b(\mathbf{x})$		Conditional Transformation Model $\theta(\mathbf{x})$		Shift Transformation Model $\beta(\mathbf{x})$		Transformation	
	N	L	T	N	L	T	Tree	Forest
Beetles Extinction Risk	0.27 (0.15, 0.28)	0.26 (0.16, 0.28)	0.18 (0.11, 0.29)	0.30 (0.29, 0.32)	0.29 (0.27, 0.31)	0.33 (0.30, 0.35)	0.03 (−0.13, 0.12)	0.30 (0.28, 0.33)
Birth Weight Prediction	1.47 (1.31, 1.55)	1.43 (1.19, 1.51)	1.25 (1.13, 1.30)	1.50 (1.32, 1.57)	1.50 (1.31, 1.56)	1.37 (1.28, 1.46)	0.39 (−0.17, 0.69)	1.32 (1.21, 1.37)
Body Fat Mass	1.17 (1.05, 1.27)	1.08 (0.96, 1.20)	0.76 (0.58, 0.87)	0.12 (0.11, 0.14)	0.12 (0.11, 0.14)	0.24 (0.22, 0.27)	0.08 (−0.61, 0.49)	1.03 (0.95, 1.09)
CAO/ARO/AIO-04 DFS	0.00 (−0.00, 0.01)	0.01 (−0.00, 0.01)	0.00 (−0.00, 0.00)	0.02 (0.01, 0.03)	0.02 (0.01, 0.03)	0.01 (0.01, 0.02)	−0.00 (−0.01, 0.00)	−0.01 (−0.02, −0.00)
Childhood Malnutrition	0.05 (0.05, 0.05)	0.05 (0.04, 0.05)	0.15 (0.14, 0.15)	0.13 (0.13, 0.14)	0.12 (0.11, 0.12)	0.14 (0.13, 0.14)	0.12 (0.12, 0.13)	0.15 (0.14, 0.15)
Head Circumference	1.09 (1.08, 1.11)	0.96 (0.94, 0.97)	1.09 (1.07, 1.10)	1.01 (1.00, 1.02)	0.96 (0.95, 0.97)	1.08 (1.07, 1.10)	1.06 (1.04, 1.08)	0.99 (0.92, 1.02)
PRO-ACT ALSFRS	0.52 (0.50, 0.55)	0.52 (0.49, 0.55)	0.48 (0.46, 0.50)	0.51 (0.47, 0.54)	0.49 (0.46, 0.52)	0.49 (0.46, 0.53)	0.38 (0.33, 0.43)	0.32 (0.30, 0.33)
PRO-ACT OS	0.04 (0.03, 0.05)	0.04 (0.03, 0.05)	0.04 (0.02, 0.05)	0.05 (0.04, 0.06)	0.05 (0.04, 0.06)	0.06 (0.05, 0.06)	0.02 (−0.02, 0.04)	0.04 (−0.01, 0.05)

Positive values indicate a superior performance of the conditional model, taking predictor variables into account, over the unconditional model. Interquartile ranges are given in parentheses. Boosting CTM Likelihoods (parameter $\theta(\mathbf{x})$) with nonlinear (N), linear (L), and tree-based (T, depth two) basis functions \mathbf{b} as well as Boosting STM Likelihoods (parameter $\beta(\mathbf{x})$) with the same basis functions were compared to transformation trees and transformation forests. The result of the best-performing model is printed in boldface

Table 3 Artificial data-generating processes (DGPs). Description of four simulation models

DGP	$\Phi^{-1}(\mathbb{P}(Y \leq y \mid \text{Group 1}, x))$	$\Phi^{-1}(\mathbb{P}(Y \leq y \mid \text{Group 2}, x))$
Linear $\beta(x)$	$h_Y(y) - 2x$	$h_Y(y) + 2 - x$
Nonlinear $\beta(x)$	$h_Y(y) - 2g(x)$	$h_Y(y) + 2 - g(x)$
Linear $\vartheta(x)$	$h_Y(y) - \beta_1(y) - \beta_2(y)x$	$h_Y(y) + \beta_1(y) - (\beta_2(y) + \beta_3(y))x$
Nonlinear $\vartheta(x)$	$h_Y(y) - \beta_1(y) - \beta_2(y)g(x)$	$h_Y(y) + \beta_1(y) - (\beta_2(y) + \beta_3(y))g(x)$

$g(x) = \sin(2\pi x)(1 + x)$, $h_Y, \beta_1, \beta_2, \beta_3$ are Bernstein polynomials of order $M = 6$ on the interval $(-4, 6)$ (see Hothorn et al. 2018)

Table 4 Artificial data-generating Processes (DGPs). Coefficients of baseline transformation h_Y and response-varying effects β_1, β_2 , and β_3

	ϑ_1	ϑ_2	ϑ_3	ϑ_4	ϑ_5	ϑ_6	ϑ_7
$h_Y(y)$	− 4.000	− 0.601	1.065	1.000	2.667	4.333	6.000
$\beta_1(y)$	0.000	− 0.518	− 1.000	− 1.414	− 1.732	− 1.932	− 2.000
$\beta_2(y)$	0.000	0.816	1.155	1.414	1.633	1.826	2.000
$\beta_3(y)$	0.000	− 0.259	− 0.500	− 0.707	− 0.866	− 0.966	− 1.000

All functions are parameterized as a Bernstein polynomial $a(y)^T \vartheta$ of order $M = 6$ on the interval $(-4, 6)$ (Hothorn et al. 2018), whose coefficients $\vartheta = (\vartheta_1, \dots, \vartheta_7)$ are given in this table

$F_Z = \text{expit}$ (standard logistic distribution) or $F_Z = \text{MEV}$ (standard minimum extreme value distribution), a too large dimension of the Bernstein basis a ($M = 12$), or $J^+ = 5$, 25 additional uninformative uniform predictor variables. The same “correct” basis functions as in Part A were used in Part B. I hypothesized a priori that models exactly matching the DGP would perform best and that tree-based boosting would outperform boosting with linear basis functions in nonlinear problems. Under misspecification, I expected the performance of all models to decrease, but this general ranking to persist.

The results for Part A presented in the top three rows of Table 5 show that the model corresponding to the underlying DGP was associated with the largest median out-of-sample log-likelihood. For linear DGPs, the performance of boosted models with nonlinear basis functions was only slightly inferior to the performance of boosted models with linear basis functions, while tree-based boosting performed substantially worse in this situation. By contrast, the signal in nonlinear DGPs was captured relatively well by tree-based boosting, whereas linear basis functions were not able to recover this signal. This shows that tree-based boosting was able to adapt to the underlying nonlinear interaction signal in the two nonlinear simulation models “Nonlinear $\beta(x)$ ” and “Nonlinear $\vartheta(x)$ ”.

The out-of-sample log-likelihoods for misspecified models presented in Table 5, Part B for $F_Z = \Phi$, follow this general pattern in that the model corresponding to the DGP performed best and tree-based boosting outperformed boosting with linear basis functions on nonlinear problems. In only two cases, which were characterized by small samples, did a linear model for $\vartheta(x)$ outperform a true linear model for $\beta(x)$ or vice versa. More frequently, the too complex nonlinear model for $\vartheta(x)$ outperformed the nonlinear model

for $\beta(x)$ slightly. Overall, Algorithms 1 and 2 seemed to be robust against overly complex basis functions a and additional noninformative predictor variables.

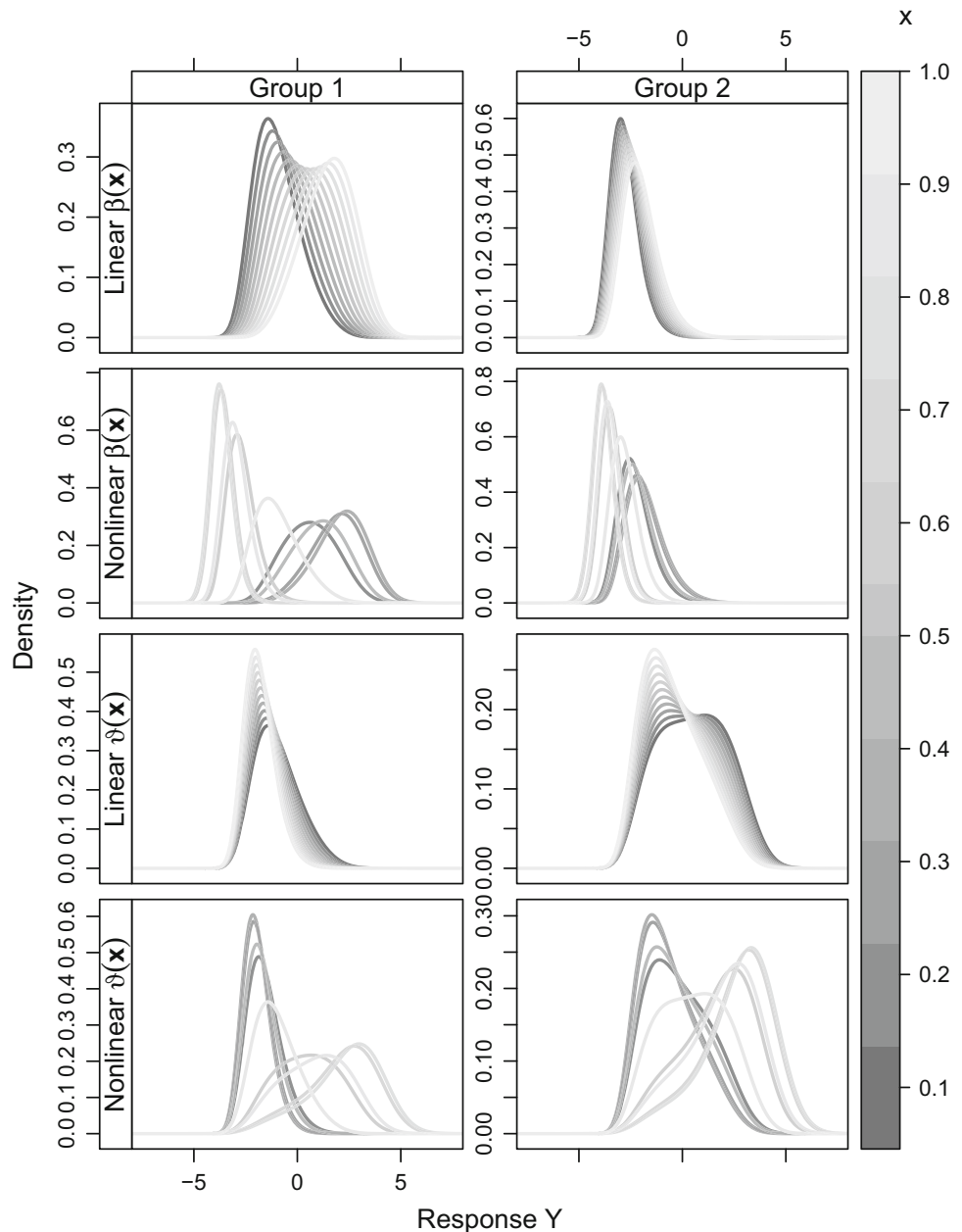
This was also true under a misspecified distribution function $F_Z = \text{expit}$ for linear shift transformation models “Linear $\beta(x)$ ”. More severe deviations occurred when an incorrect $F_Z = \text{expit}$ was used for model specification in Algorithms 1 and 2 under the nonlinear shift transformation model “Nonlinear $\beta(x)$ ”, distribution regression model “Linear $\vartheta(x)$ ”, and conditional transformation model “Nonlinear $\vartheta(x)$ ”. The absolute differences in the corresponding out-of-sample log-likelihoods were, however, marginal in most of these cases.

When the asymmetric standard minimum value distribution was used ($F_Z = \text{MEV}$), the distortions were more pronounced. The general pattern observed for $F_Z = \text{expit}$ was the same, but the centered out-of-sample log-likelihoods seemed in general smaller in this setup. Visualizations of the distributions underlying the figures in Table 5 are presented in Hothorn (2019).

6 Discussion

Models defined in terms of simple linear transformation functions up to models featuring unstructured complex transformation functions can be specified, estimated, evaluated, and compared in the unified computational framework of Algorithms 1 and 2. Data analysts are no longer limited in their freedom to define and estimate transformation models, because the strong ties between models of a certain complexity and a tailored estimation procedure (such as CTM-CPRS-boosting for additive or transformation forests

Fig. 1 Artificial Data-generating Processes (DGPs). Conditional densities given two groups (left and right panel) and $x \in [0, 1]$ (gray color coding) for four different data-generating processes



for interaction models) can be cut with the boosting algorithms presented here.

For model specification, the choice of F_Z is important in simple shift transformation models because it affects the interpretation of model parameters (log-odds ratios vs. log-hazard ratios, for example). In more complex models, a direct interpretation of parameters is hardly possible, and the estimated conditional distribution functions are insensitive to the choice of F_Z (Hothorn et al. 2018). However, one could use Algorithm 1 to estimate an unstructured log-hazard function $\mathbf{a}(y)^\top \boldsymbol{\vartheta}(\mathbf{x}) + \log(\mathbf{a}'(y)^\top \boldsymbol{\vartheta}(\mathbf{x}))$ in the model $\mathbb{P}(Y \leq y \mid \mathbf{X} = \mathbf{x}) = 1 - \exp(-\exp(\mathbf{a}(y)^\top \boldsymbol{\vartheta}(\mathbf{x})))$, with $\boldsymbol{\vartheta}(\mathbf{x})$ being, for example, the sum of B deep trees. The log-

likelihood risk function employed here, which is also able to handle time-varying covariates through appropriate truncation, avoids the technical obstacles reported by Lee and Chen (2018) when defining an appropriate nonparametric risk function for boosting in a class of models for conditional log-hazard functions.

In contrast to quantile regression, where separate models for each quantile are fitted, likelihood boosting for transformation models estimates conditional distribution functions simultaneously for all quantiles. It is interesting to note that a recently suggested Bayesian approach to simultaneous linear quantile regression (Yang and Tokdar 2017) maximizes a log-likelihood obtained from a numerical inver-

F_Z	M	J^+	N	DGP Param. b(x)	Linear $\beta(x)$				Nonlinear $\beta(x)$				Linear $\vartheta(x)$				Nonlinear $\vartheta(x)$											
					$\beta(x)$		$\vartheta(x)$		$\beta(x)$		$\vartheta(x)$		$\beta(x)$		$\vartheta(x)$		$\beta(x)$		$\vartheta(x)$									
					N	L	T	N	L	T	N	L	T	N	L	T	N	L	T	N	L	T						
Φ	6	0	75	Part	-135	-90	-373	-86	-67	-389	-189	-725	-414	-202	-799	-504	-125	-87	-217	-130	-106	-284	-143	-383	-348	-183	-422	-423
					-86	-51	-301	-47	-34	-280	-116	-694	-322	-87	-764	-364	-68	-51	-171	-91	-86	-228	-93	-356	-288	-126	-386	-329
					-55	-32	-230	-22	-17	-194	-74	-673	-250	-36	-746	-252	-51	-32	-122	-71	-68	-172	-59	-336	-224	-102	-363	-244
					-216	-149	-409	-128	-95	-491	-266	-768	-624	-295	-824	-746	-163	-115	-224	-138	-133	-293	-212	-418	-414	-250	-433	-539
					-142	-84	-315	-69	-53	-367	-185	-717	-466	-157	-782	-534	-94	-67	-176	-101	-93	-267	-134	-368	-333	-163	-390	-422
	25	75	300	Part	-286	-56	-254	-40	-27	-259	-125	-686	-370	-65	-752	-358	-63	-46	-139	-74	-73	-212	-90	-346	-266	-117	-373	-325
					-258	-203	-405	-138	-120	-523	-385	-688	-406	-862	-897	-222	-162	-219	-167	-148	-333	-287	-477	-443	-322	-442	-609	-422
					-175	-129	-344	-83	-70	-419	-235	-749	-566	-249	-803	-705	-130	-97	-179	-107	-101	-291	-170	-400	-355	-188	-404	-511
					-114	-81	-288	-53	-42	-307	-156	-701	-441	-144	-769	-481	-81	-63	-142	-81	-78	-233	-117	-360	-304	-142	-380	-384
					-170	-131	-394	-107	-95	-482	-228	-743	-468	-242	-840	-567	-123	-99	-240	-172	-150	-338	-175	-421	-364	-233	-459	-471
12	0	75	300	Part	-103	-70	-312	-63	-45	-315	-146	-706	-381	-143	-772	-384	-72	-51	-175	-98	-96	-250	-104	-367	-281	-143	-397	-363
					-65	-40	-245	-29	-22	-199	-84	-678	-306	-46	-751	-259	-44	-30	-129	-78	-71	-184	-62	-343	-230	-109	-376	-260
					-247	-186	-412	-154	-126	-601	-315	-798	-574	-337	-879	-844	-197	-147	-240	-192	-170	-334	-250	-457	-439	-284	-475	-646
					-153	-108	-345	-76	-67	-402	-213	-728	-536	-195	-803	-590	-99	-70	-177	-109	-101	-273	-137	-377	-358	-178	-412	-450
					-97	-68	-260	-47	-32	-288	-141	-691	-403	-101	-760	-371	-59	-43	-133	-82	-76	-224	-86	-349	-271	-121	-379	-336
25	75	300	Part	-309	-241	-433	-176	-147	-630	-387	-846	-765	-496	-918	-990	-231	-186	-256	-216	-188	-356	-318	-533	-495	-385	-503	-666	
				-186	-145	-362	-98	-82	-456	-252	-753	-612	-282	-827	-780	-129	-94	-188	-137	-117	-301	-151	-410	-374	-218	-414	-534	
				-123	-94	-285	-57	-46	-331	-175	-702	-456	-154	-768	-516	-80	-61	-142	-87	-80	-245	-112	-364	-294	-144	-385	-400	
				-165	-107	-415	-98	-73	-340	-198	-746	-461	-197	-813	-454	-132	-89	-212	-122	-101	-2							

Part A compares correctly specified models ($F_Z = \Phi$, $J^+ = 0$, $M = 6$) for three different sample sizes. Part B compares the performance under various forms of misspecification. The result of the best-performing model is printed in boldface

sion of the quantile function instead of using the traditional check risk minimization. In light of this approach, it seems computationally attractive to model the distribution function in the distribution regression model $F_{Y|X=x}(y | X = x) = F_Z(h_Y(y) - x^\top \beta(y))$ rather than the quantile function in a quantile regression model $Q_{Y|X=x}(\tau | X = x) = \alpha(\tau) + x^\top \delta(\tau)$ of the same complexity ($\tau \in [0, 1]$; α and δ being the probability-varying intercept and coefficient functions, respectively). Bayesian inference for the corresponding model parameters in conditional transformation models is, however, still under development (Mitrodimia and Griffin 2017).

Computational details

A reference implementation of transformation boosting machines (Algorithms 1 and 2) is available in the **tbm** package (Hothorn 2019). Analyses of all applications and simulation results can be reproduced in the dynamic document Hothorn (2019). All computations were performed using R version 3.5.2 (R Core Team 2018).

Open Access This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

References

- Athey, S., Tibshirani, J., Wager, S.: Generalized random forests. *Ann. Stat.* **47**(2), 1148–1178 (2019). <https://doi.org/10.1214/18-AOS1709>
- Box, G.E.P., Cox, D.R.: An analysis of transformations. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **26**(2), 211–252 (1964)
- Bühlmann, P., Hothorn, T.: Boosting algorithms: regularization, prediction and model fitting. *Stat. Sci.* **22**(4), 477–505 (2007). <https://doi.org/10.1214/07-STS242>. with discussion
- Bühlmann, P., Yu, B.: Boosting with the L_2 loss: regression and classification. *J. Am. Stat. Assoc.* **98**(462), 324–339 (2003). <https://doi.org/10.1198/0162145030000125>
- Cabrera, B.L., Schulz, F.: Forecasting generalized quantiles of electricity demand: a functional data approach. *J. Am. Stat. Assoc.* **112**(517), 127–136 (2017). <https://doi.org/10.1080/01621459.2016.1219259>
- Chernozhukov, V., Fernández-Val, I., Melly, B.: Inference on counterfactual distributions. *Econometrica* **81**(6), 2205–2268 (2013). <https://doi.org/10.3982/ECTA10582>
- Currie, I.D., Durban, M., Eilers, P.H.C.: Generalized linear array models with applications to multidimensional smoothing. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **68**(2), 259–280 (2006). <https://doi.org/10.1111/j.1467-9868.2006.00543.x>
- Fenske, N., Kneib, T., Hothorn, T.: Identifying risk factors for severe childhood malnutrition by boosting additive quantile regression. *J. Am. Stat. Assoc.* **106**(494), 494–510 (2011). <https://doi.org/10.1198/jasa.2011.ap09272>
- Foresi, S., Peracchi, F.: The conditional distribution of excess returns: an empirical analysis. *J. Am. Stat. Assoc.* **90**(430), 451–466 (1995). <https://doi.org/10.1080/01621459.1995.10476537>
- Fredriks, A.M., van Buuren, S., Burgmeijer, R.J.F., Meulmeester, J.F., Beuker, R.J., Brugman, E., Roede, M.J., Verloove-Vanhorick, S.P., Wit, J.: Continuing positive secular growth change in the Netherlands 1955–1997. *Pediatr. Res.* **47**(3), 316–323 (2000)
- Friedman, J.H., Hastie, T., Tibshirani, R.: Additive logistic regression: a statistical view of boosting (with discussion). *Ann. Stat.* **28**, 337–407 (2000). <https://doi.org/10.1214/aos/1016218223>
- Garcia, A.L., Wagner, K., Hothorn, T., Koebnick, C., Zunft, H.J.F., Trippo, U.: Improved prediction of body fat by measuring skinfold thickness, circumferences, and bone breadths. *Obesity* **13**(3), 626–634 (2005). <https://doi.org/10.1038/oby.2005.67>
- Garcia, T.P., Marder, K., Wang, Y.: Time-varying proportional odds model for mega-analysis of clustered event times. *Biostatistics* **20**(1), 129–146 (2019). <https://doi.org/10.1093/biostatistics/kxx065>
- Gneiting, T., Katzfuss, M.: Probabilistic forecasting. *Annu. Rev. Stat. Its Appl.* **1**(1), 125–151 (2014). <https://doi.org/10.1146/annurev-statistics-062713-085831>
- Hofner, B., Hothorn, T., Kneib, T., Schmid, M.: A framework for unbiased model selection based on boosting. *J. Comput. Graph. Stat.* **20**(4), 956–971 (2011). <https://doi.org/10.1198/jcgs.2011.09220>
- Hothorn, T.: **tbm: Transformation Boosting Machines**. R package and vignette version 0.3-0 (2019). <http://CRAN.R-project.org/package=tbm>
- Hothorn, T., Zeileis, A.: Transformation forests. Tech. rep. v2, <https://arxiv.org/abs/1701.02110> (2017)
- Hothorn, T., Kneib, T., Bühlmann, P.: Conditional transformation models. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **76**(1), 3–27 (2014). <https://doi.org/10.1111/rssb.12017>
- Hothorn, T., Möst, L., Bühlmann, P.: Most likely transformations. *Scand. J. Stat.* **45**(1), 110–134 (2018). <https://doi.org/10.1111/sjos.12291>
- Kneib, T., Hothorn, T., Tutz, G.: Variable selection and model choice in geospatial regression models. *Biometrics* **65**(2), 626–634 (2009). <https://doi.org/10.1111/j.1541-0420.2008.01112.x>
- Koenker, R.: Quantile Regression. Economic Society Monographs. Cambridge University Press, New York (2005)
- Kooperberg, C., Stone, C.J., Truong, Y.K.: Hazard regression. *J. Am. Stat. Assoc.* **90**(429), 78–94 (1995). <https://doi.org/10.1080/01621459.1995.10476491>
- Küffner, R., Zach, N., Norel, R., Hawe, J., Schoenfeld, D., Wang, L., Li, G., Fang, L., Mackey, L., Hardiman, O., Cudkowicz, M., Sherman, A., Ertaylan, G., Grosse-Wentrup, M., Hothorn, T., van Ligtengen, J., Macke, J.H., Meyer, T., Schölkopf, B., Tran, L., Vaughan, R., Stolovitzky, G., Leitner, M.L.: Crowdsourced analysis of clinical trial data to predict amyotrophic lateral sclerosis progression. *Nat. Biotechnol.* **33**, 51–57 (2015). <https://doi.org/10.1038/nbt.3051>
- Lee, D.K.K., Chen, N.: Boosting hazard regression with time-varying covariates. Tech. rep. v3, <https://arxiv.org/abs/1701.07926> (2018)
- Leorato, S., Peracchi, F.: Comparing distribution and quantile regression. Tech. Rep. 1511, Einaudi Institute for Economics and Finance, Rome, Italy (2015). <https://ideas.repec.org/p/eie/wpaper/1511.html>. Accessed 24 Nov 2018
- Li, Q., Racine, J.S.: Nonparametric estimation of conditional CDF and quantile functions with mixed categorical and continuous data. *J. Bus. Econ. Stat.* **26**(4), 423–434 (2008). <https://doi.org/10.1198/073500107000000250>
- Lu, W., Li, L.: Boosting method for nonlinear transformation models with censored survival data. *Biostatistics* **9**(4), 658–667 (2008). <https://doi.org/10.1093/biostatistics/kxn005>
- Mayr, A., Hofner, B.: Boosting for statistical modelling—a non-technical introduction. *Stat. Model.* **18**(3–4), 365–384 (2018). <https://doi.org/10.1177/1471082X17748086>

- Mayr, A., Fenske, N., Hofner, B., Kneib, T., Schmid, M.: GAMLSS for high-dimensional data—a flexible approach based on boosting. *J. R. Stat. Soc. Ser. C (Appl. Stat.)* **61**(3), 403–427 (2012). <https://doi.org/10.1111/j.1467-9876.2011.01033.x>
- Meinshausen, N.: Quantile regression forests. *J. Mach. Learn. Res.* **7**, 983–999 (2006). <http://jmlr.org/papers/v7/meinshausen06a.html>
- Mitrodimas, G., Griffin, J.E.: A Bayesian quantile time series model for asset returns. Tech. rep., SSRN, <https://doi.org/10.2139/ssrn.3050989> (2017)
- Möst, L., Hothorn, T.: Conditional transformation models for survivor function estimation. *Int. J. Biostat.* **11**(1), 23–50 (2015). <https://doi.org/10.1515/ijb-2014-0006>
- Pratola, M., Chipman, H., George, E.I., McCulloch, R.: Heteroscedastic bart using multiplicative regression trees. Tech. rep. v1, <http://arxiv.org/abs/1709.07542> (2017)
- R Core Team: R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria (2018). <http://www.R-project.org/>
- Ridgeway, G.: The state of boosting. *Comput. Sci. Stat.* **31**, 172–181 (1999)
- Rigby, R.A., Stasinopoulos, D.M.: Generalized additive models for location, scale and shape. *J. R. Stat. Soc. Ser. C (Appl. Stat.)* **54**(3), 507–554 (2005). <https://doi.org/10.1111/j.1467-9876.2005.00510.x>
- Rödel, C., Graeven, U., Fietkau, R., Hohenberger, W., Hothorn, T., Arnold, D., Hofheinz, R.D., Ghadimi, M., Wolff, H.A., Lang-Welzenbach, M., Raab, H.R., Wittekind, C., Ströbel, P., Staib, L., Wilhelm, M., Grabenbauer, G.G., Hoffmanns, H., Lindemann, F., Schlenska-Lange, A., Folprecht, G., Sauer, R.: Torsten Liersch on behalf of the German Rectal Cancer Study Group: Oxaliplatin added to fluorouracil-based preoperative chemoradiotherapy and postoperative chemotherapy of locally advanced rectal cancer (the German CAO/ARO/AIO-04 study): final results of the multicentre, open-label, randomised, phase 3 trial. *Lancet Oncol.* **16**(8), 979–989 (2015). [https://doi.org/10.1016/S1470-2045\(15\)00159-X](https://doi.org/10.1016/S1470-2045(15)00159-X)
- Rothe, C., Wied, D.: Misspecification testing in a class of conditional distributional models. *J. Am. Stat. Assoc.* **108**(501), 314–324 (2013). <https://doi.org/10.1080/01621459.2012.736903>
- Schild, R.L., Maringa, M., Siemer, J., Meurer, B., Hart, N., Goecke, T.W., Schmid, M., Hothorn, T., Hansmann, M.E.: Weight estimation by three-dimensional ultrasound in the small fetus. *Ultrasound Obstetr. Gynecol.* **32**(2), 168–175 (2008). <https://doi.org/10.1002/uog.6111>
- Schmid, M., Hothorn, T.: Flexible boosting of accelerated failure time models. *BMC Bioinform.* **9**, 269 (2008). <https://doi.org/10.1186/1471-2105-9-269>
- Schmid, M., Hothorn, T., Maloney, K.O., Weller, D.E., Potapov, S.: Geoadditive regression modeling of stream biological condition. *Environ. Ecol. Stat.* **18**(4), 709–733 (2011). <https://doi.org/10.1007/s10651-010-0158-4>
- Seibold, S., Brandl, R., Schmidl, J., Busse, J., Thorn, S., Hothorn, T., Müller, J.: Extinction risk status of saproxylic beetles reflects the ecological degradation of forests in Europe. *Conserv. Biol.* **29**(2), 382–390 (2015). <https://doi.org/10.1111/cobi.12427>
- Seibold, H., Zeileis, A., Hothorn, T.: Individual treatment effect prediction for ALS patients. *Stat. Methods Med. Res.* (2017). <https://doi.org/10.1177/0962280217693034>
- Wu, C.O., Tian, X.: Nonparametric estimation of conditional distributions and rank-tracking probabilities with time-varying transformation models in longitudinal studies. *J. Am. Stat. Assoc.* **108**(503), 971–982 (2013). <https://doi.org/10.1080/01621459.2013.808949>
- Yang, Y., Tokdar, S.T.: Joint estimation of quantile planes over arbitrary predictor spaces. *J. Am. Stat. Assoc.* **112**(519), 1107–1120 (2017). <https://doi.org/10.1080/01621459.2016.1192545>
- Yue, M., Li, J., Ma, S.: Sparse boosting for high-dimensional survival data with varying coefficients. *Stat. Med.* **37**(5), 789–800 (2017). <https://doi.org/10.1002/sim.7544>

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.