

## Research Article

# Joint Modelling Approaches to Survival Analysis via Likelihood-Based Boosting Techniques

Colin Griesbach <sup>1</sup>, Andreas Groll <sup>2</sup>, and Elisabeth Bergherr <sup>1</sup>

<sup>1</sup>Chair of Spatial Data Science and Statistical Learning, Georg August University, Germany

<sup>2</sup>Department of Statistics, TU Dortmund University, Germany

Correspondence should be addressed to Colin Griesbach; [colin.griesbach@fau.de](mailto:colin.griesbach@fau.de)

Received 30 July 2021; Accepted 15 October 2021; Published 15 November 2021

Academic Editor: Konstantin G. Arbeev

Copyright © 2021 Colin Griesbach et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Joint models are a powerful class of statistical models which apply to any data where event times are recorded alongside a longitudinal outcome by connecting longitudinal and time-to-event data within a joint likelihood allowing for quantification of the association between the two outcomes without possible bias. In order to make joint models feasible for regularization and variable selection, a statistical boosting algorithm has been proposed, which fits joint models using component-wise gradient boosting techniques. However, these methods have well-known limitations, i.e., they provide no balanced updating procedure for random effects in longitudinal analysis and tend to return biased effect estimation for time-dependent covariates in survival analysis. In this manuscript, we adapt likelihood-based boosting techniques to the framework of joint models and propose a novel algorithm in order to improve inference where gradient boosting has said limitations. The algorithm represents a novel boosting approach allowing for time-dependent covariates in survival analysis and in addition offers variable selection for joint models, which is evaluated via simulations and real world application modelling CD4 cell counts of patients infected with human immunodeficiency virus (HIV). Overall, the method stands out with respect to variable selection properties and represents an accessible way to boosting for time-dependent covariates in survival analysis, which lays a foundation for all kinds of possible extensions.

## 1. Introduction

First suggested by [1], joint models were established as a valuable tool for analysing data where event times are measured alongside a longitudinal outcome. One naive approach of evaluating such frequently occurring data structures would be separate modelling, i.e., fitting appropriate models independently for longitudinal and time-to-event data. However, separate modelling neither corrects for event-dependent dropout in longitudinal analysis nor quantifies the relation between a time-dependent covariate and the risk for an event in survival analysis [2]. Various approaches have been proposed to address these issues, one being the Andersen-Gill model [3] for time-varying covariates in survival analysis or two-stage approaches, where longitudinal model fits are included as fixed covariates in time-to-event regression. It has been shown, however, that these methods tend to produce biased

results [4, 5]. One solution therefore is combining both the survival and longitudinal models within one single joint likelihood. A wide introduction to this joint modelling framework is presented in [4] including the JM package [6]. Moreover, an evolution of joint model progression up to the year 2004 is provided in [7], and in addition, several Bayesian approaches have been carried out [8–10].

Current joint modelling estimation methods, however, lack clear concepts for proper variable selection and good performance regarding prediction. Moreover they are not feasible for high-dimensional data, in particular where the number of covariates exceeds the number of observations, i.e.,  $p > n$  problems. In order to overcome these hindrances, an algorithm was initially proposed, where joint models are fitted with gradient boosting techniques, which are known for addressing exactly these issues [11]. Evolved from machine learning as an approach to classification problems

(i) **Initialize** starting values  $\hat{\beta}_0^{[0]}, \hat{\beta}_t^{[0]}, \hat{\beta}_1^{[0]}, \hat{\lambda}^{[0]}, \hat{\beta}_s^{[0]}, \hat{\alpha}^{[0]}, \hat{\gamma}_0^{[0]}$ , and  $\hat{\gamma}_t^{[0]}$  with variance-covariance-components  $\hat{\sigma}^{2[0]}, \hat{\mathbf{Q}}^{[0]}$  according to Section 3.3. Choose  $(m_{\max,l}, m_{\max,s})$  and execute the following for every possible tuple  $(m_l, m_s)$  with  $m_l < m_{\max,l}$ ,  $m_s < m_{\max,s}$ :

**Longitudinal part**

(ii) **for**  $m = 1$  to  $m_l$  **do**

**step1: Update fixed effects**

For  $r = 1, \dots, p_l$  define  $\tilde{\beta}_r := (\hat{\beta}_0^{[m-1]}, \hat{\beta}_t^{[m-1]}, \hat{\beta}_r^{[m-1]})^T$  with  $\hat{\beta}_r^{[m-1]}$  denoting the  $r$ th component of  $\hat{\beta}_1^{[m-1]}$ . Compute the score vector and Fisher matrix

$$\mathbf{s}_r(\tilde{\beta}_r) = \partial \ell^{\text{pen}} / \partial \tilde{\beta}_r, \mathbf{F}_r(\tilde{\beta}_r) = -\mathbb{E}[\partial^2 \ell^{\text{pen}} / \partial \tilde{\beta}_r \partial \tilde{\beta}_r^T],$$

with respect to the current intercept  $\hat{\beta}_0^{[m-1]}$ , time effect  $\hat{\beta}_t^{[m-1]}$ , and the  $r$ th linear effect  $\hat{\beta}_r^{[m-1]}$ . Obtain  $p_l$  possible updates

$$\mathbf{u}_r = \mathbf{F}_r(\tilde{\beta}_r)^{-1} \mathbf{s}_r(\tilde{\beta}_r),$$

and find the best performing component  $* \in \{1, \dots, p_l\}$  according to Section 2.2, yielding the update  $\mathbf{u}_* = (u_0, u_t, u_*)^T$  containing the update for the effect  $*$  with corresponding updates  $u_0$  for intercept and  $u_t$  for the time effect. Receive  $\hat{\beta}_0^{[m]}, \hat{\beta}_t^{[m]}$ , and  $\hat{\beta}_1^{[m]}$  by updating

$$\begin{aligned} \hat{\beta}_0^{[m]} &= \hat{\beta}_0^{[m-1]} + \nu_1 u_0, \hat{\beta}_t^{[m]} = \hat{\beta}_t^{[m-1]} + \nu_1 u_t, \\ \hat{\beta}_r^{[m]} &= \begin{cases} \hat{\beta}_r^{[m-1]} & \text{if } r \neq *, \\ \hat{\beta}_r^{[m-1]} + \nu_1 u_* & \text{if } r = *, \end{cases} \quad r = 1, \dots, p_l. \end{aligned}$$

**step2: Update random effects**

Receive updates

$$\hat{\gamma}_0^{[m-1]} \longrightarrow \hat{\gamma}_0^{[m]}, \hat{\gamma}_t^{[m-1]} \longrightarrow \hat{\gamma}_t^{[m]},$$

for random intercepts and slopes in an additional Fisher scoring step on the penalized log-likelihood  $\ell^{\text{pen}}$ .

**step3: Update variance-covariance components**

Update variance-covariance components

$$\hat{\sigma}^{2[m-1]} \longrightarrow \hat{\sigma}^{2[m]}, \hat{\mathbf{Q}}^{[m-1]} \longrightarrow \hat{\mathbf{Q}}^{[m]}$$

following the description in Section 3.3.

**end for**

Proceed with estimates  $\hat{\beta}_0^{[m]}, \hat{\beta}_t^{[m]}, \hat{\beta}_1^{[m]}, \hat{\gamma}_0^{[m]}, \hat{\gamma}_t^{[m]}$  as fixed values.

**Survival part**

**for**  $m = 1$  to  $m_s$  **do**

**Update survival effects**

For  $r = 1, \dots, p_s + 1$  define  $\tilde{\beta}_r := (\hat{\lambda}^{[m-1]T}, \hat{\beta}_r^{[m-1]})^T$  with  $\hat{\beta}_r^{[m-1]}$  denoting the  $r$ th component of  $\hat{\beta}_s^{[m-1]}$  and  $\hat{\beta}_{p_s+1}^{[m-1]} = \hat{\alpha}^{[m-1]}$ . Compute the score vector and Fisher matrix

$$\mathbf{s}_r(\tilde{\beta}_r) = \partial \ell^{\text{pen}} / \partial \tilde{\beta}_r, \mathbf{F}_r(\tilde{\beta}_r) = -\mathbb{E}[\partial^2 \ell^{\text{pen}} / \partial \tilde{\beta}_r \partial \tilde{\beta}_r^T],$$

with respect to the current baseline hazard  $\hat{\lambda}^{[m-1]}$  and the  $r$ th linear effect  $\hat{\beta}_r^{[m-1]}$  or  $\hat{\alpha}^{[m-1]}$ , respectively. Obtain  $p_s + 1$  possible updates

$$\mathbf{u}_r = \mathbf{F}_r(\tilde{\beta}_r)^{-1} \mathbf{s}_r(\tilde{\beta}_r),$$

and find the best performing effect  $* \in \{1, \dots, p_s + 1\}$  according to Section 2.2, yielding the update  $\mathbf{u}_* = (\mathbf{u}_\lambda^T, u_*)^T$  containing the update for the effect  $*$  with corresponding baseline hazard update  $u_\lambda$ . Receive  $\hat{\lambda}^{[m]}, \hat{\beta}_s^{[m]}$ , and  $\hat{\alpha}^{[m]}$  by updating

$$\begin{aligned} \hat{\lambda}^{[m]} &= \hat{\lambda}^{[m-1]} + \nu_s \mathbf{u}_\lambda, \hat{\alpha}^{[m]} = \begin{cases} \hat{\alpha}^{[m-1]}, & \text{if } * \neq p_s + 1, \\ \hat{\alpha}^{[m-1]} + \nu_s u_*, & \text{if } * = p_s + 1, \end{cases} \\ \hat{\beta}_r^{[m]} &= \begin{cases} \hat{\beta}_r^{[m-1]} & \text{if } r \neq *, \\ \hat{\beta}_r^{[m-1]} + \nu_1 u_* & \text{if } r = *, \end{cases} \quad r = 1, \dots, p_s. \end{aligned}$$

**end for**

(iii) **Determine** the best performing tuple  $(m_{*,l}, m_{*,s})$  with respect to prediction based on the unpenalized joint log-likelihood  $\ell$  as explained in more detail in Section 3.3.

ALGORITHM 1: lbbJM.

originally proposed in [12], gradient boosting deals with high-dimensional data and the component-wise updating scheme offers implicit variable selection. The boosting algorithm for joint models was extended [13], but when wanting to lay more focus on the survival side of the model, gradient

boosting proved to struggle with time-varying covariates in time-to-event analysis. This has also been observed for pure survival models in [14].

Hence, this work focuses on likelihood-based boosting. First introduced in [15], likelihood-based boosting is

TABLE 1: Shrinkage and variable selection properties regarding longitudinal and survival outcomes averaged over 100 simulation runs of the low-dimensional scenario.

	$\beta_t$ (sd)	$\beta_{11}$ (sd)	$\beta_{12}$ (sd)	$\beta_{13}$ (sd)	TP	FDR	$m_1^*$
True	2	1	2	1			
JM	1.998 (0.03)	0.994 (0.07)	2.008 (0.07)	1.002 (0.07)	—	—	—
lbbJM <sup>a</sup>	1.760 (0.08)	0.914 (0.07)	1.922 (0.07)	0.923 (0.07)	1.00	0.23	108.25
lbbJM <sup>b</sup>	1.992 (0.03)	0.994 (0.07)	2.008 (0.07)	1.002 (0.07)	—	—	—
	$\alpha$ (sd)	$\beta_{s1}$ (sd)	$\beta_{s2}$ (sd)	$\beta_{s3}$ (sd)	TP	FDR	$m_s^*$
True	0.5	1	2	-2			
JM	0.457 (0.04)	0.903 (0.08)	1.807 (0.12)	-1.800 (0.12)	—	—	—
lbbJM <sup>a</sup>	0.390 (0.03)	0.728 (0.06)	1.521 (0.07)	-1.516 (0.07)	1.00	0.27	209.2
lbbJM <sup>b</sup>	0.373 (0.03)	0.713 (0.06)	1.500 (0.07)	-1.495 (0.08)	1.00	0.22	196.9
glmnet	0.427 (0.03)	0.909 (0.07)	1.833 (0.11)	-1.823 (0.10)	1.00	0.51	—

designed to directly maximize a given likelihood where concepts, which gradient boosting implicitly offers, are reproduced artificially for regular optimization methods like Newton algorithms or Fisher scoring. The method was further developed for flexible semiparametric mixed models [16] and for several classes of generalised mixed models [17–20]. The R package GMMBoost [21] covers most of these approaches. Likelihood-based boosting has also been proved useful for survival analysis with time-varying effects [22], and a general overview is given in [23]. Since the random structure plays an important role as a connector between longitudinal and time-to-event data, we additionally incorporate a novel correction step within the estimation procedure for the random effects, which was first suggested in [24, 25] and reduces possible bias arising from wrongly identified random effects.

The contribution of this work is the novel lbbJM boosting algorithm for joint models, which offers the first boosting-based regularization approach for time-dependent covariates in survival analysis and in addition new variable selection mechanics for joint models with focus on time-to-event analysis.

The remainder of this manuscript is structured as follows: Section 2 highlights the overall concepts of both joint modelling and likelihood-based boosting to give a sufficient understanding of the methods used in the following parts. Section 3 then contains a detailed description of the considered joint model together with the proposed boosting algorithm and its computational details. Sections 4 and 5 deal with applying the algorithm to different setups of simulated data as well as to the AIDS dataset [26] included in the JM

package. Results and possible extensions are discussed in the final section.

## 2. Backgrounds

Before the algorithm is presented and discussed in detail, we briefly highlight the concepts of both joint modelling and likelihood-based boosting.

*2.1. Joint Models.* In general, a joint model consists of two parts, one longitudinal and one survival submodel. A popular view on joint modelling is to choose one model as the *main model*, whereas the other model then features the analysis occurring in the main model. With the primary outcome being longitudinal data, a survival model can be used to correct for event-dependent dropout in longitudinal analysis. For time-to-event data as outcome of interest, additional longitudinal modelling reduces measurement error on the one hand and, on the other hand, extrapolates only on single time points observed longitudinal data to continuous functions which are then included in survival analysis. We will from now on focus on joint models with time-to-event data as primary outcome.

The longitudinal submodel takes the form:

$$y = \eta_1(t, \mathbf{x}) + \varepsilon, \quad (1)$$

where longitudinal outcome  $y$  is described by the longitudinal predictor function  $\eta_1$  depending on time  $t$  and a set of covariates  $\mathbf{x}$ . Although  $t$  can be included in  $\mathbf{x}$ , we will highlight it in the context of joint models, as the role of  $t$  is of greater importance. In the survival submodel, the hazard

$$\lambda(t | \mathbf{x}) = \lambda_0(t) \exp(\eta_s(\mathbf{x}) + \alpha\eta_1(t, \mathbf{x})), \quad (2)$$

is modelled by a baseline hazard  $\lambda_0(t)$  with multiplicative effects described by the survival predictor function  $\eta_s$ . In addition, the longitudinal predictor  $\eta_1$  is reappearing in the survival model, this time scaled by a factor  $\alpha$ . The parameter  $\alpha$  thus quantifies the association of the two submodels and is therefore called the *association parameter*. It can be interpreted as the impact a time-varying longitudinal covariate has on the hazard for an event.

Parameter estimation for such joint models can be done in various ways. Two common methods are two-stage and joint likelihood approaches, respectively. In the former, the longitudinal model is estimated with the estimation method of choice leading to the model fit  $\hat{\eta}_1$ , which is then carried forward as fixed covariate into survival analysis. In the latter, longitudinal and survival submodels are combined in a single joint likelihood. Let  $i = 1, \dots, n$  denote clusters and  $j = 1, \dots, n_i$  the repeated measurements. Assuming independent data generating processes for both submodels, the joint likelihood can then be written as

$$L(y, T, \delta) = \prod_{i=1}^n \left( \prod_{j=1}^{n_i} f_1(y_{ij} | \eta_1) \right) f_s(T_i, \delta_i | \alpha, \eta_1, \eta_s, \lambda_0), \quad (3)$$

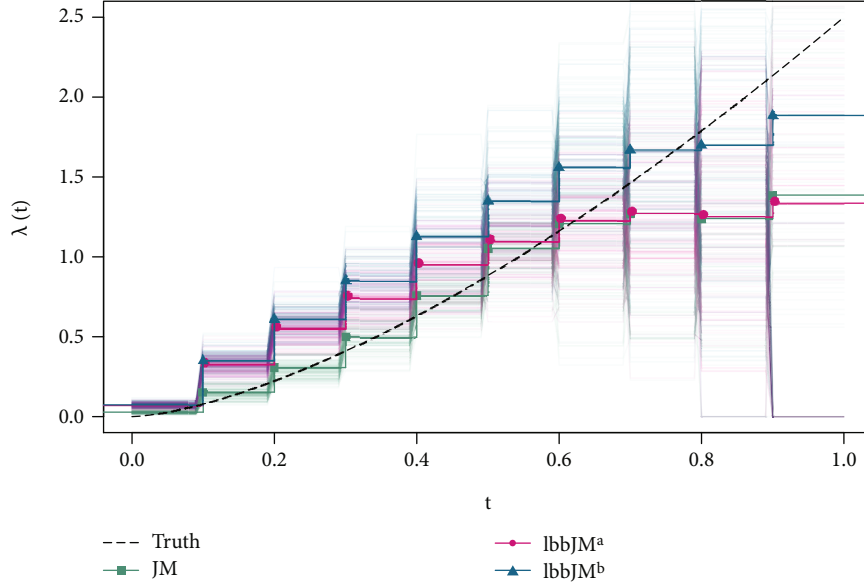


FIGURE 1: Piecewise-constant baseline hazard estimates with  $K = 10$  by JM, lbbJM<sup>a</sup> and lbbJM<sup>b</sup> averaged over 100 simulation runs of the low dimensional scenario.

TABLE 2: Shrinkage and variable selection properties regarding the longitudinal and survival outcomes averaged over 100 simulation runs of the high-dimensional scenario.

	$\beta_t$ (sd)	$\beta_{11}$ (sd)	$\beta_{12}$ (sd)	$\beta_{13}$ (sd)	TP	FDR	$m_1^*$
True	2	1	2	1			
lbbJM <sup>a</sup>	1.748 (0.20)	0.868 (0.14)	1.843 (0.14)	0.875 (0.14)	1.00	0.36	124.2
lbbJM <sup>b</sup>	1.991 (0.08)	1.008 (0.13)	1.982 (0.16)	1.011 (0.15)	—	—	—
	$\alpha$ (sd)	$\beta_{s1}$ (sd)	$\beta_{s2}$ (sd)	$\beta_{s3}$ (sd)	TP	FDR	$m_s^*$
True	0.5	1	2	-2			
lbbJM <sup>a</sup>	0.307 (0.06)	0.512 (0.12)	1.242 (0.16)	-1.216 (0.13)	1.00	0.70	136.7
lbbJM <sup>b</sup>	0.285 (0.05)	0.498 (0.13)	1.215 (0.15)	-1.191 (0.13)	1.00	0.67	127.0
glmnet	0.293 (0.08)	0.627 (0.18)	1.449 (0.30)	-1.422 (0.28)	1.00	0.83	—

with densities  $f_1$  and  $f_s$  for the longitudinal and survival sub-models and time-to-event outcome  $(T, \delta) = (T_i, \delta_i)_{i \in \mathbb{N}}$ . Regular inference is now done by maximizing (3) using appropriate maximization methods, the most prominent one being EM algorithms.

**2.2. Likelihood-Based Boosting.** We intend to give a short description of the underlying mechanics used in the following section. The overall concept of likelihood-based boosting is to create an iterative and component-wise updating scheme, which eventually converges to a maximum likelihood estimator but is stopped early in order to prevent overfitting. Let  $\beta$  model be the effect of  $p$  covariates. Likelihood-

based boosting maximizes a given log-likelihood  $l(\beta)$  by component-wise Fisher scoring in the following way:

For each covariate  $r \in \{1, \dots, p\}$  consider the subvector  $\beta_r$  containing only the coefficients referring to the  $r$ th covariate. We compute the score vector and Fisher matrix as

$$\mathbf{s}_r(\beta_r) = \frac{\partial l(\beta)}{\partial \beta_r}, \quad \mathbf{F}_r(\beta_r) = -\mathbb{E} \left[ \frac{\partial^2 l(\beta)}{\partial \beta_r \partial \beta_r^T} \right], \quad (4)$$

and obtain a possible update

$$\mathbf{u}_r := \mathbf{F}_r(\beta_r)^{-1} \mathbf{s}_r(\beta_r), \quad (5)$$

for the  $r$ th component. Now, we determine the best performing covariate with respect to likelihood maximization, i.e., find the component

$$* = \arg \max_{r=1, \dots, p} l(\tilde{\beta}_r), \quad \tilde{\beta}_r = (\beta_1, \dots, \beta_r + \mathbf{u}_r, \dots, \beta_p) \quad (6)$$

where the corresponding update yields the biggest improvement of the likelihood. One receives a new model fit by weakly updating this best performing component, i.e., by scaling with a factor  $\nu$ , the so called *step length*:

$$\beta_r^{\text{new}} = \begin{cases} \beta_r & \text{if } r \neq *, \\ \beta_r + \nu \mathbf{u}_r & \text{if } r = *, 0 < \nu \leq 1, \end{cases} \quad r = 1, \dots, p. \quad (7)$$

The step length  $\nu$  is controlling the weakness of the update to prevent overfitting and give every covariate a chance for selection. A popular choice in the literature is  $\nu = 0.1$ .

TABLE 3: Averaged computation times for one single model fit (in seconds).

Setup	JM	glmnet	lbbJM <sup>a</sup>	lbbJM <sup>b</sup>
Low	110.00	149.15	15776.16	43.76
High	—	156.44	4072.80	248.08

TABLE 4: Structure of the dataset with primary outcomes for the joint analysis in the three columns on the left.

$y$	$T$	$\delta$	$t$	Drug	Gender	AZT	prevOI	ID
10.67	16.97	0	0	ddC	Male	Intolerance	AIDS	1
8.43	16.97	0	6	ddC	Male	Intolerance	AIDS	1
9.43	16.97	0	12	ddC	Male	Intolerance	AIDS	1
6.32	19.00	0	0	ddI	Male	Intolerance	noAIDS	2
8.12	19.00	0	6	ddI	Male	Intolerance	noAIDS	2
4.58	19.00	0	12	ddI	Male	Intolerance	noAIDS	2
5.00	19.00	0	18	ddI	Male	Intolerance	noAIDS	2
3.46	18.53	0	0	ddI	Female	Intolerance	AIDS	3
3.61	18.53	0	2	ddI	Female	Intolerance	AIDS	3
6.16	18.53	1	6	ddI	Female	Intolerance	AIDS	3
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

Repeating this updating process for a sufficiently large number of iterations leads to the regular maximum likelihood estimator for  $\beta$ . But instead the algorithm is stopped early to gain better prediction performance and variable selection. The optimal amount of iterations actually is a tuning parameter of the method and can be determined via cross-validation or by focusing on information criteria like AIC or BIC [27, 28].

### 3. Boosting Joint Models

*3.1. The Model.* Before we introduce the boosting algorithm, we describe the specific joint model. With  $i = 1, \dots, n$  denoting the individual and  $j = 1, \dots, n_i$  a single specific measurement, the longitudinal submodel is given by

$$y_{ij} = \eta_1(t_{ij}, \mathbf{x}_{i1}) + \varepsilon_{ij} = \beta_0 + \beta_t t_{ij} + \beta_1^T \mathbf{x}_{i1} + \gamma_{0i} + \gamma_{ti} t_{ij} + \varepsilon_{ij}, \quad (8)$$

where  $y_{ij}$  is modelled by  $\eta_1$  depending on specific measurement times  $t_{ij}$  and longitudinal time-independent covariates  $\mathbf{x}_{i1} \in \mathbb{R}^{p_1}$ . This represents a standard linear mixed model with intercept  $\beta_0$  and fixed linear effects  $\beta_t$  and  $\beta_1$  of time and baseline covariates as well as individual specific random effects  $\gamma_{0i}$  and  $\gamma_{ti}$  with  $(\gamma_{0i}, \gamma_{ti}) \sim \mathcal{N}^{\otimes 2}(\mathbf{0}, \mathbf{Q})$ . The error terms  $\varepsilon_{ij}$  are assumed to follow a normal distribution with  $\mathbb{E}[\varepsilon_{ij}] = 0$  and  $\text{Var}(\varepsilon_{ij}) = \sigma^2 > 0$ .

Please note that the model can be additionally extended to interaction effects of time  $t_{ij}$  with baseline covariates  $\mathbf{x}_{i1} \in \mathbb{R}^{p_1}$  by including the term  $\beta_t^T \mathbf{x}_{i1} t_{ij}$  in (8). This results in

slightly adjusted integrals in the survival part and is omitted in the following for the sake of better readability.

In the survival part, the individual hazard

$$\lambda_i(t) = \lambda_0(t) \exp(\eta_s(\mathbf{x}_{si}) + \alpha \eta_1(t, \mathbf{x}_{ii})), \quad (9)$$

is modelled with the survival predictor  $\eta_s(\mathbf{x}_{si}) = \beta_s^T \mathbf{x}_{si}$  containing additional linear effects  $\beta_s$  of baseline covariates  $\mathbf{x}_{si} \in \mathbb{R}^{p_s}$ . To execute a full likelihood approach, the baseline hazard

$$\lambda_0(t) = \sum_{k=1}^K \lambda_k \mathbf{1}_{I_k}(t), \quad (10)$$

is chosen to be piecewise-constant depending on the number of segments  $K$  and their exact locations  $I_k = [t_{k-1}, t_k)$  with  $t_0 = 0$  and  $t_k = \max(\mathbf{T})$  for  $k = 1, \dots, K$ . The collection of values for the baseline hazard is denoted in  $\lambda = (\lambda_k)_{k=1, \dots, K}$ . Later, we will choose  $K$  between 7 and 10 in order to guarantee substantially more flexibility than a constant baseline hazard without becoming computationally too demanding.

Given two formulas (8) and (9), we can now calculate the joint log-likelihood. Let  $\mathbf{y} = (y_{ij})_{i=1, \dots, n, j=1, \dots, n_i}$  denote the collection of all longitudinal measurements. Assuming the time-to-event process is conditionally independent from the longitudinal random structure, the joint likelihood can be decomposed into a longitudinal and a survival term. Set  $\boldsymbol{\gamma}^T = (\boldsymbol{\gamma}_i)_{i=1, \dots, n}$  with  $\boldsymbol{\gamma}_i^T = (\gamma_{0i}, \gamma_{ti})$  and  $\vartheta_1 := (\beta_0, \beta_t^T, \beta_1^T, \boldsymbol{\gamma}^T)^T$ . Furthermore,  $\tau$  contains information on variance-covariance components  $\sigma^2$  and  $\mathbf{Q}$ . The unpenalized longitudinal log-likelihood is

$$\ell_1(\vartheta_1, \sigma^2 | \mathbf{y}) = \sum_{i=1}^n \sum_{j=1}^{n_i} \log \phi(y_{ij} | \eta_1(t_{ij}, \mathbf{x}_{i1}), \sigma^2), \quad (11)$$

where  $\phi(\cdot | m, v)$  denotes the density of a normal distribution with mean  $m$  and variance  $v$ . Laplace approximation follows [29] and then leads to an additional quadratic penalty term for the random effects yielding the penalized log-likelihood:

$$\ell_1^{\text{pen}}(\vartheta_1, \tau | \mathbf{y}) = \sum_{i=1}^n \left( \sum_{j=1}^{n_i} \log \phi(y_{ij} | \eta_1(t_{ij}, \mathbf{x}_{i1}), \sigma^2) - \frac{1}{2} \boldsymbol{\gamma}_i^T \mathbf{Q} \boldsymbol{\gamma}_i \right). \quad (12)$$

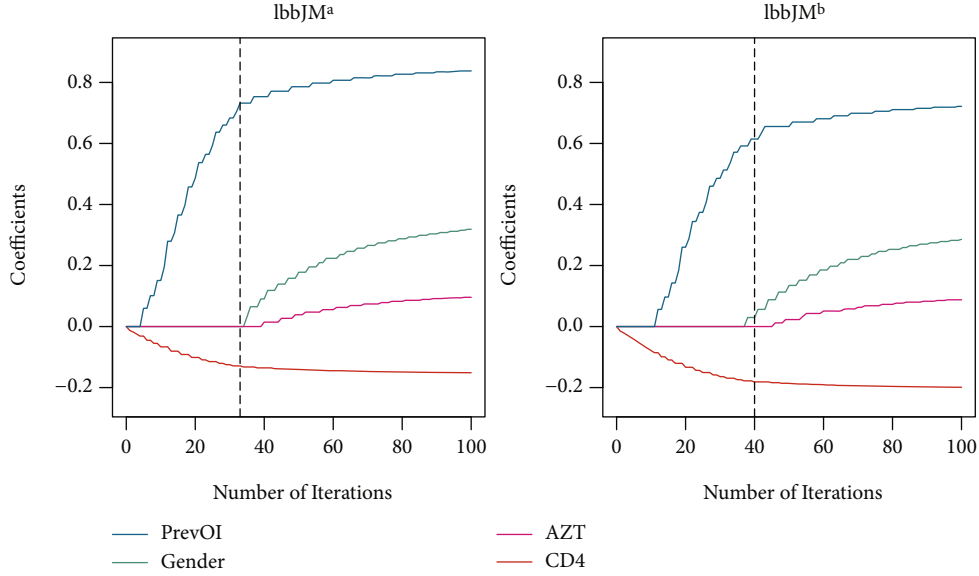
Note that for the penalized log-likelihood  $\tau$  substitutes  $\sigma^2$  as an argument, since the penalized log-likelihood additionally contains information of the variance matrix  $\mathbf{Q}$ .

On the other hand, for given survival outcome  $(\mathbf{T}, \delta) = (T_i, \delta_i)_{i=1, \dots, n}$  with event times  $T_i$  and censoring indicator  $\delta_i$ , the survival log-likelihood takes the well-known form:

$$\ell_s(\vartheta_s | \mathbf{T}, \delta) = \sum_{i=1}^n \delta_i \log \lambda_i(T_i | \eta_1, \eta_s, \alpha, \lambda) - \int_0^{T_i} \exp(\lambda_i(t | \eta_1, \eta_s, \alpha, \lambda)) dt, \quad (13)$$

TABLE 5: Shrinkage and variable selection properties by the different packages for model (22).

	$\beta_0$	$\beta_1$	$\beta_{t1}$	$\beta_{t2}$	$\alpha$	$\beta_{s1}$	$\beta_{s2}$	$\beta_{s3}$
JM	6.97	0.49	-0.18	<0.01	-0.24	0.31	0.09	0.66
lbbJM <sup>a</sup>	6.95	0.26	-0.05	0	-0.13	0	0	0.73
lbbJM <sup>b</sup>	6.95	0.48	-0.16	-0.02	-0.18	0.03	0	0.61
glmnet	—	—	—	—	-0.15	0.31	0.09	0.81
$p$ value (JM)	<0.01	0.26	<0.01	0.98	<0.01	0.23	0.61	<0.01

FIGURE 2: Coefficient progression in the survival part for lbbJM<sup>a</sup> ((a), with  $m_{*1} = 10$ ) and lbbJM<sup>b</sup> (b).

with  $\vartheta_s := (\lambda^T, \alpha, \beta_s^T)^T$ . For  $\vartheta := (\vartheta_1^T, \vartheta_s^T)^T$ , we finally receive the penalized and unpenalized joint log-likelihood:

$$\begin{aligned} \ell(\vartheta, \sigma^2 | \mathbf{y}, \mathbf{T}, \delta) &= \ell_1(\vartheta_1, \sigma^2 | \mathbf{y}) + \ell_s(\vartheta_s | \mathbf{T}, \delta), \\ \ell^{\text{pen}}(\vartheta, \tau | \mathbf{y}, \mathbf{T}, \delta) &= \ell_1^{\text{pen}}(\vartheta_1, \tau | \mathbf{y}) + \ell_s(\vartheta_s | \mathbf{T}, \delta). \end{aligned} \quad (14)$$

**3.2. The Algorithm.** The following lbbJM algorithm (likelihood-based boosting for joint models) describes a way to fit the formulated joint model by likelihood-based boosting methods discussed in Section 2.2.

**3.3. Computational Details of the Algorithm.** In general, the algorithm carries out a hybrid between a two-stage and a joint likelihood approach. For one single tuple  $(m_1, m_s)$ , the fitting procedure goes as follows: In a first step, the longitudinal submodel is boosted up to  $m_1$  iterations using the lbbLMM boosting algorithm [25]. The received estimates are carried forward into the survival model, where another boosting process up to  $m_s$  iterations takes place. This fitting process is carried out for any tuples  $(m_1, m_s)$  with  $m_1 < m_{\max,1}$ ,  $m_s < m_{\max,s}$ , where  $m_{\max,1}$  and  $m_{\max,s}$  are prespecified maximum numbers of iterations per submodel. For every of these possible combinations of stopping iterations, the corresponding estimates are evaluated based on the joint likelihood using test data, which can be achieved via cross-

validation or bootstrapping. Hence, the algorithm uses two-stage fitting but joint likelihood evaluation. We give a detailed description for both parts. Exact formulas for all appearing variants of score vectors and Fisher matrices can be found in the supplementary material. Please note that due to the component-wise updating scheme in both sub-models, the lbbJM algorithm works with arbitrarily high numbers of candidate variables and is therefore not confined to low-dimensional data structures.

For starting values, the parameters, which actually underlie the boosting process, are necessarily set to zero, thus  $\hat{\beta}_1^{[0]} = \hat{\beta}_s^{[0]} = \hat{\alpha}^{[0]} = \mathbf{0}$ . The baseline hazard is initialized with the intercept estimator  $\hat{\lambda}^{[0]} = (\sum_i \delta_i / \sum_i T_i)_{k=1, \dots, K}$ . The remaining values are extracted from a standard linear mixed model for time and random effects

$$y_{ij} = \beta_0 + \beta_t t_{ij} + \gamma_{0i} + \gamma_{ti} t_{ij} + \varepsilon_{ij}, \quad (15)$$

by using, e.g., the function `lme` from the R package `nlme`.

For boosting longitudinal fixed effects, for convenience, we omit the iteration index as well as the hat indicating estimated values in the following subsections. In the first step of the longitudinal part, the effects  $\beta_1$  follow the classical component-wise likelihood-based boosting procedure. In

each iteration, an update for every single covariate  $r$  together with intercept  $\beta_0$  and time effect  $\beta_t$  is computed, leading to  $p_1$  three-dimensional updates. Selection of the best performing component is then performed either by selecting the component yielding the optimal likelihood maximization or lowest information criteria like AIC or BIC, which minimizes the model complexity rather than residuals. The linear effect  $\beta_t$  is excluded from the selection process, as it holds a very important role in a joint model and should always be included.

For updating random effects, after the best performing fixed effect from  $\beta_1$  was updated, in every iteration, an update for the random effects is executed separately. This means that the score vector and Fisher matrix

$$\mathbf{s}_{\text{ran}}(\gamma) = \frac{\partial \ell^{\text{pen}}}{\partial \gamma}, \quad \mathbf{F}_{\text{ran}}(\gamma) = -\mathbb{E} \left[ \frac{\partial^2 \ell^{\text{pen}}}{\partial \gamma \partial \gamma^T} \right], \quad (16)$$

have to be derived in order to execute the update

$$\gamma^{\text{new}} = \gamma + \mathbf{C} \mathbf{F}_{\text{ran}}(\gamma)^{-1} \mathbf{s}_{\text{ran}}(\gamma). \quad (17)$$

The matrix  $\mathbf{C}$  is a correction matrix which prevents from potential correlations between the random effects estimates and any observed covariates [25], and its derivation can be traced in more detail in the supplementary material.

For updating variance-covariance components, the covariance matrix  $\mathbf{Q}$  of the random effects is updated with an approximate EM algorithm using the posterior curvatures  $\mathbf{F}_{ii}$  of the random effects model [30]. Receive an update by computing

$$\mathbf{Q} = \frac{1}{n} \sum_{i=1}^n (\mathbf{F}_{ii}^{-1} + \gamma_i \gamma_i^T). \quad (18)$$

The current longitudinal model error is obtained by setting  $\text{Var}(\mathbf{y} - \boldsymbol{\eta}_1)$ .

For boosting the association parameter and survival effects, once the longitudinal part was updated in up to  $m_1$  iterations, the algorithm proceeds to boost the effects  $\beta_s$  and  $\alpha$ . Although being of different structures, the association parameter  $\alpha$  is boosted alongside the effects  $\beta_s$ , meaning the algorithm decides whether the association or some baseline survival effect is updated, based on which parameter leads to the best likelihood improvement. This means, the linear effect of the whole longitudinal trajectory is also scaled by the step length  $\nu$  when being updated within the selection step, which minimizes the chance of potential overfitting also for the association parameter. An alternative method would be choosing just from the effects in  $\beta_s$  and updating  $\alpha$  in an additional step by optimizing the current likelihood. This approach was used in [11] and treats  $\alpha$  as a nuisance parameter, which might not be satisfactory with regards to the importance of  $\alpha$ . Again, only the update for  $\alpha$  and  $\beta_s$  is scaled by the step length  $\nu_s$ . The baseline hazard  $\lambda$  receives a full update.

For step lengths, apart from the stopping iterations, the step lengths are tuning parameters of the boosting algorithm. Although there is some effort in focusing on adaptive step lengths recently, we chose to set both step lengths to the constant value  $\nu_1 = \nu_s = 0.1$ . The exact choice of the step length factor is of minor importance as long as it is sufficiently small to ensure proper performance. Setting it to 0.1 is an established choice in the boosting literature [31, 32].

For stopping iterations, since the step lengths are chosen to be constant, the tuple  $(m_1, m_s)$  is the main tuning parameter of the boosting algorithm. In regular boosting with only one iteration index, it is convenient to check for every single iteration and take the estimates from the estimation count leading to the best prediction. In the present two-dimensional case, this would mean finding

$$(m_{*,1}, m_{*,s}) = \arg \max_{(m_1, m_s) \in \mathcal{M}} \ell \left( \hat{\vartheta}^{[m_1, m_s]} \mid \mathcal{X}^{\text{test}} \right), \quad (19)$$

with  $\mathcal{M} := \{1, \dots, m_{\text{max},1}\} \times \{1, \dots, m_{\text{max},s}\}$ ,  $\hat{\vartheta}^{[m_1, m_s]}$  denoting the vector of estimates received via the tuple  $(m_1, m_s)$  of total iterations and  $\mathcal{X}^{\text{test}}$  a complete set of test data for evaluation. Problem (19) is then solved via  $k$ -fold cross-validation. But since checking for every single tuple  $(m_1, m_s) \in \mathcal{M}$  implies a very high computational effort, we suggest to coarsen the grid and check for fewer possible stopping iterations in the longitudinal part, e.g.,  $m_1 \in \{10, 20, 30, \dots, m_{\text{max},1}\}$ . Because of the two-stage-approach nature of the algorithm, we still can check for every single  $m_s \in \{1, \dots, m_{\text{max},s}\}$  without gaining additional computational effort. Furthermore, parallel computing can be executed in order to minimize computational demand.

## 4. Simulations

We evaluate the lbbJM algorithm with a simulation study. The aim is to assess estimation and shrinkage characteristics in general as well as variable selection properties and performance in high dimensional, i.e.,  $p > n$  settings. The lbbJM algorithm is included in two variants. While lbbJM<sup>a</sup> executes the full approach as depicted in Section 3.2, lbbJM<sup>b</sup> performs a shortened two-stage procedure where the longitudinal sub-model is fitted in advance using regular maximum likelihood inference and does not underlie any regularization. The exact lbbJM<sup>b</sup> algorithm is depicted in detail in the supplementary material. We additionally include the JM package as state of the art for convenient estimation of joint models as well as the glmnet package, which offers elastic net penalization for start-stop-data and therefore an alternative approach for regularization of time-dependent covariates in survival analysis. None of the competitors are completely suitable for a benchmark comparison and are viewed as reference points for the specific objectives addressed by the lbbJM algorithm. Regarding glmnet, as an alternative approach to regularization of time-dependent covariates, shrinkage and variable selection properties are of interest, although it focuses solely on survival analysis. JM in addition offers unregularized effect estimates with corresponding

significance indicator but is neither suitable for high-dimensional setups nor offers variable selection.

*4.1. Setup.* The simulations are executed according to the model described in Section 3.1 with  $n \in \{100, 500\}$  and  $n_i = 5$  using inversion sampling, which is explained in detail in the supplementary material. The prespecified true parameter values are

$$\beta_0 = 1, \quad \beta_t = 2, \quad \beta_1^T = (2, 1, 2), \quad \beta_s^T = (1, 2, -1), \quad \alpha = 0.5, \quad (20)$$

with variance components

$$\sigma = 0.1, \quad \mathbf{Q} = \begin{pmatrix} 2 & 0.1 \\ 0.1 & 0.3 \end{pmatrix}. \quad (21)$$

The entries of the covariate vectors  $\mathbf{x}_{li}$  and  $\mathbf{x}_{si}$  are drawn independently from the standard normal distribution  $\mathcal{N}(0, 1)$ . In addition to the informative covariates with effects  $\beta_1$  and  $\beta_s$ , the covariate vectors  $\mathbf{x}_{li}$  and  $\mathbf{x}_{si}$  are expanded with noninformative noise variables until the chosen numbers  $p_1$  and  $p_s$  of total dimensions are reached. These additional noise variables are included to evaluate variable selection properties and robustness of the approach in case of a misspecified model. The baseline hazard is chosen as  $\lambda_0(t) = 2.5t^{1.5}$  and given the censoring mechanism described in the supplementary material, the chosen parameter values result in an average censoring rate of  $\approx 50\%$ .

Overall, we consider two scenarios. One low-dimensional setup with  $n = 500$  and  $p_1 = p_s = 9$  mimicking a more common data structure and one high-dimensional setup, where the number of covariates included in the survival submodel exceeds the number of individuals so that conventional approaches like JM fail to return results.

For the computation, JM and lbbJM use  $K = 10$  with equidistant knot placement. The grid  $\mathcal{M} := \{25, 50, \dots, 500\} \times \{1, 2, \dots, 1000\}$  is specified for possible tuples of stopping iterations and the optimal regularization parameter for glmnet is determined by the function `cv.glmnet()` based on 10-fold cross-validation.

*4.2. Results.* Since the compared estimation routines follow different approaches targeting various objectives from regular maximum likelihood estimation in joint models to regularization in pure time-to-event analysis, we focus on plain coefficient estimates averaged over 100 independent simulation runs in order to assess estimation characteristics. Variable selection properties are evaluated by considering share of true positives (TP) and false discovery rate (FDR).

For low-dimensional setup ( $n = 500$ ,  $p_s = 9$ ), Table 1 depicts the results for the low-dimensional setup. In the longitudinal submodel, lbbJM<sup>a</sup> has small shrinkage and therefore offers variable selection with a rather low false discovery rate of 0.23 but still selects all informative variables. The time effect  $\beta_t$  receives comparatively high shrinkage, since time, as a cluster-varying variable, adds more information to the model. Overall, the longitudinal submo-

del is boosted up to 108.25 iterations on average yielding rather strongly shrunk coefficient estimates. Please note that the results for lbbJM<sup>b</sup> are simply obtained by `lme()` and very similar to JM. In the survival part, both boosting approaches substantially outperform glmnet in terms of variable selection while again receiving also more shrinkage. Due to the comparatively rough baseline hazard depicted in Figure 1, all full likelihood approaches, i.e., JM and lbbJM, receive additional shrinkage which is unaffected by possible regularization. As glmnet uses the partial likelihood, there are no estimates for the baseline hazard function available and the small elastic net penalty of  $\lambda^* = 0.003$  on average also results in weaker performance regarding the rate of false positives. Overall, the lbbJM approaches yield satisfactory results regarding both regularization and variable selection. The effect estimates clearly reflect the true values approximately obtained by JM while simultaneously receiving sufficiently large shrinkage for decent performance of identifying influential covariates in both the longitudinal and survival submodels.

For high-dimensional setup ( $n = 100$ ,  $p_s = 100$ ), Table 2 depicts the results for the high-dimensional setup. As expected, estimates in the high-dimensional setup are regularized stronger and therefore experience more shrinkage. Again, the boosting approaches contained in lbbJM yield better variable selection properties and slightly more regularized coefficient estimates, although results seem to align with increasing dimensions. Note that JM is not capable of handling high-dimensional data structures and is therefore not included in the high-dimensional setup at all.

For computational effort, Table 3 shows estimates for elapsed computation time of each routine. Times are measured in seconds and depict the computation time for one single model fit which was carried out on a  $2 \times 2.66$  GHz-6-Core Intel Xeon CPU (64 GB RAM). As expected, the full boosting approach executed in lbbJM<sup>a</sup> comes with high computational costs similar to [11]. Note that the runtimes are higher in the low-dimensional scenario as the overall number of clusters is higher ( $n = 500$ ) leading to an increased burden in the already time-consuming longitudinal boosting procedure. The reduced approach lbbJM<sup>b</sup>, however, runs considerably faster and is therefore more desirable as long as research focus lies solely on the time-to-event analysis.

## 5. Application

We showcase the lbbJM algorithm by applying it to the 1994 AIDS data [26]. The study is aimed at comparing the two antiretroviral drugs, didanosine (ddI) and zalcitabine (ddC), based on a collective of 467 patients infected with human immune deficiency virus (HIV) who were either intolerant to or failed a previous treatment with Zidovudine (AZT). Alongside several baseline covariates, the square root CD4 cell count was recorded at study entry and in multiple follow-ups after 2, 6, 12, and 18 months, respectively. The CD4 cells are attacked by the virus and thus decrease over time for infected patients; hence, they are a widely used surrogate for disease progression. In addition to the



longitudinal outcome, 188 patients died during the time of the study leading to 188 observed and 279 censored events. The structure of the data is depicted in Table 4.

We formulate the joint model

$$CD4_i(t) = \eta_{li}(t) + \varepsilon_i(t) = \beta_0 + \beta_1 \text{drug}_i + \beta_{11} t + \beta_{12}(t \cdot \text{drug}_i) + \gamma_{0i} + \gamma_{1i} t + \varepsilon_i(t), \quad (22)$$

$$\lambda_i(t) = \lambda_0(t) \exp(\beta_{s1} \cdot \text{gender}_i + \beta_{s2} \cdot \text{AZT}_i + \beta_{s3} \cdot \text{prevOI}_i + \alpha \eta_{li}(t)). \quad (23)$$

The CD4 cell count  $CD4_i(t)$  for  $i = 1, \dots, 467$  is described by a linear mixed model with random intercepts, random slopes for time  $t$  and linear effects of time, drug and an additional interaction between time and drug. The *true*, i.e. modelled by  $\eta_{li}(t)$ , underlying profile of the CD4 cell count is then included together with the remaining baseline covariates in the Cox full likelihood model, thus the model error  $\varepsilon(t)$  is eliminated. Here, drug is a dummy for ddI, gender for female gender, AZT for failure of Zidovudine therapy, and prevOI for prevalence of AIDS. The number of segments for the baseline hazard was chosen to be  $K = 7$ .

We fit model (22) with the same methods as already used in the simulation study. The tuning parameters of the boosting algorithm were chosen to be  $\nu_1 = \nu_s = 0.1$  for step lengths and  $\mathcal{M}^a := \{5, 10, \dots, 100\} \times \{1, 2, \dots, 250\}$  for the grid of possible tuples of stopping iterations for lbbJM<sup>a</sup> and  $\mathcal{M}^b := \{1, 2, \dots, 250\}$  for lbbJM<sup>b</sup>. Again, glmnet was tuned using the cv.glmnet() function and all regularization approaches are based on 10-fold cross validation.

For lbbJM<sup>a</sup>,  $m_{*,1} = 10$  and  $m_{*,s} = 33$  formed the best performing tuple of stopping iterations. The two-stage approach lbbJM<sup>b</sup> used  $m_{*,s} = 40$  and the optimal penalization parameter for glmnet turned out as  $\lambda_{\text{pen}}^* = 0.0014$ . The corresponding coefficient estimates are shown in Table 5. Overall, the results reflect what was already observed in the simulation study. While glmnet shows quite conservative shrinkage properties where every variable is included in the final model, the lbbJM approaches tend to stop rather early yielding bigger shrinkage and in addition effects, which did not get selected at all. In order to give some point of reference regarding variable selection properties, the  $p$  values computed by JM are included in Table 5. The selected effects align quite nicely with being significant according to JM. While lbbJM<sup>a</sup> only includes the variable drug with only half the effect size, lbbJM<sup>b</sup> additionally sees a quite tiny impact of female gender.

The corresponding coefficient progressions for the survival submodel are visualized in Figure 2. Both algorithms update the coefficients referring to the longitudinal CD4 cell profile and the variable prevOI right away and therefore see a strong connection between the risk for death and the CD4 cell count as well as whether or not a patient has AIDS. Due to early stopping, lbbJM<sup>a</sup> selects neither of the two remaining covariates into the final model. lbbJM<sup>b</sup> includes one variable more, gender, which however only has a very small effect.

## 6. Outlook and Discussion

Overall, the lbbJM algorithm introduces a novel boosting-based regularization scheme to joint models focusing on survival analysis as well as to Cox models with time-dependent covariates. The method fits in well among alternative routines and especially stands out with respect to variable selection properties. Due to its clear advantage regarding computational effort as depicted in Table 3, lbbJM<sup>b</sup> is the preferred routine when research interest clearly lies on time-to-event analysis, whereas lbbJM<sup>a</sup> is capable of regularizing both submodels simultaneously. Besides the good results regarding variable selection, it is also expected that the proposed boosting methods can improve the predictive power of a joint model, since boosting algorithms are, due to their model tuning based on test errors, primarily used for prediction. A thorough investigation of improving and evaluating prediction performance of a joint model by several boosting techniques remains an interesting task.

Still, the presented foundation is of a comparatively simple nature and possible extensions include more flexible modelling, e.g., based on P-splines which allow smooth effects of time-dependent as well as time-independent covariates and can additionally include time-varying effects as well as possibly time-varying association structures. As the current algorithm is confined to a time-constant association parameter  $\alpha$ , an extension to  $\alpha(t)$  would increase the flexibility of the model. While it is usually difficult to disentangle a time-dependent association parameter from the longitudinal trajectory for parameter estimation, the lbbJM algorithm could potentially avoid these identification issues due to its clear separation in a longitudinal and a survival boosting process.

Similar regularization-based approaches have been proved useful for time-to-event settings [22, 33], and it can be assumed that the presented method could benefit from these concepts as well.

Moreover, the presented work only lays a foundation to several extensions addressing known issues for both boosting and joint modelling. It represents an accessible way to boosting for time-dependent covariates in survival analysis. Gradient boosting has known limitations in this matter, and although efforts for a framework to overcome these issues are made [34], things rather quickly become technical and the proven robustness of likelihood-based boosting represents a flexible and far more intuitive alternative. Further developments in this direction could include multiple time-dependent covariates based on the two-stage approach of lbbJM<sup>b</sup>, where additional effort is necessary in order to provide a fair competition between time-dependent and time-independent covariates within the selection procedure.

Furthermore, the component-wise updating process is capable of including allocation mechanisms into the algorithm. While the lbbJM algorithm is due to its two-stage fitting process fairly robust to estimate models, where one candidate variable is assigned to both submodels, this kind of specification is not advised in general as identification issues may arise and a proper interpretation of the resulting model might be challenging. However, it is usually tricky to

decide, whether a covariate should be included in the longitudinal or the survival part of the model and these decisions are often made using prior knowledge. An allocation routine based on likelihood maximization could therefore greatly improve joint model inference and would additionally eliminate the two-stage nature of the lbbJM algorithm. This would not only decrease the computational burden but also provide a far more flexible algorithm allowing for regularization, variable selection, and allocation.

## Data Availability

The R code data used to support the findings of this article are included within the supplementary information files.

## Conflicts of Interest

The authors declare that there is no conflict of interest.

## Acknowledgments

The work on this article was supported by the DFG (Deutsche Forschungsgemeinschaft; Projekt WA 4249/2-1) and the Volkswagen Foundation. We further acknowledge support by the Open Access Publication Funds of the Göttingen University.

## Supplementary Materials

Supplementary Materials: a more detailed discussion on the computational details as well as the lbbJM<sup>b</sup> routine and the simulation algorithm can be found in the web-appendix. Additionally, the corresponding R code is available in the supplementary files. (*Supplementary Materials*)

## References

- [1] M. S. Wulfsohn and A. A. Tsiatis, "A joint model for survival and longitudinal data measured with error biometrics," *Biometrics*, vol. 53, p. 330, 1997.
- [2] X. Guo and B. P. Carlin, "Separate and joint modeling of longitudinal and event time data using standard computer packages," *The American Statistician*, vol. 58, pp. 16–24, 2004.
- [3] P. K. Andersen and R. D. Gill, vol. 10, pp. 1100–1120, 1982.
- [4] D. Rizopoulos, *Joint Models for Longitudinal and Time-to-Event Data: With Applications in R; 6 of Chapman & Hall/CRC Biostatistics Series*, CRC Press, Boca Raton, 2012.
- [5] M. J. Sweeting and S. G. Thompson, "Joint modelling of longitudinal and time-to-event data with application to predicting abdominal aortic aneurysm growth and rupture," *Biometrical Journal*, vol. 53, no. 5, pp. 750–763, 2011.
- [6] D. Rizopoulos, "JM: an R package for the joint modelling of longitudinal and time-to-event data," *Journal of Statistical Software*, vol. 35, no. 9, pp. 1–33, 2010.
- [7] A. A. Tsiatis and M. Davidian, "Joint modeling of longitudinal and time-to-event data: an overview," *Statistica Sinica*, vol. 14, pp. 809–834, 2004.
- [8] M. Köhler, N. Umlauf, A. Beyerlein, C. Winkler, A. G. Ziegler, and S. Greven, "Flexible Bayesian additive joint models with an application to type 1 diabetes research," *Biometrical Journal*, vol. 59, pp. 1144–1165, 2017.
- [9] R. Martins, G. L. Silva, and V. Andreozzi, "Joint analysis of longitudinal and survival AIDS data with a spatial fraction of long-term survivors: a Bayesian approach," *Biometrical Journal*, vol. 59, pp. 1166–1183, 2017.
- [10] M. Rué, E.-R. Andrinopoulou, D. Alvares, C. Armero, A. Forte, and L. Blanch, "Bayesian joint modeling of bivariate longitudinal and competing risks data: an application to study patient-ventilator asynchronies in critical care patients," *Biometrical Journal*, vol. 59, no. 6, pp. 1184–1203, 2017.
- [11] E. Waldmann, D. Taylor-Robinson, N. Klein et al., "Boosting joint models for longitudinal and time-to-event data," *Biometrical Journal*, vol. 59, no. 6, pp. 1104–1121, 2017.
- [12] Y. Freund and R. E. Schapire, "Experiments with a new boosting algorithm," *Proceedings of the Thirteenth International Conference on Machine Learning Theory*, Morgan Kaufmann, San Francisco, 1996.
- [13] C. Griesbach, A. Mayr, and E. Waldmann, "Extension of the Gradient Boosting Algorithm for Joint Modeling of Longitudinal and Time-to-Event Data Ar Xiv e-Prints," *Cornell University*, 2018.
- [14] B. Hofner, "Variable selection and model choice in survival models with time-varying Effects diploma thesis Ludwig-Maximilians-Universität München," *Open Access LMU*, 2008.
- [15] G. Tutz and H. Binder, "Generalized additive modeling with implicit variable selection by likelihood-based boosting," *Biometrics*, vol. 62, no. 4, pp. 961–971, 2006.
- [16] G. Tutz and F. Reithinger, "A boosting approach to flexible semiparametric mixed models," *Statistics in Medicine*, vol. 26, no. 14, pp. 2872–2900, 2007.
- [17] G. Tutz and A. Groll, *Generalized Linear Mixed Models Based on Boosting in Statistical Modelling and Regression Structures*, K. Thomas and T. Gerhard, Eds., Springer-Verlag, Berlin Heidelberg, 2010.
- [18] G. Tutz and A. Groll, "Regularization for generalized additive mixed models by likelihood-based boosting," *Methods of Information in Medicine*, vol. 51, no. 2, pp. 168–177, 2012.
- [19] G. Tutz and A. Groll, "Likelihood-based boosting in binary and ordinal random effects models," *Journal of Computational and Graphical Statistics*, vol. 22, no. 2, pp. 356–378, 2013.
- [20] A. Groll, *Variable selection by regularization methods for generalized mixed models*, Ludwig-Maximilians-Universität München Munich, Germany, 2011.
- [21] A. Groll, *GMMBoost: Likelihood-Based Boosting for Generalized Mixed Models*, p. 1.1.3, 2020, *R package version*.
- [22] B. Hofner, T. Hothorn, and T. Kneib, "Variable selection and model choice in structured survival models," *Computational Statistics*, vol. 28, no. 3, pp. 1079–1101, 2013.
- [23] A. Mayr, H. Binder, O. Gefeller, and M. Schmid, "The evolution of boosting algorithms," *Methods of Information in Medicine*, vol. 53, no. 6, pp. 419–427, 2014.
- [24] C. Griesbach, B. Säfken, and E. Waldmann, "Gradient boosting for linear mixed models," *The International Journal of Biostatistics*, 2021.
- [25] C. Griesbach, A. Groll, and E. Bergherr, "Addressing cluster-constant covariates in mixed effects models via likelihood-based boosting techniques," *PLoS One*, vol. 16, no. 7, article e0254178, 2021.
- [26] D. I. Abrams, A. I. Goldman, C. Launer et al., "A comparative trial of didanosine or zalcitabine after treatment with zidovudine in patients with human immunodeficiency virus

- infection,” *New England Journal of Medicine*, vol. 330, no. 10, pp. 657–662, 1994.
- [27] H. Akaike, *Information Theory and the Extension of the Maximum Likelihood Principle in Second International Symposium on Information Theory: 267-281*, 1973.
- [28] G. Schwarz, “Estimating the dimension of a model,” *Annals of Statistics*, vol. 6, pp. 461–464, 1978.
- [29] N. E. Breslow and D. G. Clayton, “Approximate inference in generalized linear mixed models,” *Journal of the American Statistical Association*, vol. 88, pp. 9–52, 1993.
- [30] L. Fahrmeir and G. Tutz, *Multivariate Statistical Modelling Based on Generalized Linear Models*, Springer-Verlag, New York, 2001.
- [31] J. Friedman, T. Hastie, and R. Tibshirani, “Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors),” *The annals of statistics*, vol. 28, no. 2, pp. 337–407, 2000.
- [32] P. Bühlmann and T. Hothorn, “Boosting algorithms: regularization, prediction and model fitting,” *Statistical Sciences*, vol. 27, pp. 477–505, 2007.
- [33] A. Groll, T. Hastie, and G. Tutz, “Selection of effects in Cox frailty models by regularization methods,” *Biometrics*, vol. 73, no. 3, pp. 846–856, 2017.
- [34] D. K. Lee, N. Chen, and H. Ishwaran, “Boosted nonparametric hazards with time-dependent covariates,” *Annals of Statistics*, vol. 49, no. 4, 2021.