

Heteroskedastic BART and time-to-event outcomes

Rodney Sparapani

Associate Professor of Biostatistics
Medical College of Wisconsin

September 17, 2025

*This research funded, in part, by the US Office of Naval
Research*

Heteroskedastic BART (HBART)

Pratola, Chipman, George & McCulloch 2020 JCGS

$$y_i = f(x_i) + s(x_i)\epsilon_i \quad \epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$

$$f \stackrel{\text{prior}}{\sim} \text{BART}(H, \mu, \kappa, \tau, \alpha, \beta)$$

$$s^2 \stackrel{\text{prior}}{\sim} \text{HBART}(\tilde{H}, \tilde{\lambda}, \tilde{\nu}, \tilde{\alpha}, \tilde{\beta})$$

$$s^2(x_i) \equiv \prod_{h=1}^{\tilde{H}} g(x_i; \tilde{\mathcal{T}}_h, \tilde{\mathcal{M}}_h) \quad \tilde{H} \approx H/5$$

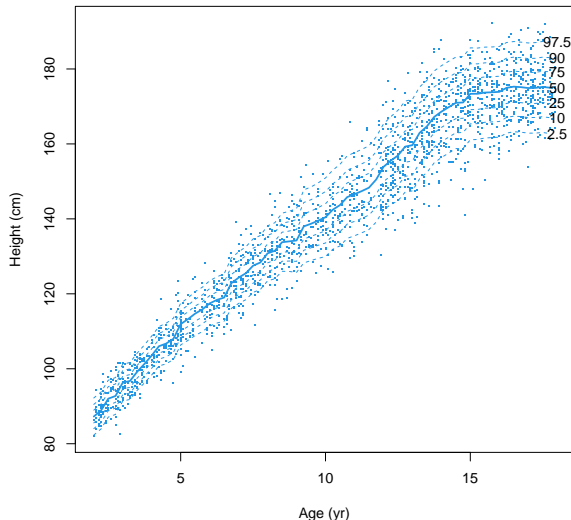
$$\sigma_{hl}^2 | \tilde{\mathcal{T}}_h \stackrel{\text{prior}}{\sim} \lambda \nu \chi^{-2}(\nu) \text{ leaves of } \tilde{\mathcal{T}}_h \quad \lambda = \tilde{\lambda}^{1/\tilde{H}}$$

$$\stackrel{\text{prior}}{\sim} \text{Gamma}^{-1}(\nu/2, \lambda\nu/2) \quad \mathbb{E}[\sigma_{hl}^2] = \lambda\nu/(\nu-2)$$

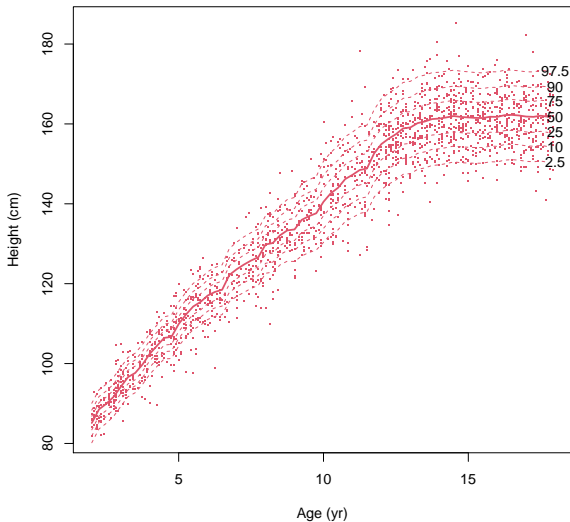
$$\in \tilde{\mathcal{M}}_h$$

$$\nu = 2 \left[1 - \left(1 - \frac{2}{\tilde{\nu}} \right)^{1/\tilde{H}} \right]^{-1}$$

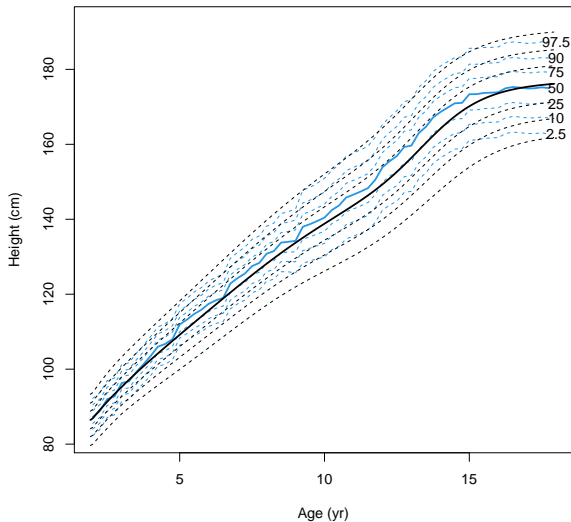
Marginal effect of age: HBART predictions for **M**
 $H = 300, \tilde{H} = 60, \text{numcut} = 200$



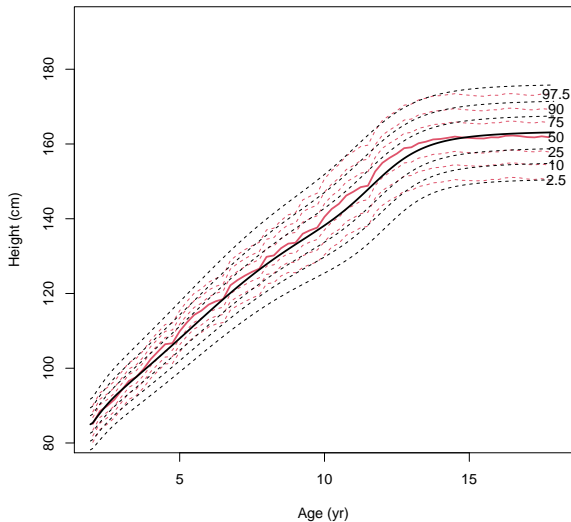
Marginal effect of age: HBART predictions for **F**



Marginal effect of age: HBART vs. CDC for M



Marginal effect of age: HBART vs. CDC for F



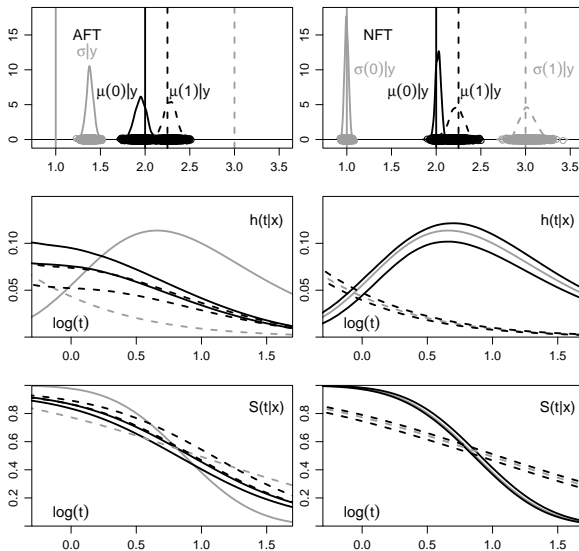
Personalized Hematopoietic Stem Cell Transplant (HSCT)

- HSCT is a treatment for white blood/bone marrow cancers
- Here we are concerned with unrelated donors that are human leukocyte antigen (HLA) 8/8 matched to the recipients transplanted from 2016:2019
- Goal: optimal donor matching for better recipient outcomes
- The outcome here is time to an event, i.e., event-free survival with both right and left censoring
- Events include death, relapse, graft failure/rejection or moderate/severe chronic graft vs. host disease (GVHD): whichever comes first
- There are $P = 45$ covariates that may have an impact
- 5 are donor-related characteristics: age, sex/childbearing, HLA DPB1 match, HLA DQB1 match and CMV match
- We wanted to *learn* the (likely complex) functional relationship between these covariates and the outcome with BART
- The cohort has 10016 for training and 1802 for validation
- A bit too large for our Discrete Time BART approach
- For this application, we developed NFT BART methodology

Methodological/Computational Pros and Cons

	Comparison of BART survival analysis methods				
Property	Hierarchical	Discrete Time	AFT	Modulated	NFT
Flexible assumptions	Con	Pro	Con	Pro	Pro
Non-parametric	Con	Pro	Pro	Pro	Pro
Left-censoring	Con	Con	Pro	Con	Pro
Time-dep. covariates	Con	Pro	Con	Pro	Con
Friendly to compute	Pro	Con	Pro	Con	Pro
First-author Year	Bonato 2011	Sparapani 2016	Henderson 2018	Linero 2021	Sparapani 2023

Two groups: **AFT BART** vs. **NFT BART**



Bayesian Additive Regression Trees (BART)

NFT notation

Sparapani, Logan, Laud & McCulloch 2023 *Biometrics*

$$\mu \stackrel{\text{prior}}{\sim} \text{BART} (a = 0.95, b = 2, H = 200, \kappa = 2, \tilde{\mu} = \bar{y})$$

$$y_i = \mu(x_i) + \epsilon_i \text{ where } \epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$

$$\mu(x_i) \equiv \tilde{\mu} + \sum_h g(x_i; \mathcal{T}_h, \mathcal{M}_h)$$

the **BART** prior implies the following priors (among others)

$$\mu_{hl} | \mathcal{T}_h \stackrel{\text{prior}}{\sim} N\left(0, \frac{0.25 \text{ range}(y)^2}{H \kappa^2}\right) \text{ leaves of } \mathcal{T}_h$$

$$\in \mathcal{M}_h$$

$$\sigma^2 \stackrel{\text{prior}}{\sim} \lambda \nu \chi^{-2}(\nu)$$

Heteroskedastic BART (HBART)

NFT notation

Pratola, Chipman, George & McCulloch 2019 *JCGS*

$$\mu \stackrel{\text{prior}}{\sim} \text{BART} (a, b, H = 200, \kappa = 5, \tilde{\mu})$$

$$\sigma^2 \stackrel{\text{prior}}{\sim} \text{HBART} (\tilde{a} = 0.95, \tilde{b} = 2, \tilde{H} = 40, \tilde{\lambda}, \tilde{\nu})$$

$$y_i = \mu(x_i) + \epsilon_i \text{ where } \epsilon_i \stackrel{\text{iid}}{\sim} \text{N}(0, \sigma^2(x_i))$$

$$\sigma^2(x_i) \equiv \prod_{h=1}^{\tilde{H}} g(x_i; \tilde{\mathcal{T}}_h, \tilde{\mathcal{M}}_h) \text{ where } \tilde{H} \approx H/5$$

the HBART prior implies the following priors (among others)

$$\begin{aligned} \sigma_{hl}^2 | \tilde{\mathcal{T}}_h &\stackrel{\text{prior}}{\sim} \lambda \nu \chi^{-2}(\nu) \text{ leaves of } \tilde{\mathcal{T}}_h \\ &\in \tilde{\mathcal{M}}_h \end{aligned}$$

The Accelerated Failure Time (AFT) model: part 1

- Time-to-event data notation: (t_i, δ_i) $i = 1, \dots, N$ subjects
if $\delta_i = 0$, then t_i is a right censoring time
else if $\delta_i = 1$, then a failure time
else if $\delta_i = 2$, then left censoring
- How is failure time explained by a vector of covariates x_i ?
- take logarithms $y_i = \log t_i$ and use a **linear model (Con)**
 $y_i = [1, x_i']\beta + \sigma\epsilon_i = \beta_0 + x_i'\beta_x + \sigma\epsilon_i$
where β and σ are unknown coefficients to be estimated
with $\epsilon_i \stackrel{\text{iid}}{\sim} F_\epsilon(\mu_\epsilon = 0, \sigma_\epsilon^2 = 1)$
which is typically **parametric (Con)**

The Accelerated Failure Time (AFT) model: part 2

- Consider a *baseline* survival function for a *standard* subject where the covariates are all zero, i.e., $S_0(t) = S(t|x=0)$.
- We can define the survival function for any given subject with a standard subject by accelerating, or decelerating, failure time

$$\begin{aligned} S(t|x_i) &= P[s_i > t|x_i] = P[y_i > \log t|x_i] \\ &= P[\beta_0 + x_i' \beta_x + \sigma \epsilon_i > \log t|x_i] \\ &= P[\beta_0 + \sigma \epsilon_i > \log t - x_i' \beta_x|x_i] \\ &= S_0(t \exp\{-x_i' \beta_x\}) \end{aligned}$$

- however, AFT is a precarious **restrictive assumption (Con)**
 $S(t|x) = P[\log s > \log t] = 1 - F_\epsilon(\log t; \mathbf{x}'\boldsymbol{\beta}, \sigma^2)$
the covariates can only explain a log-linear location shift

Survival analysis with AFT BART

NFT notation

Henderson, Louis et al. 2018 *Biostatistics*

- $y_i = \mu(x_i) + \epsilon_i$ where $\epsilon_i | \mu_i \sim N(\mu_i, \sigma^2)$: Pro $\mu^{\text{prior}} \sim \text{BART}$
- To ensure identifiability, constrain $\frac{1}{N} \sum_i \mu_i = 0$
- $\mu_i | G \sim G$
 $G | \alpha^{\text{prior}} \sim \text{DP}(\alpha, F_0)$
- $S(t, x) = 1 - \frac{1}{N} \sum_i \Phi\left(\frac{\log t - \mu_i - \mu(x)}{\sigma}\right)$

Con: the covariates still only explain a log-linear location shift

Survival analysis with NFT BART

Sparapani et al. 2023 *Biometrics*

- $y_i = \mu(x_i) + \epsilon_i$ where $\epsilon_i | (\mu_i, \sigma_i) \sim N(\mu_i, \sigma_i^2 \sigma^2(x_i))$: Pro

$\mu \stackrel{\text{prior}}{\sim} \text{BART}$

$\sigma^2 \stackrel{\text{prior}}{\sim} \text{HBART}$

- To ensure identifiability: $\frac{1}{N} \sum_i \mu_i = 0$ and $\frac{1}{N} \sum_i \sigma_i^2 = 1$
- if $\delta_i = 1$, then $y_i = \log t_i$
else draw

$$y_i \sim N(\mu_i + \mu(x_i), \sigma_i^2 \sigma^2(x_i)) \begin{cases} I(\log t_i, \infty) & \text{if } \delta_i = 0 \\ I(-\infty, \log t_i) & \text{if } \delta_i = 2 \end{cases}$$

- $(\mu_i, \sigma_i) | G \sim G$

$G | \alpha \stackrel{\text{prior}}{\sim} \text{DP}(\alpha, F_0)$

- $S(t, x) = 1 - \frac{1}{N} \sum_i \Phi\left(\frac{\log t - \mu_i - \mu(x)}{\sigma_i \sigma(x)}\right)$

Pro: the covariates can explain a location shift and rescaling!

Dirichlet Process Mixtures (DPM): infinite mixtures

Ferguson 1973 & Antoniak 1974 *Annals of Statistics*;

Escobar & West 1995 *JASA*; Neal 2000 *JCGS*

DPM-like finite mixture clustering: Miller & Harrison 2017 *JASA*

$$y_i | \theta_i \sim F(\theta_i) \quad \text{usual notation}$$

where $i = 1, \dots, N$

$$y_i | \theta_{c_i}^* \sim F(\theta_{c_i}^*) \quad \text{ephemeral clusters}$$

where $c_i \in \{1, \dots, k\}$ k is random

$$\theta_i | G \sim G \quad \text{nonparametric (Pro)}$$

$$G | \alpha \stackrel{\text{prior}}{\sim} \text{DP}(\alpha, F_0) \quad G \text{ "centered" on } F_0$$

$$\alpha \stackrel{\text{prior}}{\sim} \text{Gamma}(a, b) \quad \text{concentration parameter}$$

$$\propto k$$

$$\theta_1 \sim F_0 \quad \text{integrating over } G$$

$$\theta_2 | \theta_1 \sim \frac{1}{1 + \alpha} \delta_K(\theta_1) + \frac{\alpha}{1 + \alpha} F_0 \quad \text{mixture}$$

Constrained DPM

Yang, Dunson & Baird 2010

Computational Statistics & Data Analysis

- How do we constrain $\frac{1}{N} \sum_i \mu_i = 0$?
- Simply sample $(\tilde{\mu}_i, \tilde{\sigma}_i) | G \sim G$ as usual
Let $\tilde{\mu}_0 = \frac{1}{N} \sum_i \tilde{\mu}_i$
And $\mu_i = \tilde{\mu}_i - \tilde{\mu}_0$
- Similarly, if we need to constrain $\frac{1}{N} \sum_i \sigma_i^2 = 1$
Let $\tilde{\sigma}_0 = \sqrt{\frac{1}{N} \sum_i \tilde{\sigma}_i^2}$
And $\sigma_i = \tilde{\sigma}_i / \tilde{\sigma}_0$

Low Information Omnibus (LIO)

Dirichlet Process Mixtures prior hierarchy

Shi, Martens, Banerjee, Laud 2018 *Bayesian Analysis*

Sparapani et al. 2023 *Biometrics*

- With either DPM or Constrained DPM
- For convenience, re-parameterize in terms of $\tau_i = \sigma_i^{-2}$
 $F_0(\mu_0, \mathbf{k}_0, a_0, \mathbf{b}_0)$ is a Normal-Gamma prior
 $[\mu_i, \tau_i | \mathbf{k}_0, \mathbf{b}_0] = [\tau_i | \mathbf{b}_0] [\mu_i | \tau_i, \mathbf{k}_0]$
with $\mu_i | \tau_i, \mathbf{k}_0 \stackrel{\text{prior}}{\sim} N(\mu_0, (\tau_i \mathbf{k}_0)^{-1})$
and $\tau_i | \mathbf{b}_0 \stackrel{\text{prior}}{\sim} \text{Gamma}(a_0, \mathbf{b}_0)$
- NFT LIO prior parameter settings
 $\mu_0 = 0, \mathbf{k}_0 \stackrel{\text{prior}}{\sim} \text{Gamma}(1.5, 7.5)$
 $a_0 = 3, \mathbf{b}_0 \stackrel{\text{prior}}{\sim} \text{Gamma}(2, 1)$

NFT model: prediction intervals

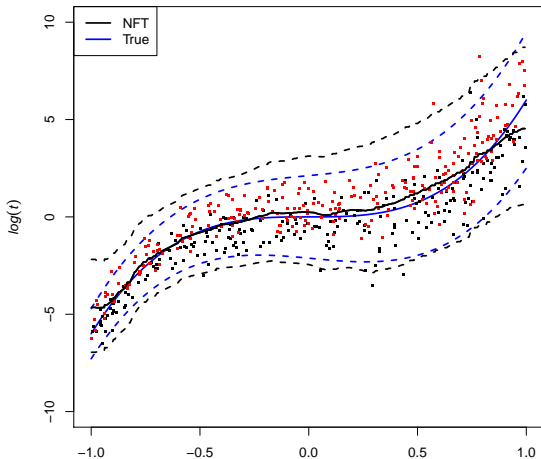
- $\log t_i = y_i = \mu(x_i) + \epsilon_i$ where $\epsilon_i \sim N(\mu_i, \sigma_i^2 \sigma^2(x_i))$
To ensure identifiability: $\frac{1}{N} \sum_i \mu_i = 0$ and $\frac{1}{N} \sum_i \sigma_i^2 = 1$
- $F_\epsilon = \frac{1}{N} \sum_i N(\mu_i, \sigma_i^2)$: nonparametric mixture of Normals
- $(1 - \alpha) \times 100\%$ Prediction Interval
 $(\mu(x) + c_{\alpha/2} \sigma(x), \mu(x) + c_{1-\alpha/2} \sigma(x))$
where $c_\pi = F_\epsilon^{-1}(\pi)$

NFT scenario $t(\mathbf{16})$: $N = 500$ with 50% censoring

$$f(x) = 6x^3, \quad s(x) = \exp 0.5x,$$

$$\log t = f(x) + s(x)\epsilon \text{ where } \epsilon \sim t(\mathbf{16})$$

and $x \sim U(-1,1)$: $R^2 = 84.8\%$ uncensored, $R^2 = 85.1\%$ censored

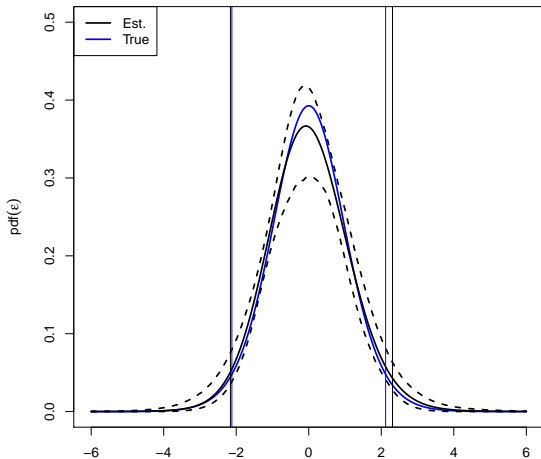


NFT scenario $t(\mathbf{16})$: $N = 500$ with 50% censoring

$$f(x) = 6x^3, \quad s(x) = \exp 0.5x,$$

$$\log t = f(x) + s(x)\epsilon \text{ where } \epsilon \sim t(\mathbf{16})$$

and $x \sim U(-1, 1)$: $R^2 = 84.8\%$ uncensored, $R^2 = 85.1\%$ censored

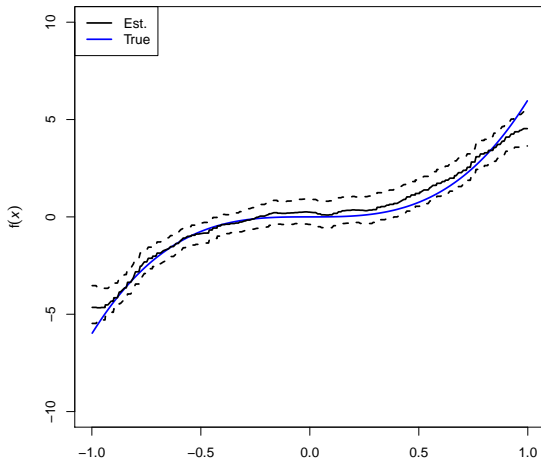


NFT scenario $t(16)$: $N = 500$ with 50% censoring

$$f(x) = 6x^3, \quad s(x) = \exp 0.5x,$$

$$\log t = f(x) + s(x)\epsilon \text{ where } \epsilon \sim t(16)$$

and $x \sim U(-1, 1)$: $R^2 = 84.8\%$ uncensored, $R^2 = 85.1\%$ censored

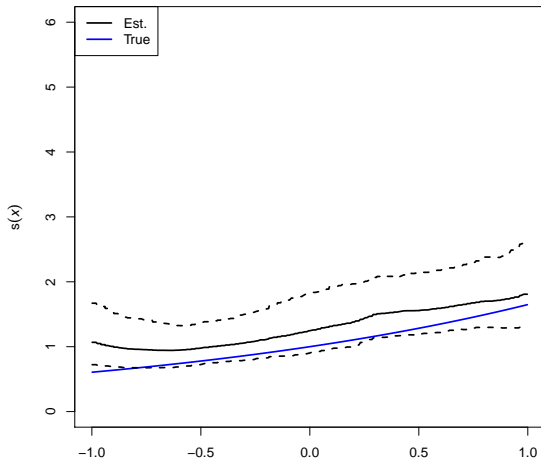


NFT scenario $t(\mathbf{16})$: $N = 500$ with 50% censoring

$$f(x) = 6x^3, \quad s(x) = \exp 0.5x,$$

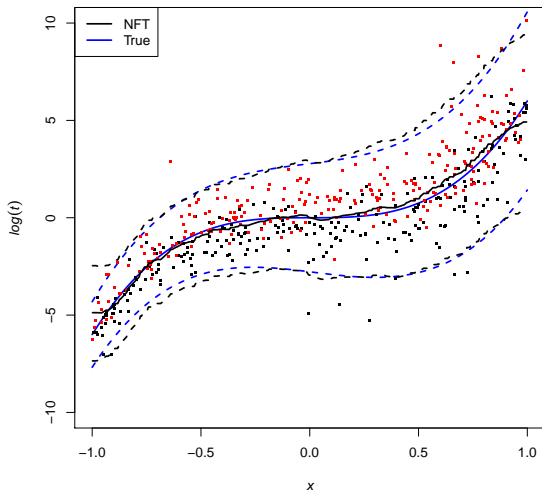
$$\log t = f(x) + s(x)\epsilon \text{ where } \epsilon \sim t(\mathbf{16})$$

and $x \sim U(-1,1)$: $R^2 = 84.8\%$ uncensored, $R^2 = 85.1\%$ censored



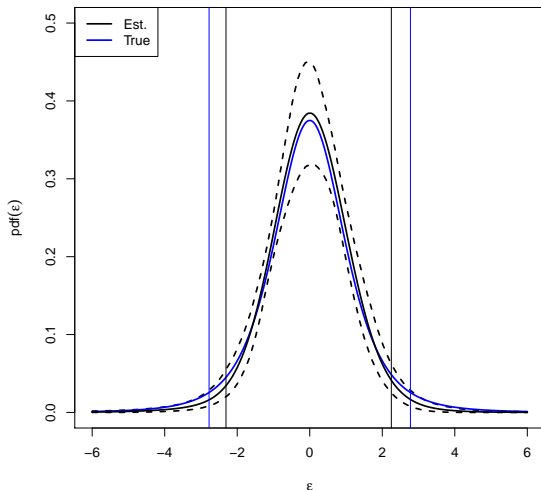
NFT scenario $t(4)$: $N = 500$ with 50% censoring

$f(x) = 6x^3$, $s(x) = \exp 0.5x$, $\log t = f(x) + s(x)\epsilon$ where $\epsilon \sim t(4)$
and $x \sim U(-1, 1)$: $R^2 = 80.7\%$ uncensored, $R^2 = 78.3\%$ censored



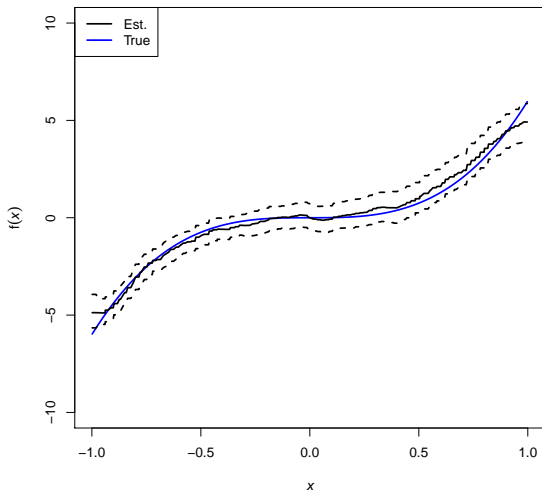
NFT scenario $t(4)$: $N = 500$ with 50% censoring

$f(x) = 6x^3$, $s(x) = \exp 0.5x$, $\log t = f(x) + s(x)\epsilon$ where $\epsilon \sim t(4)$
and $x \sim U(-1,1)$: $R^2 = 80.7\%$ uncensored, $R^2 = 78.3\%$ censored



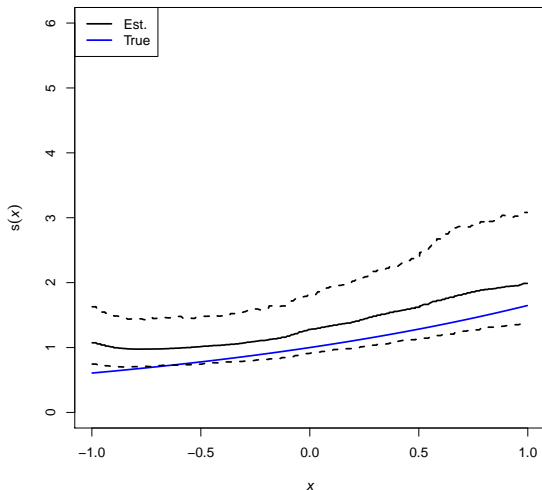
NFT scenario $t(4)$: $N = 500$ with 50% censoring

$f(x) = 6x^3$, $s(x) = \exp 0.5x$, $\log t = f(x) + s(x)\epsilon$ where $\epsilon \sim t(4)$
and $x \sim U(-1,1)$: $R^2 = 80.7\%$ uncensored, $R^2 = 78.3\%$ censored



NFT scenario $t(4)$: $N = 500$ with 50% censoring

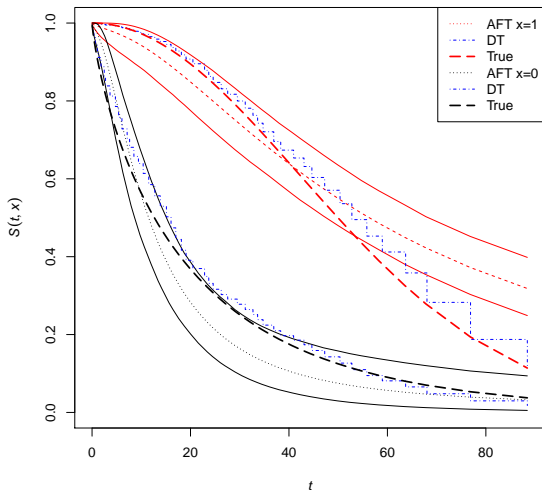
$f(x) = 6x^3$, $s(x) = \exp 0.5x$, $\log t = f(x) + s(x)\epsilon$ where $\epsilon \sim t(4)$
and $x \sim U(-1,1)$: $R^2 = 80.7\%$ uncensored, $R^2 = 78.3\%$ censored



Neither AFT nor NFT scenario: **AFT failure!**

$N = 500$ with 50% censoring

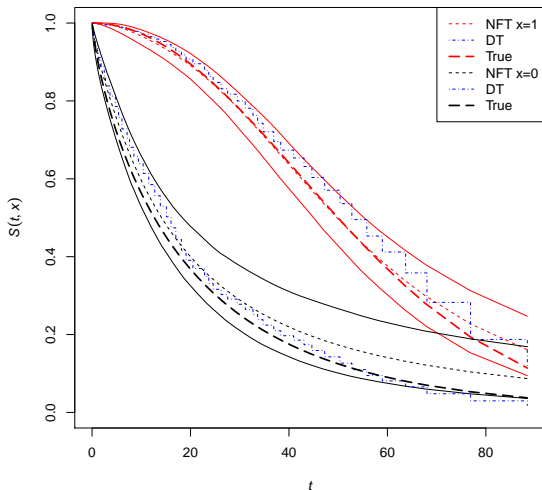
Wei $(0.8 + 1.2x, 20 + 40x)$ where $x \sim B(0.5)$



Neither AFT nor NFT scenario: **NFT success!**

$N = 500$ with 50% censoring

Wei $(0.8 + 1.2x, 20 + 40x)$ where $x \sim B(0.5)$



NFT BART posterior inference: the survival and distribution functions

$$S_m(t|x) = 1 - F_m(t|x)$$

$$\begin{aligned} F_m(t|x) &= \int \Phi \left\{ \frac{\log t - \mu_* - \mu_m(x)}{\sigma_* \sigma_m(x)} \right\} G_m(d\mu_*, d\sigma_*) \\ &= \sum_{j=1}^{\infty} \omega_j \Phi \left\{ \frac{\log t - \mu_j^* - \mu_m(x)}{\sigma_j^* \sigma_m(x)} \right\} \\ &\approx \sum_{j=1}^{K_m} \omega_{jm} \Phi \left\{ \frac{\log t - \mu_{jm}^* - \mu_m(x)}{\sigma_{jm}^* \sigma_m(x)} \right\} \end{aligned}$$

where $(\mu_{jm}^*, \sigma_{jm}^*)$ are from the training set

NFT BART posterior inference: the survival function

$$\widehat{S}(t|x) = M^{-1} \sum_m S_m(t|x)$$

$1 - 2\pi$ level credible intervals from π and $1 - \pi$ quantiles
($\widehat{S}_\pi(t|x), \widehat{S}_{1-\pi}(t|x)$) such that $\widehat{S}_p(t|x) = S_{m_p}(t|x)$
where m_p corresponds to the $p = \pi$ or $p = 1 - \pi$

NFT BART posterior inference: the hazard and density functions

$$h_m(t|x) = f_m(t|x)/S_m(t|x)$$

$$\begin{aligned} f_m(t|x) &= \int \frac{\phi\left\{\frac{\log t - \mu_* - \mu_m(x)}{\sigma_* \sigma_m(x)}\right\}}{t \sigma_* \sigma_m(x)} G_m(d\mu_*, d\sigma_*) \\ &= \sum_{j=1}^{\infty} \frac{\omega_j \phi\left\{\frac{\log t - \mu_j^* - \mu_m(x)}{\sigma_j^* \sigma_m(x)}\right\}}{t \sigma_j^* \sigma_m(x)} \\ &\approx \sum_{j=1}^{K_m} \frac{\omega_{jm} \phi\left\{\frac{\log t - \mu_{jm}^* - \mu_m(x)}{\sigma_{jm}^* \sigma_m(x)}\right\}}{t \sigma_{jm}^* \sigma_m(x)} \end{aligned}$$

where $(\mu_{jm}^*, \sigma_{jm}^*)$ are from the training set

NFT BART posterior inference: marginal effects by Friedman's partial dependence function

Friedman 2001 *Annals of Statistics*

- The covariates of interest are fixed at settings of interest: a single setting denoted x_A
- The complement take on the observed values found in the training data set denoted x_{iB} for subject i
- So the setting for all covariates denoted as (x_A, x_{iB})

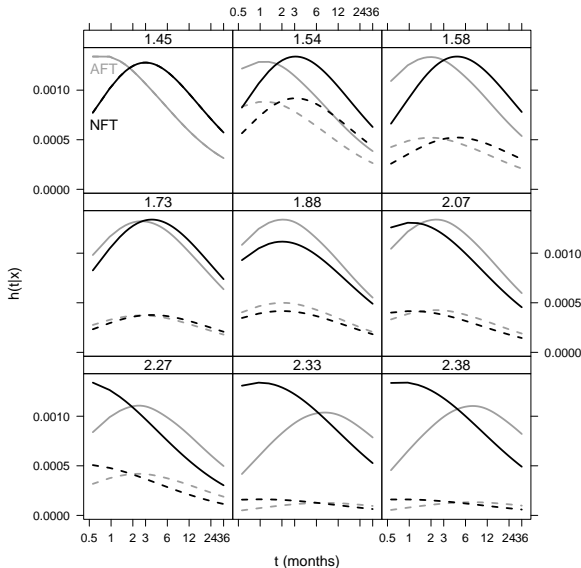
$$F_{Am}(t|x_A) = N^{-1} \sum_i \Phi \left(\frac{\log t - \mu_{im} - \mu_m(x_A, x_{iB})}{\sigma_{im} \sigma_m(x_A, x_{iB})} \right)$$

where (μ_{im}, σ_{im}) are from the training set

$$\widehat{S}_A(t|x_A) = 1 - M^{-1} \sum_m F_{Am}(t|x_A)$$

Real data example: **AFT BART** vs. **NFT BART**

%-iles of $\hat{\sigma}(x_i)$: 1, 5, 10, 30, 50, 70, 90, 95, 99



Thompson Sampling Variable Selection (TSVS)

Liu & Rockova 2023 *JASA*

Set \mathbf{H} small: 10, 20 or 40; smaller numbers engender more sparsity ($\tilde{\mathbf{H}} \approx \mathbf{H}/5$). TSVS is an iterative process: $k = 1, \dots, K$

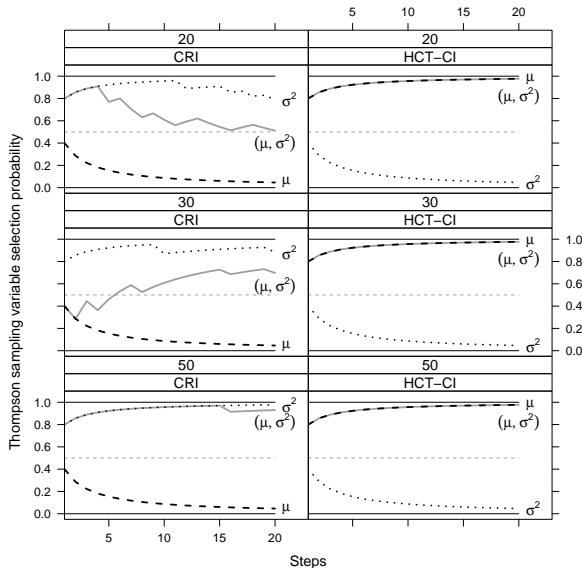
Pseudo-Bayesian prior parameter defaults: $a_{j0} = 1$ and $b_{j0} = 0.5$

- a. For $j = 1, \dots, P$: draw $\eta_{jk} \sim \text{Beta}(a_{j,k-1}, b_{j,k-1})$
- b. Set $\mathbf{B}_k = \{j : \eta_{jk} \geq 0.5\}$: covariate subset selected at step k
- c. Fit an NFT BART model with covariates x_{ij} where $j \in \mathbf{B}_k$
- d. For $j = 1, \dots, P$: do each sub-step
 - (i) Reward: if $j \notin \mathbf{B}_k$, then $\gamma_{jk} = 0$,
else $\gamma_{jk} = \mathbf{I}(u_{jkM} + v_{jkM} > 0)$ where u_{jkM} and v_{jkM}
are the number of branches for variable $x_{.j}$ in step k
from μ and σ^2 , respectively, at posterior draw M
 - (ii) Update via the reward: $a_{jk} = a_{j,k-1} + \gamma_{jk}$
and $b_{jk} = b_{j,k-1} + 1 - \gamma_{jk}$
 - (iii) Calculate inclusion probabilities: $\pi_{jk} = \frac{a_{jk}}{a_{jk} + b_{jk}}$
- e. If $k < K$, then return to a. and increment k

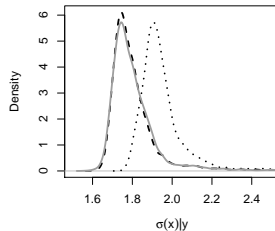
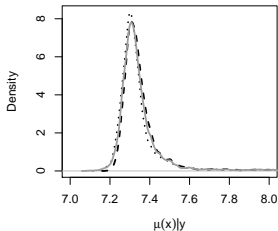
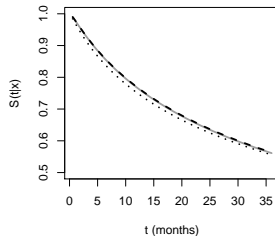
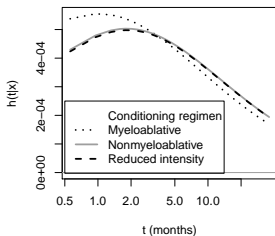
Important variables have trajectories of π_{jk} exceeding 0.5 by K

Real data example: TSVS with $H = 20, 30, 50$

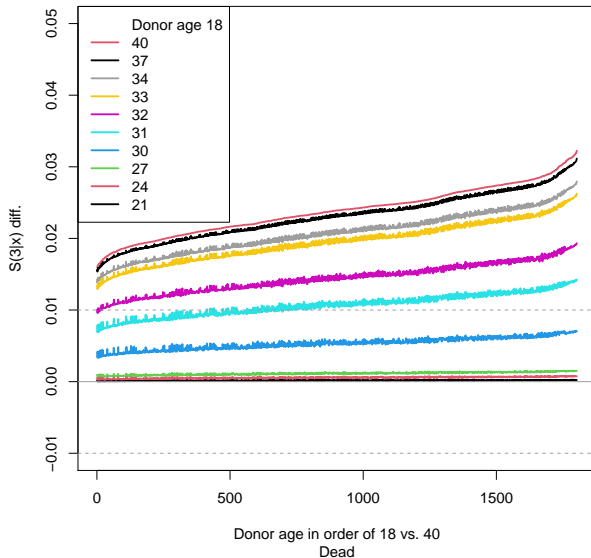
Conditioning regimen (CRI) and Comorbidity (HCT-CI)



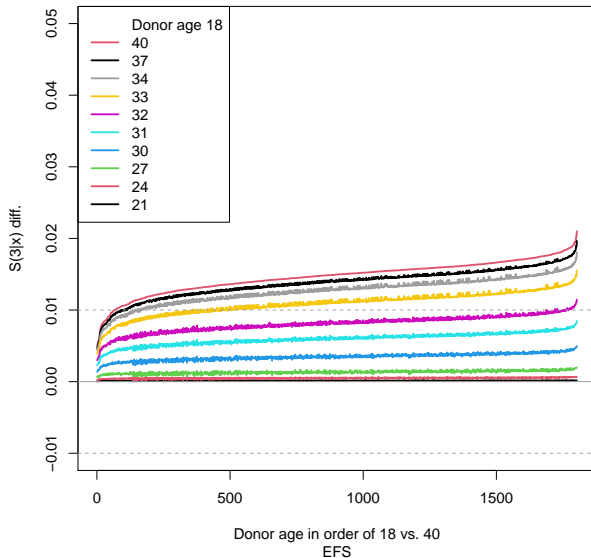
Real data example: Heteroskedasticity of the Conditioning Regimen Intensity (CRI)



Real data example: Death and donor age



Real data example: EFS and donor age



Conclusions: part 1

- We constructed our new **Nonparametric Failure Time (NFT)** approach from robust **Bayesian Nonparametric** building blocks
 - Bayesian Additive Regression Trees (BART) and Heteroskedastic BART (HBART)
 - Constrained Dirichlet Process Mixtures (DPM) with the Low Information Omnibus (LIO) prior hierarchy
- along with the **nftbart** v2.1 R package available on the Comprehensive R Archive Network (CRAN)

Conclusions: part 2

- NFT has desirable properties
 - **computationally friendly** via MCMC
 - **very flexible model** which does not resort to precarious restrictive assumptions
 - **default prior parameter settings** that work well without computationally expensive cross-validation
 - natural extensions to
 - variable selection** via Thompson Sampling and
 - marginal effects** by Friedman's partial dependence function
- Personalized Hematopoietic Stem Cell Transplant (HSCT)
 - For Event-free Survival of HSCT recipients
 - younger male** donors likely result in better outcomes