Nonparametric Failure Time: Time-to-event Machine Learning with Heteroskedastic Bayesian Additive Regression Trees and Low Information Omnibus Dirichlet Process Mixtures

Rodney Sparapani, Brent Logan, Prakash Laud and Rob McCulloch (Biometrics 2023)

RS, BL, PL and RM supported, in part, by the US Office of Naval Research Award Number: N00014-18-1-2888

Outline

- Motivation: a clinical application in Personalized Hematopoietic Stem Cell Transplant (HSCT) requires a new time-to-event BART methodology that scales better
- Pros and Cons of BART survival analysis methods
- ▶ BART and Heteroskedastic BART (HBART)
- Accelerated Failure Time (AFT) and AFT BART
- ► Nonparametric Failure Time (NFT) BART
- Dirichlet Process Mixtures (DPM), Constrained DPM and the Low Information Omnibus (LIO) DPM prior hierarchy
- Simulated data sets and methodology comparisons
- nftbart v1.6 R package on CRAN: nftbart v1.7 on my github
- ▶ Growth charts re-visited in nftbart v1.7: demo/bmx.R

Personalized Hematopoietic Stem Cell Transplant (HSCT)

- ► HSCT is a standard treatment for 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
- ► For this application, we developed NFT BART methodology

Methodological and Computational Pros and Cons

	Published BART survival analysis methods				
	Hier.	Discrete	AFT	Mod.	NFT
Property		Time (DT)			
Restrictive	Con	Pro	Con	Pro	Pro
Assump.					
Nonparam.	Con	Pro	Pro	Pro	Pro
Left-censor	Con	Con	Pro	Con	Pro
Comp.	Pro	Con	Pro	Con	Pro
Complexity					
First-author	Bonato	Sparapani	Henderson	Linero	Sparapani
Year	2011	2016	2018	2021	2023

Bayesian Additive Regression Trees (BART) NFT notation

Sparapani, Logan, Laud & McCulloch 2023 Biometrics

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

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

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

Heteroskedastic BART (HBART) NFT notation

Pratola, Chipman, George & McCulloch 2019 JCGS

$$y_{i} = \mu(x_{i}) + \epsilon_{i} \text{ where } \epsilon_{i} \stackrel{\text{iid}}{\sim} N(0, \sigma^{2}(x_{i}))$$
$$\mu \stackrel{\text{prior}}{\sim} \text{BART} (a, b, H = 200, \kappa = 5, \tilde{\mu})$$
$$\sigma^{2} \stackrel{\text{prior}}{\sim} \text{HBART} (\tilde{a} = 0.95, \tilde{b} = 2, \tilde{H} = 40, \tilde{\lambda}, \tilde{\nu})$$
$$\sigma^{2}(x_{i}) \equiv \prod_{h=1}^{\tilde{H}} g(x_{i}; \tilde{\mathcal{T}}_{h}, \tilde{\mathcal{M}}_{h}) \text{ where } \tilde{H} \approx H/5$$

The Accelerated Failure Time (AFT) model: part 1

- Time-to-event data notation: (t_i, δ_i) i = 1,..., N subjects if δ_i = 0, then t_i is a right censoring time else if δ_i = 1, then a failure time else if δ_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 centered, 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

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

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

Survival analysis with AFT BART NFT notation

Henderson et al. 2018 Biostatistics

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

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

Survival analysis with NFT BART

- ► $y_i = \mu(x_i) + \epsilon_i$ where $\epsilon_i | (\mu_i, \sigma_i) \sim N(\mu_i, \sigma_i^2 \sigma^2(x_i))$: Pro $\mu \sim^{\text{prior}} BART$ $\sigma^2 \sim^{\text{prior}} HBART$
- 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

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

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

►
$$S(t,x) = 1 - \frac{1}{N} \sum_{i} \Phi\left(\frac{\log t - \mu_{i} - \mu(x)}{\sigma_{i}\sigma(x)}\right)$$

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

Dirichlet Process Mixtures (DPM)

Ferguson 1973 & Antoniak 1974 *Annals of Statistics*; Escobar & West 1995, Geng et al. 2018 *JASA*; Neal 2000 *JCGS*

 $|\mathbf{y}_i|_{\boldsymbol{\theta}_i} \sim F(\boldsymbol{\theta}_i)$ usual notation where $i = 1, \dots, N$ $y_i | \theta_{c_i}^* \sim F(\theta_{c_i}^*)$ ephemeral random clusters where $c_i \in \{1, ..., k\}$ $k \in \{1, ..., N\}$ k is random parametric (Con): $\theta_i \stackrel{\text{prior}}{\sim} F_0$ $\theta_i | G \sim G$ nonparametric (Pro) $G \mid \alpha \stackrel{\text{prior}}{\sim} DP(\alpha, F_0)$ G "centered" on F_0 $\alpha \stackrel{\text{prior}}{\sim} \text{Gamma}(a, b)$ concentration parameter $\propto k$ $\theta_1 \sim F_0$ integrating over G $\theta_2|\theta_1 \sim \frac{1}{1+\alpha}\delta_K(\theta_1) + \frac{\alpha}{1+\alpha}F_0$ mixture

Constrained DPM

Yang et al. 2010 Computational Statistics and 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)

prior hierarchy

Shi, Martens, Banerjee, Laud 2018 Bayesian Analysis

- ► With either DPM or Constrained DPM
- For convenience, re-parameterize in terms of $\tau_i = \sigma_i^{-2}$ $F_0(\mu_0, k_0, a_0, b_0)$ is NoGa $[\mu_i, \tau_i | \mathbf{k}_0, \mathbf{b}_0] = [\tau_i | \mathbf{b}_0] [\mu_i | \tau_i, \mathbf{k}_0]$ with $\tau_i | b_0 \stackrel{\text{prior}}{\sim} \text{Gamma}(a_0, b_0)$ and $\mu_i | \tau_i, k_0 \stackrel{\text{prior}}{\sim} N(\mu_0, (\tau_i k_0)^{-1})$ ► LIO prior parameter settings: (standardized, unstandardized; finite variance of errors for NFT) $(0, m_0; 0)$ elicited median of the data (1, s_0 ; 1) elicited $\frac{1}{2}$ distance from the median to the 95%-ile if no other recourse, then m_0 and s_0 can be set empirically $\mu_0 = (0, m_0/s_0; 0)$ $k_0 \stackrel{\text{prior}}{\sim} \text{Gamma}(1.5, (7.5, 7.5/s_0^2; 7.5))$ $a_0 = (1.5, 2; 3)$ $b_0 \stackrel{\text{prior}}{\sim} \text{Gamma}\left((0.5, 0.5; 2), (1, 1/s_0^2; 1)\right)$

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 α Prediction Interval $(\mu(x) + c_{\alpha/2}\sigma(x), \ \mu(x) + c_{1-\alpha/2}\sigma(x))$ where $c_{\pi} = F_{\epsilon}^{-1}(\pi)$

$$f(x) = 6x^3$$
, $s(x) = \exp 0.5x$,
log $t = f(x) + s(x)\epsilon$ where $\epsilon \sim t(16)$
and $x \sim U(-1, 1)$: $R^2 = 84.8\%$ uncensored, $R^2 = 85.1\%$ censored



$$f(x) = 6x^3$$
, $s(x) = \exp 0.5x$,
log $t = f(x) + s(x)\epsilon$ where $\epsilon \sim t(16)$
and $x \sim U(-1,1)$: $R^2 = 84.8\%$ uncensored, $R^2 = 85.1\%$ censored



$$f(x) = 6x^3$$
, $s(x) = \exp 0.5x$,
log $t = f(x) + s(x)\epsilon$ where $\epsilon \sim t(16)$
and $x \sim U(-1, 1)$: $R^2 = 84.8\%$ uncensored, $R^2 = 85.1\%$ censored



16/31

$$f(x) = 6x^3$$
, $s(x) = \exp 0.5x$,
log $t = f(x) + s(x)\epsilon$ where $\epsilon \sim t(16)$
and $x \sim U(-1, 1)$: $R^2 = 84.8\%$ uncensored, $R^2 = 85.1\%$ 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 $x \sim B(0.5)$



Event-free Survival: TSVS



Steps

Event-free Survival: MDS disease5



Event-free Survival: Recipient Age



Event-free Survival: Donor Age



Event-free Survival: Donor Age Waterfall Plot



Event-free Survival: Donor Sex/Child-birth Parity



Event-free Survival: Donor Sex/Child-birth Parity



Event-free Survival: Donor Sex/Child-birth Parity

