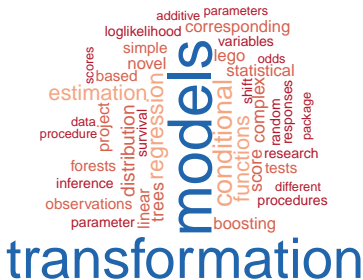




Transformation Models: Pushing the Boundaries

Torsten Hothorn



On the Menu Today

- What are transformation models?
- What are they capable of?
- What is their connection to trees and forests?

Transformation Models

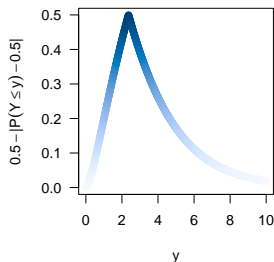
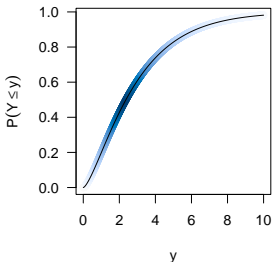
Models for Distributions, not Means

Regression:

$$\mathbb{P}(Y \leq y \mid \mathbf{X} = \mathbf{x})$$

not only

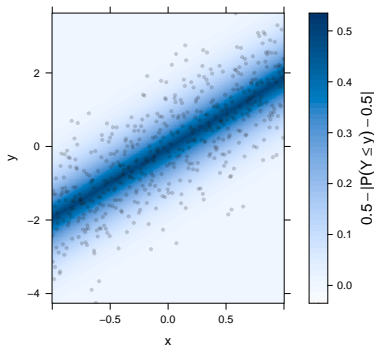
$$\mathbb{E}(Y \mid \mathbf{X} = \mathbf{x})$$



The Normal Linear Regression Model

$$Y = \alpha + \mathbf{x}^\top \boldsymbol{\gamma} + \sigma Z, \quad Z \sim N(0, 1)$$

- everything but “normal”
- most special case
- $Y | \mathbf{x} \sim N(\alpha + \mathbf{x}^\top \boldsymbol{\gamma}, \sigma^2)$
- no way escaping normal land



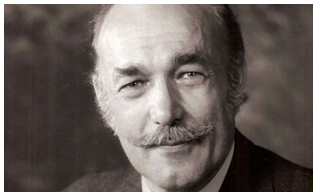
Generalisation I

$$\begin{aligned} Y &= \alpha + \mathbf{x}^\top \boldsymbol{\gamma} + \sigma Z, & Z &\sim N(0, 1) \\ \iff \mathbb{E}(Y \mid \mathbf{X} = \mathbf{x}) &= \alpha + \mathbf{x}^\top \boldsymbol{\gamma} & , & Y \mid \mathbf{X} = \mathbf{x} \sim N(,) \\ \hookrightarrow g(\mathbb{E}(Y \mid \mathbf{X} = \mathbf{x})) &= \alpha + \mathbf{x}^\top \boldsymbol{\gamma} & , & Y \mid \mathbf{X} = \mathbf{x} \sim \text{ExpFam}(,) \end{aligned}$$

Generalized Linear Models

By J. A. NELDER and R. W. M. WEDDERBURN

Rothamsted Experimental Station, Harpenden, Herts



Generalisation II

$$\begin{aligned} Y &= \alpha + \mathbf{x}^\top \boldsymbol{\gamma} + \sigma Z, \quad Z \sim N(0, 1) \\ \iff \frac{Y - \alpha}{\sigma} &= \mathbf{x}^\top \frac{\boldsymbol{\gamma}}{\sigma} + Z, \quad Z \sim N(0, 1) \\ \iff \frac{Y - \alpha}{\sigma} &= \mathbf{x}^\top \boldsymbol{\beta} + Z, \quad Z \sim N(0, 1) \\ \hookrightarrow h(Y) &= \mathbf{x}^\top \boldsymbol{\beta} + Z, \quad Z \sim N(0, 1) \\ \hookrightarrow h(Y) &= \mathbf{x}^\top \boldsymbol{\beta} + Z, \quad Z \sim \end{aligned}$$

Transformation models, $Z \in \mathbb{R}$ with absolute continuous log-concave density f_Z , $h: \mathbb{R} \rightarrow \mathbb{R}$ nondecreasing

An Analysis of Transformations (1964)



An Analysis of Transformations

By G. E. P. Box and D. R. COX
University of Wisconsin *Birkbeck College, University of London*

[Read at a RESEARCH METHODS MEETING of the SOCIETY, April 8th, 1964,
Professor D. V. LINDLEY in the Chair]

Conceptually more powerful but, at the time, hard to compute and thus restricted to

$$h(y | \lambda) = \begin{cases} \frac{y^{\lambda+1}}{\lambda} & \lambda > 0 \\ \log(y) & \lambda = 0 \end{cases} \quad \text{with } Z \sim N(0, \sigma^2)$$

“Box-Cox” power transformation

Conditional Distribution Functions

$$h(Y) = \mathbf{x}^\top \boldsymbol{\beta} + Z$$
$$\Leftrightarrow \mathbb{P}(Y \leq y \mid \mathbf{X} = \mathbf{x}) = F_Z(h(y) - \mathbf{x}^\top \boldsymbol{\beta})$$

also allows discrete models via step-function h

Linear transformation models: Proportional hazards, proportional odds, ...

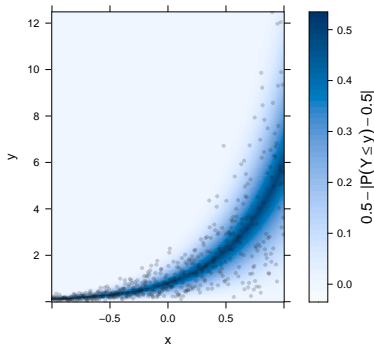
aka Probabilistic index models:

$$\mathbb{P}(Y_1 < Y_2 \mid \mathbf{x}_1, \mathbf{x}_2) = m_Z((\mathbf{x}_1 - \mathbf{x}_2)^\top \boldsymbol{\beta})$$

Cox Models

$$\mathbb{P}(Y \leq y \mid \mathbf{X} = \mathbf{x}) = 1 - \exp(-\exp(h(y) + \mathbf{x}^\top \boldsymbol{\beta}))$$

- the most prominent transformation model
- h is baseline log-cumulative hazard
- $\mathbf{x}^\top \boldsymbol{\beta}$ is log-hazard ratio
- partial likelihood profiles out h
- np score test: log-rank
- no censoring in *model*, only in *likelihood*

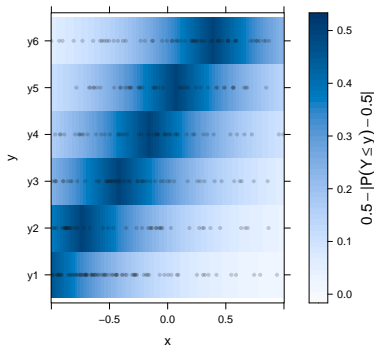


Proportional-odds Models

Ordinal outcome at categories $y_1 < y_2 < \dots < y_K$

$h(y_k) = \vartheta_k, k = 1, \dots, K - 1$ with $F_Z = \text{expit}$

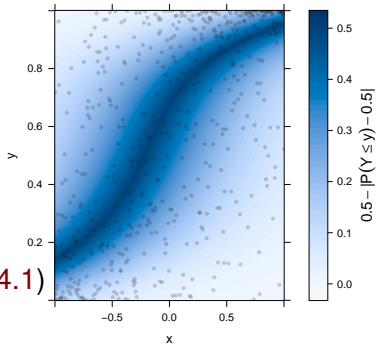
- Proportional-odds model
- Popular for ordinal data analysis
- $\mathbf{x}^\top \boldsymbol{\beta}$ is log-odds ratio
- *Simultaneous* estimation of h (via ϑ_k) and $\boldsymbol{\beta}$



Conditional Outcome Logistic Regression / Ordinal Regression Model

$$\mathbb{P}(Y \leq y \mid \mathbf{X} = \mathbf{x}) = \text{expit}(h(y) - \mathbf{x}^\top \beta)$$

- Continuous proportional-odds model
- h is baseline log-odds function
- $\mathbf{x}^\top \beta$ is log-odds ratio
- np score test: Wilcoxon
- parametric: Colr
([10.12688/f1000research.12934.1](#))
- nonparametric: orm
([10.1002/sim.7433](#))

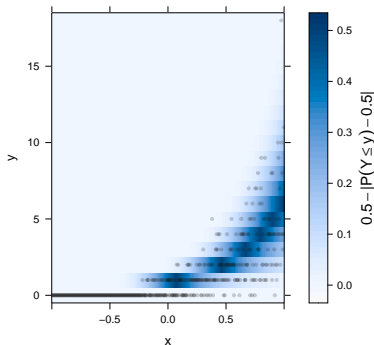


Count Transformation Models

Count outcome $Y \in \{0, 1, 2, \dots\}$

$h(k) := h(\lfloor y \rfloor) \quad \forall k \leq y < k + 1$ with $h : \mathbb{R} \rightarrow \mathbb{R}$

- Smooth $h(\cdot \mid \vartheta)$
evaluated discretely
- More flexible than
Poisson/NB
- Discrete count likelihood
- 10.1111/2041-
210X.13383



Distribution Regression

Quantile regression:

$$Q(\tau | \mathbf{X} = \mathbf{x}) = \alpha(\tau) + \mathbf{x}^\top \boldsymbol{\delta}(\tau)$$

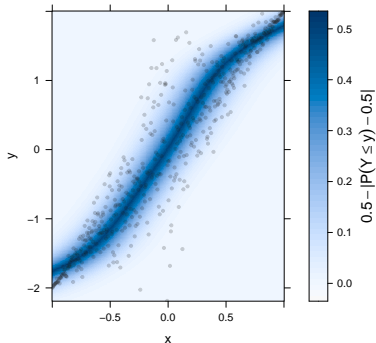
Distribution regression:

$$F_Z^{-1}(\mathbb{P}(Y \leq y | \mathbf{X} = \mathbf{x})) = h(y) - \mathbf{x}^\top \boldsymbol{\beta}(y)$$

i.e. quantile regression on the log-cumulative hazard / odds
or probit scale

Distribution Regression

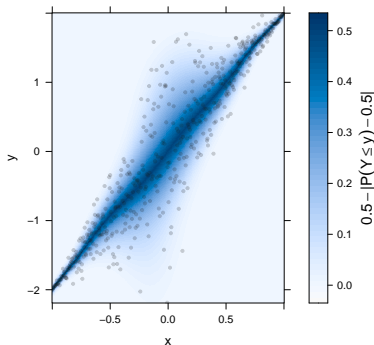
- Often (re-)discovered
- $\text{logit}(\mathbb{E}(I(Y \leq c))) = \alpha_c + \mathbf{x}^\top \beta_c$
- Full likelihood possible
- Splines for $h(y)$ and $\beta(y)$



Conditional Transformation Models

$$\mathbb{P}(Y \leq y \mid \mathbf{X} = \mathbf{x}) = F_Z(h(y \mid \mathbf{x}))$$

- Practically no assumptions
- Tensor-product spline bases for y and \mathbf{x}
- Full likelihood possible
- 10.1111/rssb.12017



Transformation Likelihood

Log-likelihoods

Observed $(\underline{y}, \bar{y}] \subset \mathbb{R}$:

$$\log[F_Z\{h(\bar{y} | \mathbf{x})\} - F_Z\{h(\underline{y} | \mathbf{x})\}]$$

This includes discrete and censored observations and, via $(Y_{(k)}, Y_{(k+1)}]$, the nonparametric likelihood.

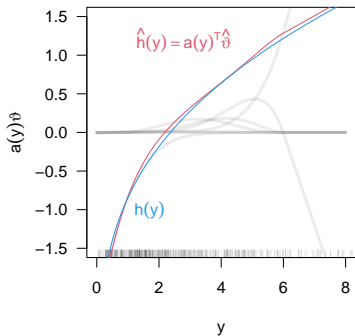
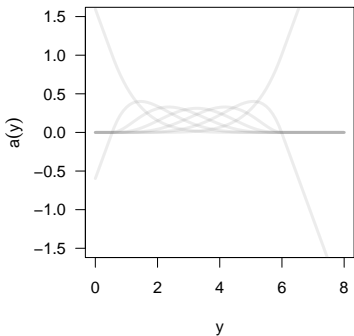
Observed $Y \in \mathbb{R}$:

$$\approx \log[f_Z\{h(y | \mathbf{x})\}] + \log\{h'(y | \mathbf{x})\}$$

10.1111/sjos.12291

Parameterisation

- h typically swept under the carpet
- More fun: parameterise $h(y | \vartheta) = \mathbf{a}(y)^\top \vartheta$ in terms of $\vartheta \in \Theta$ (a la [10.1080/15598608.2013.772835](#))
- Estimate all parameters *simultaneously* (via ML)
- Does it hurt? Not really: [10.1002/sim.8425](#)



Model Baking

- **Ingredients:** Take F_Z , parameterise $h(y | \vartheta)$, define impact of \mathbf{x} via $\mathbf{x}^\top \beta$, $\mathbf{x}^\top \beta(y)$, $h(y | \mathbf{x})$
- **Mix:** Data defines likelihood function (handles discrete and continuous observations, censoring)
- **Oven:** Optimise, get $\hat{\vartheta}, \hat{\beta}$ + limiting distribution
- **Serve:** Interpret/interrogate fully specified model

And now for some new recipes...

Pushing the Boundaries

Correlated Observations

Mixed-effects transformation models via conditional distribution

$$\mathbb{P}(Y \leq y \mid \mathbf{X} = \mathbf{x}, \mathbf{U} = \mathbf{u}, \xi) = F_Z(h(y \mid \vartheta) - \mathbf{x}^\top \boldsymbol{\beta} - \mathbf{u}^\top \boldsymbol{\xi})$$

with normal random effects $\boldsymbol{\xi} \sim N_q(\mathbf{0}, \Sigma)$

Integrate conditional likelihood wrt random effects, using the fabulous **TMB package**

[10.32614/RJ-2021-075](https://doi.org/10.32614/RJ-2021-075) [10.1093/biostatistics/kxab045](https://doi.org/10.1093/biostatistics/kxab045)

High Dimensions

$$\mathbb{P}(Y \leq y \mid \mathbf{X} = \mathbf{x}) = F_Z(h(y \mid \vartheta) - \mathbf{x}^\top \beta)$$

with $\beta \in \mathbb{R}^p$, with p being large

Add L_1, L_2, \dots penalty for β to likelihood, using the fabulous **CVXR package**

Transformation ridge, lasso, elastic net,
etc. **10.32614/RJ-2021-054**

Additive Models

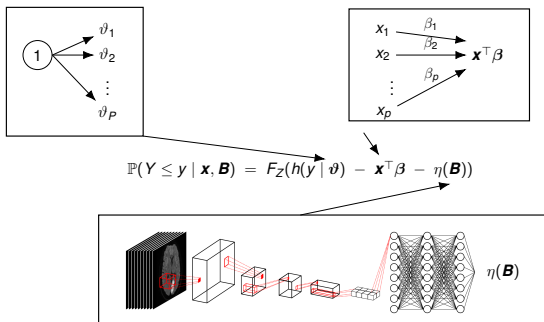
$$\mathbb{P}(Y \leq y \mid \mathbf{X} = \mathbf{x}) = F_Z \left(h(y \mid \vartheta) - \sum_{j=1}^J f_j(\mathbf{x}) \right)$$

aka “transform-both-sides”

Use connection to mixed models and leverage **mgcv** infrastructure; in the making

Alternative: Boosting via **mboost**,
[10.1007/s11222-019-09870-4](https://doi.org/10.1007/s11222-019-09870-4)

Unstructured Information



“Deep” transformations, [10.1016/j.patcog.2021.108263](https://arxiv.org/abs/10.1016/j.patcog.2021.108263)
[10.1007/978-3-030-86523-8_1](https://arxiv.org/abs/10.1007/978-3-030-86523-8_1)

Multivariate Transformation Models

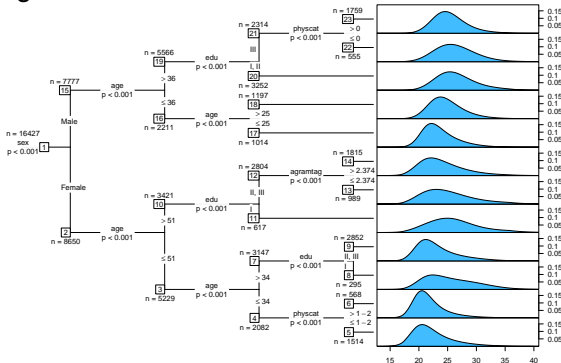
$$\mathbb{P} \left(\bigcap_{j=1}^J Y_j \leq y_j \mid \mathbf{X} = \mathbf{x} \right) = F_{\mathbf{Z}}(h_1(y_1 \mid \mathbf{x}) + h_2(y_2 \mid \mathbf{x}) + \lambda_{21}(\mathbf{x})h_1(y_1 \mid \mathbf{x}) + \dots + h_J(y_J \mid \mathbf{x}) + \sum_{j=1}^{J-1} \lambda_{Jj}(\mathbf{x})h_j(y_j \mid \mathbf{x}))$$

with $\mathbf{Z} \sim N_J(\mathbf{0}, \mathbf{I})$. Correlations depend on \mathbf{x} through $\lambda_{jj'}(\mathbf{x})$

Connection to Gaussian copulas and normalising flows,
[10.1111/sjos.12501](https://doi.org/10.1111/sjos.12501)

Trees and Forests

Model-based recursive partitioning (MOB) based on transformation models, trees and forests for distributional regression



10.1515/ijb-2019-0063, 10.1177/0962280219862586

Random Forest

What is a random forest (in general, not only B&C)?

Classical:

$$\mathbb{E}(Y | \mathbf{X} = \mathbf{x}) = f(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}$$

Here:

$$\mathbb{P}(Y \leq y | \mathbf{X} = \mathbf{x}) = \mathbb{P}_{Y|\mathbf{X}=\mathbf{x}}(y) = f(y | \mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}$$

RF often understood as being “nonparametric” but life is easier when we have a *model*.

Parametric (!) Setup

Unconditional model for response

$$\mathbb{P}_{Y,\Theta} = \{\mathbb{P}_{Y,\vartheta} \mid \vartheta \in \Theta\}$$

Assumption: Regression model belongs to this family:

$$\mathbb{P}_{Y|\mathbf{X}=\mathbf{x}} = \mathbb{P}_{Y,\vartheta(\mathbf{x})}$$

Task: Estimate ϑ function

Likelihood Contributions

“Learning” data $(y_i, \mathbf{x}_i), i = 1, \dots, N$ plus family $\mathbb{P}_{Y, \Theta}$ defines likelihood function

$$\ell_i : \Theta \rightarrow \mathbb{R}$$

$\ell_i(\vartheta(\mathbf{x}_i))$ gives the likelihood for observation i with candidate parameters $\vartheta(\mathbf{x}_i)$

Handle censoring and truncation appropriately here

Adaptive Local Likelihood Estimators

$$\hat{\vartheta}^N(\mathbf{x}) := \arg \max_{\vartheta \in \Theta} \sum_{i=1}^N w_i^N(\mathbf{x}) \ell_i(\vartheta)$$

Conditioning works via weight functions $w_i^N(\mathbf{x})$ only

Unconditional Maximum Likelihood

$$\hat{\boldsymbol{\vartheta}}_{\text{ML}}^N := \arg \max_{\boldsymbol{\vartheta} \in \Theta} \sum_{i=1}^N \ell_i(\boldsymbol{\vartheta})$$

Trees

$$\mathcal{X} = \dot{\bigcup}_{b=1, \dots, B} \mathcal{B}_b$$

$$w_{\text{Tree},i}^N(\mathbf{x}) := \sum_{b=1}^B I(\mathbf{x} \in \mathcal{B}_b \wedge \mathbf{x}_i \in \mathcal{B}_b)$$

$$\hat{\vartheta}_{\text{Tree}}^N(\mathbf{x}) := \arg \max_{\vartheta \in \Theta} \sum_{i=1}^N w_{\text{Tree},i}^N(\mathbf{x}) \ell_i(\vartheta)$$

Forests

$$\chi = \dot{\bigcup}_{b=1, \dots, B_t} \mathcal{B}_{tb} \text{ for } t = 1, \dots, T \text{ trees}$$

$$w_{\text{Forest},i}^N(\mathbf{x}) := \sum_{t=1}^T \sum_{b=1}^{B_t} I(\mathbf{x} \in \mathcal{B}_{tb} \wedge \mathbf{x}_i \in \mathcal{B}_{tb})$$

$$\hat{\vartheta}_{\text{Forest}}^N(\mathbf{x}) := \arg \max_{\vartheta \in \Theta} \sum_{i=1}^N w_{\text{Forest},i}^N(\mathbf{x}) \ell_i(\vartheta)$$

OK, Done! Really?

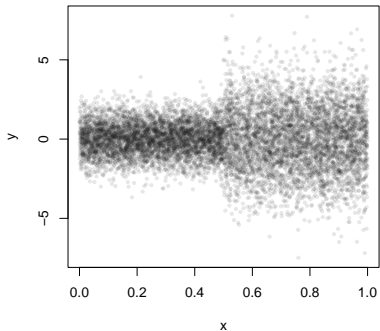
These “nearest neighbor weights” have been used before, first in

- “bagging survival trees” (2004), in
- “conditional inference forests” (**party(kit)**, since 2005) and in
- “quantile regression forests” (**quantregForest**, since 2006)

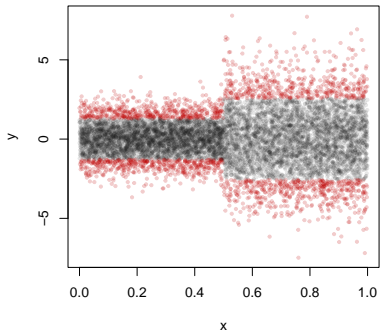
with *standard* trees (CART- or CTree-like).

Unfortunately, there is a catch.

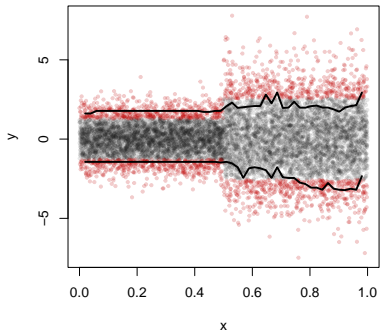
The Problem



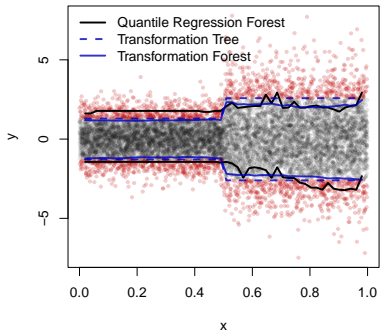
The Problem



The Problem



The Problem



The Solution

We need splits sensitive to *distributional* and not just *mean* changes.

Generic approach (“Distribution trees and forests”):

$$\mathbb{P}(Y \leq y \mid \mathbf{X} = \mathbf{x}) = \mathbb{P}_{Y, \vartheta(\mathbf{x})}(y)$$

Here: Use transformation model

$$\mathbb{P}(Y \leq y \mid \mathbf{X} = \mathbf{x}) = F_Z(\mathbf{a}(y)^\top \vartheta(\mathbf{x}))$$

10.1080/10618600.2021.1872581

Parameterisation

Transformation trees and forests based on parameterisation

$$\mathbb{P}(Y \leq y \mid \mathbf{X} = \mathbf{x}) = F_Z(\mathbf{a}_{\text{Bs},d}(y)^\top \boldsymbol{\vartheta}(\mathbf{x}))$$

- $\mathbf{a}_{\text{Bs},d}(y)^\top \boldsymbol{\vartheta}(\mathbf{x})$ is a smooth, monotonic Bernstein polynomial of degree d
- $d = 1$ with $F_Z = \Phi$ means $\mathbb{P}_{Y|\mathbf{X}=\mathbf{x}} = \mathcal{N}(\mu(\mathbf{x}), \sigma^2(\mathbf{x}))$
- $d = 5$ is surprisingly flexible

Model-based Recursive Partitioning (MOB)

Core idea

- Fit parameters $\hat{\vartheta}_{\text{ML}}$ in *unconditional* model $\mathbb{P}_{Y, \vartheta}$
- Compute individual gradient contributions (“scores”)

$$\mathbf{s}_i = \left. \frac{\partial \ell_i(\vartheta)}{\partial \vartheta} \right|_{\vartheta = \hat{\vartheta}_{\text{ML}}}$$

- Select predictor from \mathbf{x} with strongest parameter instability as indicated by highest association to $\mathbf{s}_i, i = 1, \dots, N$
- Find “best” binary split; repeat recursively

Implemented for many models, including (G)LM(M)s, parametric survival, β -regression, spatial lag, Bradley-Terry-Luce, various Item Response Theory models, subgroup analyses, etc.

Transformation Trees (TTree)

- Start with $\hat{\vartheta}_{ML}^N$
- Search for parameter instabilities in $\hat{\vartheta}_{ML}^N$ as a function of \mathbf{x} using (a beefed-up version) of MOB
- Potentially find changes in the mean AND higher moments
- Forests: Aggregate these trees via adaptive local likelihood estimation

Transformation Forests (TForest)

$$\hat{\mathbb{P}}(Y \leq y \mid \mathbf{X} = \mathbf{x}) = F_Z(\mathbf{a}_{\text{Bs},d}(y))^\top \hat{\boldsymbol{\vartheta}}_{\text{Forest}}^N(\mathbf{x})$$

makes the forest “parametric” (one model for each \mathbf{x}) with

- Forest likelihood
- Prediction intervals
- Likelihood-based variable importance
- Parametric bootstrap
- ...

and applicable to censored and truncated data.

More Complex Models

For example: Subgroup analysis, stratified / personalised medicine, ...

Conditional transformation model

$$\mathbb{P}(Y \leq y \mid \text{treatment}, \mathbf{X} = \mathbf{x}) = F_Z(\mathbf{a}_{\text{BS},d}(y)^\top \boldsymbol{\vartheta}(\mathbf{x}) - \beta(\mathbf{x})I(\text{treated}))$$

- Both the “intercept function” $\mathbf{a}_{\text{BS},d}(y)^\top \boldsymbol{\vartheta}(\mathbf{x})$ and
- the treatment effect $\beta(\mathbf{x})$ may depend on \mathbf{x}
- $F_Z() = 1 - \exp(-\exp())$ makes β a log-hazard ratio
- Include $\hat{\beta}$ in search for parameter instabilities

Stratified Medicine

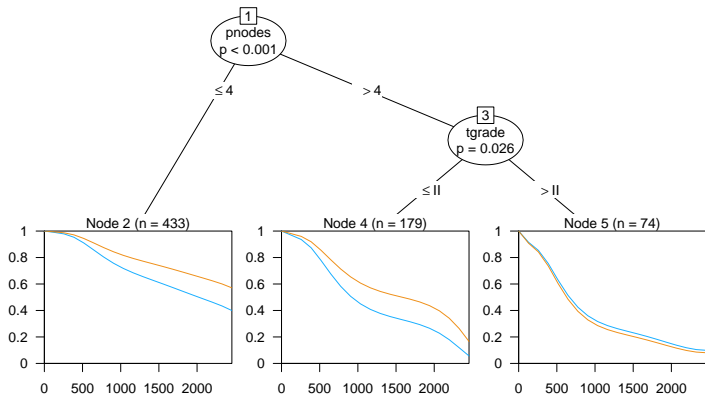
Partition log-hazard ratio β from a fully parametric Cox model

$$\mathbb{P}(T > t \mid \text{treatment}) = \exp(-\exp(\mathbf{a}_{\text{BS},d}(t)^\top \boldsymbol{\vartheta} - \beta I(\text{treated})))$$

for a randomised controlled clinical trial on hormonal treatment of breast-cancer patients

```
library("tram")
cmod <- Coxph(ctime ~ horTh, data = GBSG2)
library("trtf")
tmod <- trafotree(cmod,
                  formula = ctime ~ horTh | .,
                  data = GBSG2)
```

Stratified Medicine



Survival Forests

Log-rank splitting implicitly assumes proportional hazards model

$$\mathbb{P}(T > t \mid \mathbf{X} = \mathbf{x}) = \exp(-\exp(h(y) - \beta(\mathbf{x})))$$

⇒ `cforest`, `ranger`, `randomForestSRF` are insensitive to non-proportional hazards effects.

Switching to transformation forests based on

$$\mathbb{P}(T > t \mid \mathbf{X} = \mathbf{x}) = \exp(-\exp(\mathbf{a}(y)^\top \boldsymbol{\vartheta}(\mathbf{x})))$$

relaxes this restriction, [10.1177/0962280219862586](https://doi.org/10.1177/0962280219862586).

Ordinal Transformation Forests

Proportional-odds models assume

$$\mathbb{P}(Y \leq y_k \mid \mathbf{X} = \mathbf{x}) = \text{expit}(\vartheta_k - \beta(\mathbf{x}))$$

Transformation models can be set-up as estimators for $\beta(\mathbf{x})$ or, without assuming proportional odds, for more general models

$$\mathbb{P}(Y \leq y_k \mid \mathbf{X} = \mathbf{x}) = \text{expit}(\vartheta_k(\mathbf{x}))$$

10.1089/neu.2020.7407

R Add-on Packages

- **mlt**: Basic infrastructure
- **tram**: Model interfaces, multivariate models
- **cotram**: Count models
- **tramnet**: Penalisation
- **tramvs**: Variable selection
- **tramME**: Mixed-effects
- **tbm**: Boosting
- **trtf**: Trees and Forests

A word cloud centered around the word "models" in a large, bold, blue font. Below it, the word "transformation" is written in a large, blue, sans-serif font. The word cloud consists of various terms in shades of brown and orange, including: additive, parameters, loglikelihood, corresponding, variables, simple, novel, based, scores, estimation, project, data, procedure, forests, inference, observations, parameter, linear, trees, regression, survival, distribution, conditional, functions, shift, complex, random, responses, package, research, tests, different, procedures, boosting, odds, statistical, and lego.

models

transformation

<https://ctm.R-forge.R-project.org/>