


# Dynamic estimation with random forests for discrete-time survival data

Hooraa MORADIAN<sup>1</sup>, Weichi YAO<sup>2</sup>, Denis LAROCQUE<sup>1\*</sup> , Jeffrey S. SIMONOFF<sup>2</sup>, and Halina FRYDMAN<sup>2</sup>

<sup>1</sup>Department of Decision Sciences, HEC Montréal, Montréal, Québec, Canada

<sup>2</sup>Department of Technology, Operations, and Statistics, Stern School of Business, New York University, New York, New York, USA

*Key words and phrases:* Discrete-time survival analysis; landmark analysis; random forests; survival forests; time-varying covariates.

*MSC 2020:* Primary 62N02; secondary 62G05.

*Abstract:* Time-varying covariates are often available in survival studies, and estimation of the hazard function needs to be updated as new information becomes available. In this article, we investigate several different easy-to-implement ways that random forests can be used for dynamic estimation of the survival or hazard function from discrete-time survival data. Results from a simulation study indicate that all methods can perform well, and that none dominates the others. In general, situations that are more difficult from an estimation point of view (such as weaker signals and less data) favour a global fit, pooling over all time points, while situations that are easier from an estimation point of view (such as stronger signals and more data) favour local fits. *The Canadian Journal of Statistics* 00: 000–000; 2021 © 2021 Statistical Society of Canada

*Résumé:* Dans une analyse de survie, il arrive fréquemment que des variables explicatives dont la valeur changeant dans le temps soient disponibles. Lorsque c'est le cas, les estimations doivent être mise à jour au fur et à mesure que de nouvelles informations s'ajoutent. Dans cet article, nous étudions plusieurs manières, qui peuvent être mises en œuvre facilement, d'utiliser les forêts aléatoires pour obtenir des estimations de la fonction de risque de façon dynamique avec des données de survie à temps discret. Les résultats d'une simulation montrent que toutes les méthodes étudiées performant bien, et qu'aucune ne domine les autres. En général, les situations qui sont plus difficiles du point de vue de l'estimation (un signal plus faible et une taille d'échantillon plus petite) favorisent un modèle global qui regroupe tous les points temporels, tandis que celles plus faciles (un signal plus fort une taille d'échantillon plus grande) favorisent les modèles locaux. *La revue canadienne de statistique* 00: 000–000; 2021 © 2021 Société statistique du Canada

## 1. INTRODUCTION

Survival analysis studies with time-to-event data have applications in many research areas. It is common in practice that the actual time until the occurrence of an event of interest is observed for only some of the subjects and partial information about the time is available for other subjects, for example, because the study ended before all subjects experienced the event, or because some of them were lost during the study. This concept is known as censoring (Klein & Moeschberger, 2003). Right-censoring, i.e., when only a lower bound on the actual time is observed, is the most common situation and will be the main focus of this article. A comprehensive introduction to

---

Additional Supporting Information may be found in the online version of this article at the publisher's website.

\* Corresponding author: [denis.larocque@hec.ca](mailto:denis.larocque@hec.ca)

modelling time-to-event data can be found in Kleinbaum & Klein (2005) and Hosmer, Lemeshow & May (2011).

Many of the traditional methods for analyzing continuous time-to-event data rely on some parametric (e.g., Weibull) or semiparametric (e.g., Cox) assumptions about the link between the covariates and the time response, which may result in poor performance in real-world applications. Recently, more flexible models and adapted machine learning algorithms that use data to find relevant structures, instead of imposing them a priori, have been developed in the survival analysis domain (Wang, Li & Reddy, 2019). One class of such models is tree-based methods, which are the focus of this article.

Tree-based methods were first developed for a categorical or continuous outcome. Breiman et al. (1984) is the earliest monograph about trees and details the classification and regression tree (CART) paradigm. Gordon & Olshen (1985) extended the tree paradigm to survival data and introduced survival trees (Segal, 1988; Leblanc & Crowley, 1993). However, it is well known that ensembles of trees often provide better estimation performance than a single tree. One popular and efficient ensemble method is the random forest, introduced by Breiman (2001) and extended to model right-censored survival data (Ishwaran et al., 2004, 2008; Hothorn et al., 2006; Zhu & Kosorok, 2012). There is a vast literature on survival trees and forests. Bou-Hamad, Larocque & Ben-Ameur (2011b) present a general overview.

In many studies, an estimate of the hazard function for a subject is obtained at time 0 using only the baseline covariate information. However, when time-varying covariates are present, it is often preferable to update the estimates of hazard probabilities as new longitudinal information becomes available. This is the topic of “dynamic estimation,” an area of growing interest. There are primarily three approaches to building dynamic estimates in this context: (1) landmark analysis, (2) joint modelling and (3) a counting process approach. The idea of landmark analysis (Anderson, Cain & Gelber, 1983; Madsen, Hougaard & Gilpin, 1983) is to build models, usually Cox, at different landmark times  $t$  using the covariate information available up to  $t$  from those subjects who are still at risk of experiencing the event at  $t$ . Comprehensive treatments of this approach are given in van Houwelingen (2007) and van Houwelingen & Putter (2011). The second approach uses joint modelling of the time-varying covariate processes and the event time data process (Henderson, Diggle & Dobson, 2000). This approach depends on the correct specification of the model for the time-varying covariate trajectories, and this problem amplifies as the number of time-varying covariates increases. The main idea of the third approach is to partition the follow-up information for each individual into multiple segments on nonoverlapping intervals (Bacchetti & Segal, 1995). This is used to accommodate time-varying covariates in the tree-building process (Bertolet, Brooks & Bittner, 2016; Fu & Simonoff, 2017b). Survival forest algorithms based on this same counting process approach can then be developed to provide dynamic estimation of hazards or survival probabilities (Wongvibulsin, Wu & Zeger, 2020; Yao et al., 2020).

Most of the research, including the work cited above, assumes that the time to event is measured continuously when, in fact, it is measured on a discrete scale in many cases. This can happen with binned data where the event occurs in an interval of time and the intervals are not necessarily of the same length. For example, the Framingham Heart Study<sup>1</sup> requires the participants to return to the study approximately every 2–6 years in order for their medical history data to be collected and physical exams and laboratory tests done. Another example of binned data is term insurance, or any other annual contract with churn (lack of renewal of the contract) being the event of interest. Alternatively, the observed time may come from a truly discrete process, such as the number of elapsed time units or trials before reaching a specific

---

<sup>1</sup><https://www.nhlbi.nih.gov/science/framingham-heart-study-fhs>.

goal (e.g., the number of cycles until pregnancy). Although traditional modelling approaches for continuous-time survival data can also be applied to discrete-time survival data, Tutz & Schmid (2016) explain the advantages of using statistical methods that are specifically designed for discrete event times. They point out that the hazard functions derived in the discrete case are more easily interpretable than those for continuous survival time data, since the hazards can be formulated as conditional probabilities. Moreover, discrete models do not have any problems dealing with ties. Therefore, in this article we focus on methods specifically designed for discrete-time survival data.

Survival trees and forests designed specifically for discrete-time responses were developed by Bou-Hamad et al. (2009), Bou-Hamad, Larocque & Ben-Ameur (2011a), Schmid et al. (2016, 2020) and Berger et al. (2019). Section 2.1 provides a description of some of these methods since they are central to this article. Elgmati et al. (2015) propose a penalized Aalen additive model for dynamic estimation of the hazard function for discrete-time recurrent event data, but their method is limited to one-step-ahead estimation while we also explore multistep-ahead estimation.

From the above discussion, we see that no tree-based methods have addressed the problem of dynamic estimation with discrete survival responses. In this article, we investigate different ways that random forests can be used for dynamic estimation of hazard function with discrete-time survival response data.

The rest of the article is organized as follows. Section 2 describes the data setting and the proposed methods. The results from a simulation study are presented in Section 3. Section 4 provides conclusions and directions for future work. More details about the simulation study and a real example using bankruptcy data can be found in the Supplementary Material.

## 2. DESCRIPTION OF THE METHODS

Suppose we have data on  $N$  independent subjects. For subject  $i$ , observations are in the form of  $(\tau_i, \delta_i, \mathbf{x}_i)$  where  $\tau_i \in \{1, 2, \dots, T\}$  is the discrete time to event;  $T$  is the maximum observed time in the dataset;  $\delta_i$  is the censoring index, which takes a value of 0 if the observation for the subject  $i$  is right-censored and a value of 1 if subject  $i$  has experienced the event of interest; and  $\mathbf{x}_i$  is a set of covariates, some of which can be time-varying and some time-invariant. We will denote by  $x_{ki}(t)$  the value of the  $k$ th covariate,  $k \in \{1, 2, \dots, p\}$ , at time  $t \in \{0, 1, \dots, T-1\}$  for subject  $i$ . Hence,  $\mathbf{x}_i(0)$  gives the baseline covariate values. For simplicity, we will use this notation for all covariates, time-varying or not. Hence  $x_{ki}(t)$  remains constant for all  $t$  for a time-invariant covariate. The values of the actual time to event and the censoring times for subject  $i$  are denoted by  $U_i$  and  $V_i$ , respectively. Hence we have  $\tau_i = \min(U_i, V_i)$ , and we assume that  $U_i$  and  $V_i$  are independent given  $\mathbf{x}_i$ . The hazard function for subject  $i$  is denoted by  $h_i(t) = P(U_i = t \mid U_i \geq t)$  for simplicity, but it is obvious that  $\tau_i, \delta_i, U_i$  and  $V_i$  depend on  $\mathbf{x}_i$ . Similarly, the survival function for subject  $i$  is  $S_i(t) = P(U_i > t)$ , and the probability that the event occurs at time  $t$  is  $\pi_i(t) = P(U_i = t)$ . These two functions can be obtained from the hazard function with the recursive formulae  $S_i(t) = S_i(t-1)(1 - h_i(t))$  and  $\pi_i(t) = S_i(t-1) - S_i(t)$ , with  $S_i(0) = 1$ . Hence, it is sufficient to model the hazard function (or any one of the other two functions) to recover the other ones.

### 2.1. Description of Existing Methods for Discrete-Time Survival Data

The existing methods for dynamic estimation based on time-varying covariate data use a counting process approach to reformat the data. To fix ideas, a generic dataset of 10 observations with two covariates,  $X_1$  being time-varying and  $X_2$  being time-invariant, is given in Table 1. For instance, the first subject experienced the event at the second time point and thus values of the time-varying covariate  $X_1(t)$  are available only up to  $t = 1$ , with NA's for the other time points. Note that we do not assume that the covariate values at the event or censoring time are available

TABLE 1: A generic dataset with 10 observations and two covariates, with  $X_1$  being time-varying and  $X_2$  time-invariant.

id	$\tau$	$\delta$	$X_1(0)$	$X_1(1)$	$X_1(2)$	$X_1(3)$	$X_1(4)$	$X_2$
1	2	1	$x_{11}(0)$	$x_{11}(1)$	NA	NA	NA	$x_{21}$
2	4	1	$x_{12}(0)$	$x_{12}(1)$	$x_{12}(2)$	$x_{12}(3)$	NA	$x_{22}$
3	3	0	$x_{13}(0)$	$x_{13}(1)$	$x_{13}(2)$	NA	NA	$x_{23}$
4	1	0	$x_{14}(0)$	NA	NA	NA	NA	$x_{24}$
5	4	1	$x_{15}(0)$	$x_{15}(1)$	$x_{15}(2)$	$x_{15}(3)$	NA	$x_{25}$
6	4	0	$x_{16}(0)$	$x_{16}(1)$	$x_{16}(2)$	$x_{16}(3)$	NA	$x_{26}$
7	2	0	$x_{17}(0)$	$x_{17}(1)$	NA	NA	NA	$x_{27}$
8	4	1	$x_{18}(0)$	$x_{18}(1)$	$x_{18}(2)$	$x_{18}(3)$	NA	$x_{28}$
9	3	1	$x_{19}(0)$	$x_{19}(1)$	$x_{19}(2)$	NA	NA	$x_{29}$
10	4	0	$x_{110}(0)$	$x_{110}(1)$	$x_{110}(2)$	$x_{110}(3)$	NA	$x_{210}$

(e.g., the event or censoring may occur before the observation of the covariates). This process is repeated for each of the subjects in the dataset. The reformatted dataset is often called the “person–period” dataset.

We describe first the existing approaches for estimating the hazard of a subject at the  $u$ th discrete time point that use the last available values of the time-varying covariates. One widely used method is the discrete-time proportional odds (DTPO) model, also known as the continuation ratio model

$$\log \left( \frac{h_i(u)}{1 - h_i(u)} \right) = \alpha_1 D_{1i}(u) + \cdots + \alpha_T D_{Ti}(u) + \beta_1 X_{1i}(u-1) + \cdots + \beta_p X_{pi}(u-1), \quad (1)$$

for  $i = 1, 2, \dots, n$  and  $u = 1, 2, \dots, T$ , where the  $D_{ri}(u)$ 's are indicator variables indexing the  $r$ th discrete time point, which are defined by  $D_{ri}(u) = 1$  if  $r = u$  and 0 otherwise. The intercept parameters  $\alpha_1, \alpha_2, \dots, \alpha_T$  define the baseline hazard at each time point, and the  $\beta$  coefficients describe the effects of covariates on the baseline hazard function. Applying the counting process approach to reformat the generic dataset gives the person–period data in Table 2. The model parameters in (1) can then be estimated by fitting a logistic regression to the reformatted dataset. More detail can be found in Willett & Singer (1993), p. 171.

Bou-Hamad, Larocque & Ben-Ameur (2011a) were the first to propose building trees and forests using the person–period dataset with  $y$  as the response and a likelihood-based splitting criterion. Schmid et al. (2016) proposed a classification tree by applying the CART algorithm based on the Gini impurity measure (Breiman et al., 1984) to the same dataset, again with  $y$  as the response. Schmid et al. (2020) proposed building discrete-time random survival forests using Hellinger’s distance criterion (Cieslak et al., 2012) as the splitting rule. This criterion was also implemented in a classification tree approach for the modelling of competing risks in discrete time (Berger et al., 2019). Numerical results given in Schmid et al. (2020) suggest that node-splitting by Hellinger’s distance improves the performance when compared to skew-sensitive split criteria such as the Gini impurity. This is consistent with the results of simulations performed here, and therefore we investigate forest methods using only Hellinger’s distance criterion. The time point

TABLE 2: Person–period dataset using the counting process approach for DTPO model. Only the first two subjects (up to id = 2) are shown to save space. It has one row of observation for each discrete time point  $u$  in which the subject is at risk of experiencing the event and the response  $y$  equals 1 if the event occurred at that time and 0 otherwise.

id	$y$	$u$	$D_1$	$D_2$	$D_3$	$D_4$	$X_1$	$X_2$
1	0	1	1	0	0	0	$x_{11}(0)$	$x_{21}$
1	1	2	0	1	0	0	$x_{11}(1)$	$x_{21}$
2	0	1	1	0	0	0	$x_{12}(0)$	$x_{22}$
2	0	2	0	1	0	0	$x_{12}(1)$	$x_{22}$
2	0	3	0	0	1	0	$x_{12}(2)$	$x_{22}$
2	1	4	0	0	0	1	$x_{12}(3)$	$x_{22}$

TABLE 3: The 10 different estimating problems when  $T = 4$ . For instance, at time point  $t = 2$ , given a subject who has survived up to this time point, we are interested in estimating its hazard function at the future time points  $u = 3, 4$ .

$t$	$u$			
Value	1	2	3	4
0	✓	✓	✓	✓
1		✓	✓	✓
2			✓	✓
3				✓

$u$  itself is also included as an ordinal covariate (Schmid et al., 2016; Berger et al., 2019; Schmid et al., 2020). To fix ideas, with the dataset in Table 2, this means building a classification forest with  $y$  as the response using the three covariates  $X_1, X_2$ , and the time point  $u$ . Using the time point as a predictor implies that the subjects can be split apart in the person–period data even if no time-varying covariates are present among the original covariates, since the time point itself is a time-varying covariate. In a terminal node, the estimate of the hazard is the proportion of 1 (events) in the node.

### 2.2. Description of the Set-up for Dynamic Estimation

In line with the purpose of dynamic estimation, where we want to estimate future risks, at the current time point  $t$  the goal is to estimate the hazard of a subject at some future time point  $u$  for  $u = t + 1, t + 2, \dots, T$ . We assume that measurements for all covariates are available at  $0, 1, 2, \dots, t$ , and the methods are entitled to use all of that information. Hence, all covariate information up to time  $t$  can be used to estimate the hazard function at  $u$ . Table 3 illustrates the possible combinations of  $t$  and  $u$  with  $T = 4$  as an example. One can also see that, for a given value of  $T$ , the total number of possible estimation problems is  $T(T + 1)/2$  ( $= 10$  when  $T = 4$ ). For the following discussion,  $t$  always denotes the current time point,  $u$  always denotes the future time point we are interested in for estimation, and  $u > t$  by definition.

For simplicity of the presentation, we will use only the last available value of the time-varying covariates to build the models. However, without loss of generality, we can assume that any past information we also want to use is already incorporated into the covariates at the current time point  $t$ . For example, if we want to use the lag of a time-varying covariate, say  $X_1(t-1)$ , we can simply define a new covariate at time  $t$  to represent this lag, that is,  $\tilde{X}_1(t) = X_1(t-1)$ .

We investigate different methods to solve the hazard function estimation problem for each pair  $(t, u)$  as illustrated in Table 3. These methods can be divided into three main approaches to address the same estimation problem based on how they make use of the information provided in the generic dataset, i.e., how they construct the training datasets.

Given the estimation problem for a specific pair  $(t^*, u^*)$ , the first approach is to use only corresponding local information to train the model. More precisely, to construct the training dataset to estimate the hazard for the given pair  $(t^*, u^*)$ , we consider only the subjects that are still alive and not censored at time point  $u^* - 1$ , as these subjects are still at risk of experiencing the event at time point  $u^*$ . Moreover, the training dataset contains only their covariate information at the current time point  $t^*$ . For a subject with covariate information available up to time  $t^*$ , this approach builds separate models to estimate the hazard function at each future time point. Using separate models might be effective if the hazards at different time points are related to different covariate patterns, but this approach will likely lose efficiency when the hazards are related to similar covariate patterns because of the variability induced by using separate models.

The second approach solves the estimation problems for all future time points at once, from a given time point  $t^*$ . In this case, for a given  $t^*$ , we construct a single training dataset that pools the local information  $(t^*, u)$  from all possible values of  $u$ . This can reduce the variability when the hazards at a given time point are related to similar covariate patterns. All the covariates are used, and the future time point  $u$  itself is also considered as a covariate. The model trained with this dataset is then used to estimate all future hazards for any subject, with its current covariate information at the given time  $t^*$ . The Schmid et al. (2020) method presented in the last section builds the forest based on this idea.

The third approach is inspired by the so-called “supermodel” based on stacked data used in landmark analysis, presented by van Houwelingen (2007) and van Houwelingen & Putter (2011). Instead of pooling the information from the different estimation horizons only for a given  $t^*$ , as in the second approach, we can go a step further and pool all the information for all combinations of  $(t, u)$  together. The idea is to borrow information from different values of  $t$ , in addition to that of different estimation horizons for a given  $t^*$ . This results in a super person–period training dataset which is created by stacking the training datasets from all values of  $t$  used in the Schmid et al. (2020) method described above. The model trained on this super person–period dataset is then used to estimate hazard probabilities for a subject at any future time points with covariate information available at any current time point. This time, both the estimation horizon  $u$  and the value of  $t$  are potential covariates, in addition to the other covariates.

Table 4 provides an illustration of the training dataset used for all three approaches to solve each of the 10 estimation problems given in Table 3. The person–period dataset is reformatted based on the generic dataset given in Table 1. Each subject has one row for each pair value of  $(t, u)$  where it was still at risk of experiencing the event, i.e., neither its event time nor censoring time has yet occurred at  $u - 1$ . Only the first three subjects (up to  $\text{id} = 3$ ) are shown in the table to save space. For example, to solve the estimation problem for the pair  $(t^*, u^*) = (1, 2)$ , i.e., to estimate the hazard probability for any subject at time point 2 with its covariate information at time point 1, the training dataset used for the separate method would be the one given in rows 10–12 in Table 4. Note that only the subjects whose event time and censoring time have not yet occurred at  $u^* - 1 = 1$  are included. The outcome  $y$  has a value of 1 if the event occurred

TABLE 4: Training dataset used for the three approaches to solve each of the estimation problems given in Table 3 ( $T = 4$ ): (i) the first approach—Separate; (ii) the second approach—the Schmid et al. (2020) method (iii) the third approach—super person–period. Only the first three subjects (up to id = 3) are shown to save space.

row	id	y	Available covariates				Box of data used to train a given Method (Covariates used) to estimate hazards for which value of $(t, u)$ .		
			$X_1$	$X_2$	$u$	$t$	Separate $(X_1, X_2)$	Schmid et al. (2020) $(X_1, X_2, u)$	Super person-period $(X_1, X_2, u, t)$
1	1	0	$x_{11}(0)$	$x_{21}$	1	0	(0, 1)	All possible combinations of $(t, u): t < u, t = 0, 1, \dots, 4, u = 1, 2, \dots, 4.$	
2	2	0	$x_{12}(0)$	$x_{22}$	1	0			
3	3	0	$x_{13}(0)$	$x_{23}$	1	0			
4	1	1	$x_{11}(0)$	$x_{21}$	2	0	(0, 2)		
5	2	0	$x_{12}(0)$	$x_{22}$	2	0			
6	3	0	$x_{13}(0)$	$x_{23}$	2	0			
7	2	0	$x_{12}(0)$	$x_{22}$	3	0	(0, 3)		
8	3	0	$x_{13}(0)$	$x_{23}$	3	0			
9	2	0	$x_{12}(0)$	$x_{22}$	4	0			
10	1	1	$x_{11}(1)$	$x_{21}$	2	1	(1, 2)		
11	2	0	$x_{12}(1)$	$x_{22}$	2	1			
12	3	0	$x_{13}(1)$	$x_{23}$	2	1			
13	2	0	$x_{12}(1)$	$x_{22}$	3	1	(1, 3)		
14	3	0	$x_{13}(1)$	$x_{23}$	3	1			
15	2	0	$x_{12}(1)$	$x_{22}$	4	1			
16	2	0	$x_{12}(2)$	$x_{22}$	3	2	(2, 3)		
17	3	0	$x_{13}(2)$	$x_{23}$	3	2			
18	2	0	$x_{12}(2)$	$x_{22}$	4	2			
19	2	0	$x_{12}(3)$	$x_{22}$	4	3	(3, 4)		

at time point  $u^* = 2$ , and 0 otherwise. Two covariates are used for this method,  $X_1$  and  $X_2$ . For the same problem, the Schmid et al. (2020) method uses the training dataset as given in rows 10–15 in Table 4 and adds  $u$  as a covariate. The third approach uses  $X_1, X_2, u$  and  $t$  as covariates. Its training dataset consists of all rows of the person–period data. One can see that to produce 10 estimated hazard probabilities, one for each estimation problem as given in Table 3, the first approach builds 10 models (one for each pair of  $(t, u)$ ), the second approach builds 4 models (one for each  $t$ ), and the third approach builds only 1 model (one for all pairs  $(t, u)$ ).



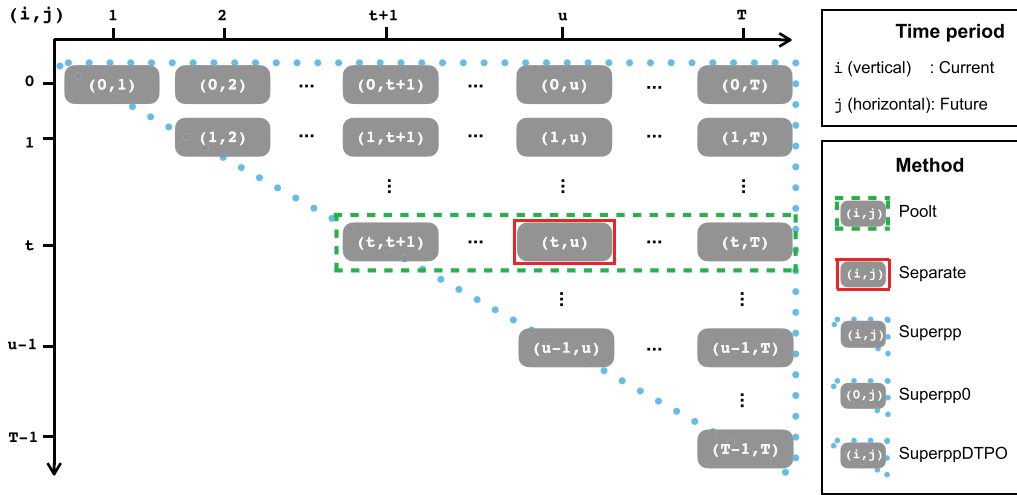


FIGURE 1: Graphical overview of the methods for dynamic estimation.

In the simulations summarized in the next section, we investigate these three approaches applied to random forest methods: separate forests, forests using the Schmid et al. (2020) method, and a forest built on the super person–period dataset, which will be referred to as “Separate,” “Poolt” and “Superpp,” respectively.

We also compare the performance of these three methods to the following two methods in the simulation study:

- (1) Super person–period forest with baseline information only. That is, Superpp using only the covariate information at  $t = 0$ . This method will be referred to as “Superpp0.”
- (2) DTPO model using the super person–period construction. This method will be referred to as “SuperppDTPO.”

Note that SuperppDTPO targets the log-linear survival relationship. Superpp0 is a nonparametric method but never updates the information from the initial status. These two methods serve as benchmark parametric and nonparametric methods, respectively, as we investigate the performance of the three methods under different model set-ups.

Figure 1 provides a graphical overview of the methods for dynamic estimation. Consider the set of time points  $\{0, 1, \dots, T\}$ . Each entry  $(i, j)$  contains the subjects that are still alive and not censored at time  $j - 1$ , and it gives the information available at time  $i$  from each subject in that cell. Suppose we are at the current time point  $t$  and want to estimate the hazard function for some future time point  $u (> t)$ . The red box (solid line) contains all the subject information that Separate uses for training the forest model, the green box (dashed line) contains all that Poolt uses, and the blue triangular region (dotted line) contains all that Superpp uses. Note that SuperppDTPO uses the same subject information as Superpp, and Superpp0 also uses the same subjects but with all  $(i, j)$  replaced by  $(0, j)$ , indicating it uses the baseline information only. There is in total one blue region,  $T$  green boxes and  $T \times (T + 1)/2$  red boxes, implying that the construction of one Superpp forest, one Superpp0 forest, one SuperppDTPO model,  $T$  Poolt forests and  $T \times (T + 1)/2$  Separate forests are used to construct estimates of hazards for all combinations of  $(t, u)$  for a given value of  $T$ . Note that Table 4 is a specific example of this construction where there are only three subjects and  $T = 4$ .



### 3. SIMULATION STUDY

R (R Core Team, 2020) was used to perform the simulations. The package `ranger` (Wright, Wager & Probst, 2020) was used to build the forests with the Hellinger splitting rule for the methods `Separate`, `Poolt`, `Superpp` and `Superpp0`, that is, all methods that require a classification forest. The number of trees in all forests is 500. `SuperppDTPO` was implemented using logistic regression on the `Superpp` dataset.

#### 3.1. Simulation Design

The data generating process (DGP) is a discretized version of the continuous-time survival data generated from the model used in the simulation study in Yao et al. (2020). We consider the following factors for different variations of DGPs:

- (1) Different combinations of numbers of time-invariant and time-varying covariates in the true generating model (Scenario).
- (2) Different matrices to generate covariates' values with autocorrelation for the time-varying variables (labelled as "Strong" and "Weak"). Note that stronger autocorrelation would imply that covariate values from earlier time points would tend to be more similar to those in later time points, making future estimation easier.
- (3) Different signal-to-noise ratios (SNRs) labelled as "High" and "Low," constructed by choosing different coefficients in the model.
- (4) Different survival distributions: Exponential, Weibull and Gompertz.
- (5) Different survival relationships between the hazards and covariates: a log-linear one, a log-nonlinear one and a log-interaction model.
- (6) Different censoring rates: 10% and 50%.
- (7) Different training sample sizes:  $n = 200, 1000$  and  $5000$ .
- (8) Different total numbers of time points:  $T = 4$  and  $8$ .

Each model is fitted with a training sample of size 1000. The performance of the fitted models is then evaluated with  $T$  independent test sets of size 1000 each. The  $k$ th test set ( $k = 1, 2, \dots, T$ ) includes only the subjects that are still at risk at  $u = k$ , so it can be used when  $(t, u) = (j, k)$  for all  $j = 0, 1, \dots, k - 1$ . Each simulation is repeated 500 times. See Section S1.1 in the Supplementary Material for more details of the simulation design.

#### 3.2. Simulation Results

We consider three criteria for evaluating the accuracy of the methods: absolute difference (ADIST), absolute log odds ratio (ALOR) and concordance index (C-index) for hazard rates. Let  $\hat{h}$  and  $h$  be the estimated and the true hazards. ADIST is defined by

$$\text{ADIST}(h, \hat{h}) = |\hat{h} - h|,$$

and ALOR by

$$\text{ALOR}(h, \hat{h}) = |\ln((\hat{h}(1 - h))/((1 - \hat{h})h))|.$$

Both ADIST and ALOR take a minimum value of 0 when  $\hat{h} = h$ , while ALOR also takes the magnitude of  $h$  and  $\hat{h}$  into account. The C-index computes the proportion of concordant pairs over all possible evaluation pairs:

$$C = \frac{\sum_{i \neq j} \mathbb{I}(h_i > h_j) \cdot \mathbb{I}(\hat{h}_i > \hat{h}_j)}{\sum_{i \neq j} \mathbb{I}(h_i > h_j)},$$

where the indices  $i$  and  $j$  refer to pairs of hazards in the test sample for a given combination of  $(t, u)$ . It is designed to estimate the concordance probability  $\mathbb{P}(\hat{h}_i > \hat{h}_j | h_i > h_j)$ , which compares the rankings of two independent pairs of hazard rates  $h_i, h_j$  and estimates  $\hat{h}_i, \hat{h}_j$ . Concordance probability evaluates whether the values of  $\hat{h}_i$  are directly associated with the values of  $h_i$ . Note that while both ADIST and ALOR measure the distance between the true hazard and its estimate, the C-index is a rank-based metric that evaluates whether the true and estimated values are ordered similarly; a high value does not necessarily imply that the estimated values are close to the true ones.

Extensive simulation studies show that the total number of time points  $T$  in the true model does not affect the general conclusions. In the following discussion, we focus on the cases where  $T = 4$  (see Table S1.2 in the Supplementary Material for performance comparison between  $T = 4$  and  $T = 8$ ).

Boxplots for the 500 simulation runs of each method for each combination  $(t, u)$  based on the evaluation of ADIST and C-index are provided in Section S1.2 in the Supplementary Material. Boxplot results from ALOR are not reported since the conclusions are essentially the same as those from ADIST (ALOR results for performance comparison are still provided in summary tables in Section S1.3 in the Supplementary Material). Figures 2 and 3 give an example of the boxplots for ADIST and C-index, respectively, when the training sample size is 1000, the censoring rate is 10%, and the data are generated following a Weibull distribution with an interaction survival relationship in the scenario 2TI+4TV (two time-invariant and four time-varying covariates), with high SNR and strong autocorrelation, and with a total number of time points of  $T = 4$ . In general, for a given  $t$  (i.e., for a given plot), the performance of the methods usually worsens as  $u$  increases. This is expected because it is more difficult to estimate the hazard for horizons further away.

From the boxplots based on ADIST evaluation, the parametric SuperppDTPO method works well as expected when the underlying survival relationship is linear. In most other cases, it is outperformed by the nonparametric forest methods. Superpp always gives better performance for dynamic estimation than Superpp0, which is again expected as the latter uses only the baseline covariate values. In general, the three forest methods that use all the covariate information, namely Separate, Poolt and Superpp, perform the best compared to the other two simpler methods, presumably because the hazard estimates from the three forests are less biased in general due to the flexibility of the estimators.

Note that the boxplot results for evaluation from ADIST and those from C-index do not always agree with each other. In particular, the C-index tends to favour SuperppDTPO in general. For example, Figure 3 shows that SuperppDTPO outperforms Separate when  $(t, u) = (1, 4)$  and dominates the other methods when  $(t, u) = (2, 4)$ , while in Figure 2 it gives the worst performance among all methods in both cases. As noted, ADIST is a calibration metric, whereas C-index is a rank-based metric. Bias is more important for accurate estimation of hazards, while variance is more important for accurate ordering of hazards. This results in favourable performance for forests using the time-varying information for the ADIST criterion, and sometimes a favourable performance for the parametric and the simpler forest that uses only the baseline information for the C-index criterion.

We now focus on the three forest methods, i.e., Separate, Poolt and Superpp. Summary tables that provide the ranking of these three methods for performance comparison using ADIST, ALOR and C-index for each factor separately are given in Section S1.3 in the Supplementary Material. In each situation, the Poolt method always ranks between Separate and Superpp, so we focus on the comparison between Separate and Superpp. Specifically, the comparison is carried out under two situations separately, i.e., when the estimation horizon  $(u - t)$  is equal to 1 and when it is larger than 1. We give  $T = 4$  as an example. In each

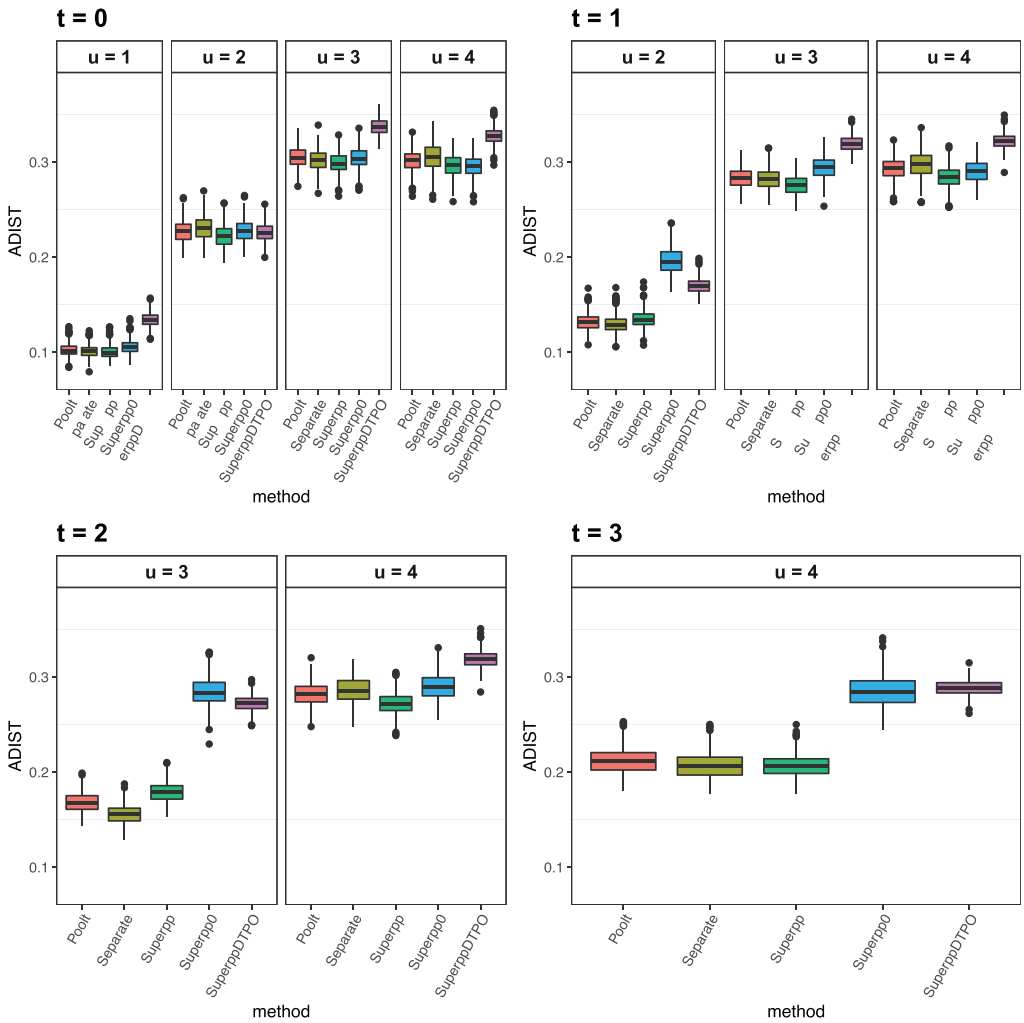


FIGURE 2: Simulation results comparing the distribution of ADIST on test sets across methods for each pair of  $(t, u)$ , trained on sample data of size 1000, with 10% censoring rate, generated following a Weibull distribution with an interaction survival relationship in the scenario 2TI+4TV, with high SNR and strong autocorrelation. The total number of time points is  $T = 4$ .

situation, using factorial designs, we study the difference of a given measure between Separate and Superpp under the effects of the following factors: autocorrelation, censoring rate, survival distribution, survival relationship, training sample size, scenario and SNR. The effects are estimated on the basis of an analysis of variance model fit with these factors as main effects.

Figures 4 and 5 provide the main effect plots for the difference between Separate and Superpp under all three measurements for  $(u - t) = 1$  and  $(u - t) > 1$ , respectively. In both cases, for each given effect, the general pattern of the change in difference resulting from varying the level of the effect is the same for ADIST and ALOR, and opposite for the C-index. Recall that low values of ADIST and ALOR and high values of C-index reflect better comparative performance of Separate over Superpp. Superpp is always the best performer for  $(u - t) > 1$ , although the effects

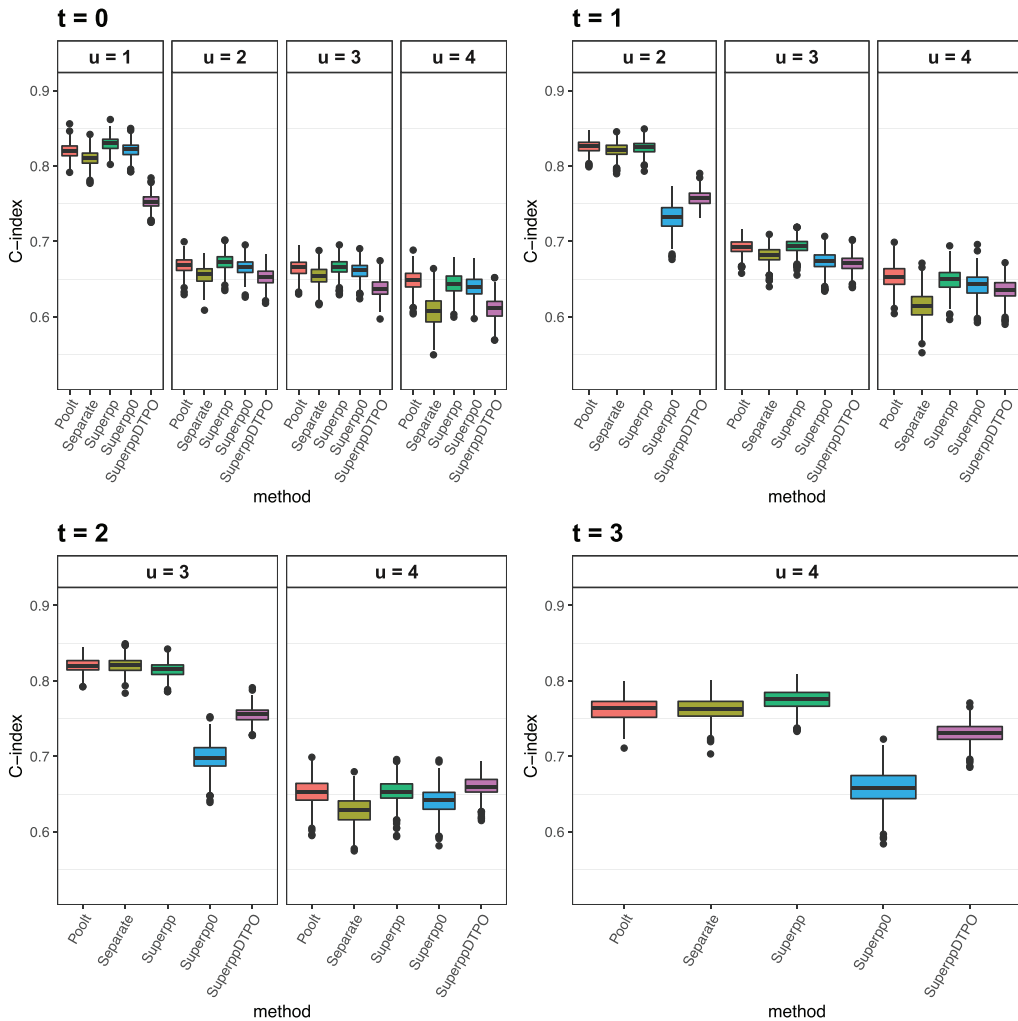


FIGURE 3: Simulation results comparing the distribution of C-index on test sets across methods for each pair of  $(t, u)$ , trained on sample data of size 1000, with 10% censoring rate, generated following a Weibull distribution with an interaction survival relationship in the scenario 2TI+4TV, with high SNR and strong autocorrelation. The total number of time points is  $T = 4$ .

are weaker than for  $(u - t) = 1$ , reflecting the difficulties of predicting farther in the future. We therefore focus on the estimation horizon  $(u - t) = 1$  in the following discussion.

We first examine the results based on ADIST. The overall centre of location is positive, highlighting that Superpp performs generally better than Separate. However, Separate can improve relative to Superpp with changes in factors. The larger the training sample size, the higher the SNR, or the smaller the censoring rate, the stronger the ability of any method to estimate the underlying survival relationship. In that situation, the flexibility of the Separate method is advantageous, while the stability of pooling is advantageous when the underlying relationship is more difficult to estimate. It is clear that the difference between the number of time-invariant (TI) and the number of time-varying (TV) covariates is driving the scenario effect.

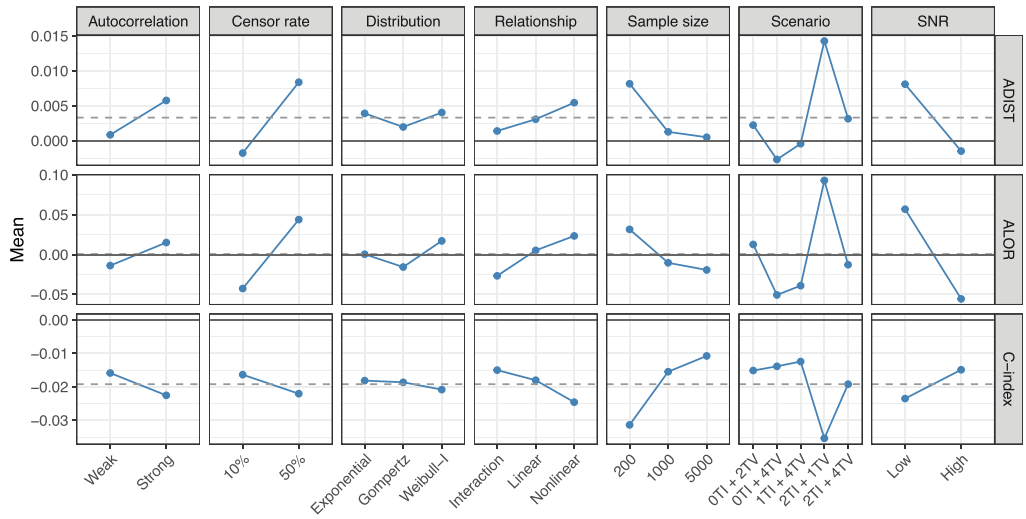


FIGURE 4: Main effect plot of difference for each measurement between Separate and Superpp method for the estimation horizon  $(u - t) = 1$ , that is, one-step-ahead estimation, when  $T = 4$ . Given any measurement  $m$ , the difference is computed as  $m_{\text{Separate}} - m_{\text{Superpp}}$ . The solid line gives the zero value and the dashed line gives the mean value over all effects for reference.

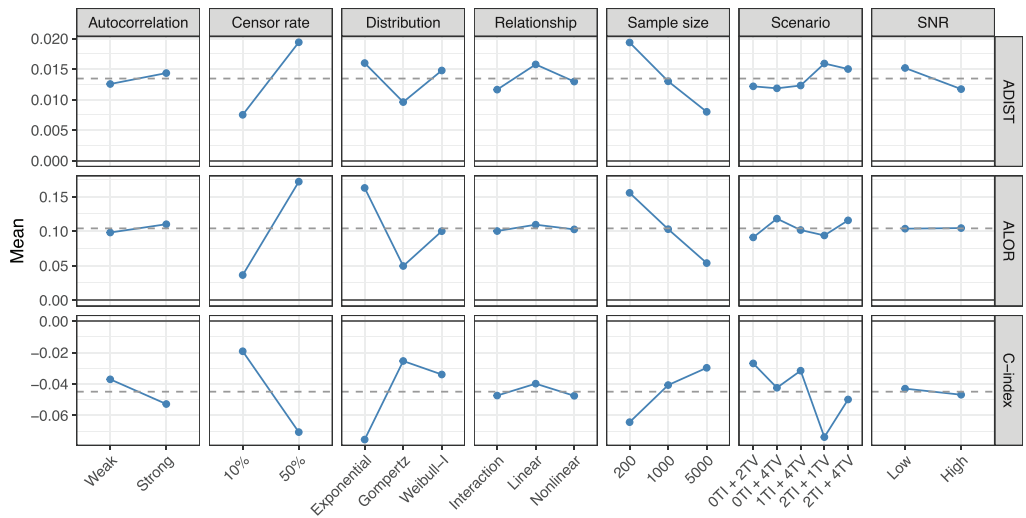


FIGURE 5: Main effect plot of difference for each measurement between Separate and Superpp method for the estimation horizon  $(u - t) > 1$ , that is, more than one-step-ahead estimation, when  $T = 4$ . Given any measurement  $m$ , the difference is computed as  $m_{\text{Separate}} - m_{\text{Superpp}}$ . The solid line gives the zero value and the dashed line gives the mean value over all effects for reference.

When  $\#TI - \#TV = 1$ , Superpp is the big winner; when  $\#TI - \#TV = -2$ , Superpp still wins, but by a smaller margin; when  $\#TI - \#TV = -3$ , Separate wins; and when  $\#TI - \#TV = -4$ , Separate wins by the largest margin. Presumably, this reflects that the Separate method is more sensitive to local time-varying effects, while pooling benefits from the stability associated with time-invariant covariates.

Separate performs better relative to Superpp when using ALOR as the measure of accuracy (sometimes beating it), reflecting that it can estimate extreme hazards more effectively. This is caused by the pooling underlying Superpp shrinking the estimated hazards away from the extremes; see the corresponding plot and discussion in Section S1.4 of the Supplementary Material.

The relative performance of Separate and Superpp using C-index is similar to that using ADIST, with Superpp being most effective. This may be explained by the fact that pooling reduces the variance and thus makes Superpp superior when we evaluate the performance with C-index.

Overall, weaker autocorrelation in covariates, higher censoring rate, smaller training sample size, smaller portion of covariates being time-varying, lower SNR and estimation further in the future, all reflect more difficult estimation tasks, and the less flexible but more stable pooling approach dominates. Conversely, in the opposite situations where signals are stronger and noise is less extreme, the more flexible but more variable Separate approach is more effective.

#### 4. CONCLUDING REMARKS

This article investigated different discrete-time survival forest methods for dynamic estimation with time-varying covariates. All methods investigated can be easily implemented using existing *R* packages. The results show that all methods perform well and none dominates the others. As a general rule, situations that are more difficult from an estimation point of view (such as weaker signals and less data) favour a global fit, pooling over all time points and taking advantage of reduced variance, while situations that are easier from an estimation point of view (such as stronger signals and more data) favour local fits, taking advantage of increased flexibility.

It should be noted that all the methods discussed here assume that censoring is uninformative; that is, subjects are censored for reasons unrelated to the time to event being examined. This is potentially an issue in the bankruptcy data examined in Section S2 in the Supplementary Material, as it is possible that companies that are in danger of declaring bankruptcy stop filing financial disclosures in order to hide their precarious financial position. A common parametric approach to this problem is the use of joint modelling, in which the assumed parametric forms for longitudinal predictors and the time to event are linked through shared random effects (Rizopoulos, 2012). It is possible that such models could be generalized to the discrete survival situation to allow tree-based structures on the joint distribution, perhaps based on recently developed tree-based methods for longitudinal data such as those described in Hajjem, Bellavance & Larocque (2011, 2014), Sela & Simonoff (2012) and Fu & Simonoff (2015).

In this article, we have limited ourselves to an event that is incomplete due only to right-censoring. Other reasons that the actual time to event is hidden are possible, such as left-truncation and interval censoring. Generalization of the methods discussed here would be useful future work, perhaps based on the tree- and forest-based methods for continuous time-to-event data discussed in Fu & Simonoff (2017a, 2017b) and Yao, Frydman & Simonoff (2021).

Presumably, all information of the time-varying covariates is available up to the given time for the estimation of the hazard function at a future time point. In this article, we implemented the forest methods based only on the current (latest) values of the time-varying covariates without including any lagged values. Future work can be done to investigate how to use the available lags efficiently, including the associated variable selection problems.

## ACKNOWLEDGEMENTS

We would like to thank the associate editor and three anonymous reviewers for their interesting and constructive comments that led to an improved version of this article. Denis Larocque acknowledges the financial support of The Natural Sciences and Engineering Research Council of Canada (NSERC) and Fondation HEC Montréal.

## DATA AVAILABILITY STATEMENT

The datasets generated and analyzed in the simulation study are available from the github repository, <https://github.com/ElainaYao/DynamicEstimationDTSD>, including R scripts for reproducibility of the simulations, and the Supplementary Material mentioned in the text.

## REFERENCES

- Anderson, J. R., Cain, K. C., & Gelber, R. D. (1983). Analysis of survival by tumor response. *Journal of Clinical Oncology*, 1, 710–719.
- Bacchetti, P. & Segal, M. R. (1995). Survival trees with time-dependent covariates: Application to estimating changes in the incubation period of AIDS. *Lifetime Data Analysis*, 1, 35–47.
- Berger, M., Welchowski, T., Schmitz-Valckenberg, S., & Schmid, M. (2019). A classification tree approach for the modeling of competing risks in discrete time. *Advances in Data Analysis and Classification*, 13, 965–990.
- Bertolet, M., Brooks, M. M., & Bittner, V. (2016). Tree-based identification of subgroups for time-varying covariate survival data. *Statistical Methods in Medical Research*, 25, 488–501.
- Bou-Hamad, I., Larocque, D., Ben-Ameur, H., Mâsse, L. C., Vitaro, F., & Tremblay, R. E. (2009). Discrete-time survival trees. *Canadian Journal of Statistics*, 37, 17–32.
- Bou-Hamad, I., Larocque, D., & Ben-Ameur, H. (2011a). Discrete-time survival trees and forests with time-varying covariates application to bankruptcy data. *Statistical Modelling*, 11, 429–446.
- Bou-Hamad, I., Larocque, D., & Ben-Ameur, H. (2011b). A review of survival trees. *Statistics Surveys*, 5, 44–71.
- Breiman, L. (2001). Random forests. *Machine Learning*, 45, 5–32.
- Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. (1984). *Classification and Regression Trees*. Wadsworth and Brooks, Monterey, CA.
- Cieslak, D. A., Hoens, T. R., Chawla N., & Kegelmeyer, W. P. (2012). Hellinger distance decision trees are robust and skew-insensitive. *Data Mining and Knowledge Discovery*, 24, 136–158.
- Elgmati, E., Fiaccone, R. L., Henderson, R., & Matthews, J. N. S. (2015). Penalised logistic regression and dynamic prediction for discrete-time recurrent event data. *Lifetime Data Analysis*, 21, 542–560.
- Fu, W. & Simonoff, J. S. (2015). Unbiased regression trees for longitudinal and clustered data. *Computational Statistics and Data Analysis*, 88, 53–74.
- Fu, W. & Simonoff, J. S. (2017a). Survival trees for interval-censored survival data. *Statistics in Medicine*, 36, 4831–4842.
- Fu, W. & Simonoff, J. S. (2017b). Survival trees for left-truncated and right-censored data, with application to time-varying covariate data. *Biostatistics*, 18, 352–369.
- Gordon, L. & Olshen, R. A. (1985). Tree-structured survival analysis. *Cancer Treatment Reports*, 69, 1065–1069.
- Hajjem, A., Bellavance, F., & Larocque, D. (2011). Mixed effects regression trees for clustered data. *Statistics and Probability Letters*, 81, 451–459.
- Hajjem, A., Bellavance, F., & Larocque, D. (2014). Mixed-effects random forest for clustered data. *Journal of Statistical Computation and Simulation*, 83, 1313–1328.
- Henderson, R., Diggle, P., & Dobson, A. (2000). Joint modelling of longitudinal measurements and event time data. *Biostatistics*, 1, 465–480.
- Hosmer, D. W., Lemeshow, S., & May, S. (2011). *Applied Survival Analysis: Regression Modeling of Time to Event Data*. Wiley, New York.
- Hothorn, T., Bühlmann, P., Dudoit, S., Molinaro, A., & Van Der Laan, M. J. (2006). Survival ensembles. *Biostatistics*, 7, 355–373.



- Ishwaran, H., Blackstone, E. H., Pothier, C., & Lauer, M. S. (2004). Relative risk forests for exercise heart rate recovery as a predictor of mortality. *Journal of the American Statistical Association*, 99, 591–600.
- Ishwaran, H., Kogalur, U. B., Blackstone, E. H., & Lauer, M. S. (2008). Random survival forests. *The Annals of Applied Statistics*, 2, 841–860.
- Klein, J. P. & Moeschberger, M. L. (2003). Survival analysis: Techniques for censored and truncated data. *Statistics for Biology and Health*, Springer, New York.
- Kleinbaum, D. G. & Klein, M. (2005). Survival analysis: A self-learning text. *Statistics for Biology and Health*, Springer, New York.
- Leblanc, M. & Crowley, J. (1993). Survival trees by goodness of split. *Journal of the American Statistical Association*, 88, 457–467.
- Madsen, E. B., Hougaard, P., & Gilpin, E. (1983). Dynamic evaluation of prognosis from time-dependent variables in acute myocardial infarction. *The American Journal of Cardiology*, 51, 1579–1583.
- R Core Team. (2020). *R: R Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Rizopoulos, D. (2012). Joint models for longitudinal and time-to-event data: With applications in R. *Chapman and Hall/CRC Biostatistics Series*, CRC Press, Boca Raton, CA.
- Schmid, M., Kchenhoff, H., Hoerauf, A., & Tutz, G. (2016). A survival tree method for the analysis of discrete event times in clinical and epidemiological studies. *Statistics in Medicine*, 35, 734–751.
- Schmid, M., Welchowski, T., Wright, M. N., & Berger, M. (2020). Discrete-time survival forests with Hellinger distance decision trees. *Data Mining and Knowledge Discovery*, 34, 812–832.
- Segal, M. R. (1988). Regression trees for censored data. *Biometrics*, 44, 35–47.
- Sela, R. J. & Simonoff, J. S. (2012). RE-EM trees: A data mining approach for longitudinal and clustered data. *Machine Learning*, 86, 169–207.
- Tutz, G. & Schmid, M. (2016). Modeling discrete time-to-event data. *Springer Series in Statistics*, Springer, Switzerland.
- van Houwelingen, H. (2007). Dynamic prediction by landmarking in event history analysis. *Scandinavian Journal of Statistics*, 34, 70–85.
- van Houwelingen, H. & Putter, H. (2011). *Dynamic Prediction in Clinical Survival Analysis*. CRC Press, Boca Raton, CA.
- Wang, P., Li, Y., & Reddy, C. K. (2019). Machine learning for survival analysis: A survey. *ACM Computing Surveys*, 51. Article 110.
- Willett, J. B. & Singer, J. D. (1993). Investigating onset, cessation, relapse, and recovery: Why you should, and how you can, use discrete-time survival analysis to examine event occurrence. *Journal of Consulting and Clinical Psychology*, 61, 952–965.
- Wongvibulsin, S., Wu, K. C., & Zeger, S. L. (2020). Clinical risk prediction with random forests for survival, longitudinal, and multivariate (RF-SLAM) data analysis. *BMC Medical Research Methodology*, 20, 1.
- Wright, M. N., Wager, S., & Probst, P. (2020). *Ranger: A Fast Implementation of Random Forests*. R package version 0.12.1.
- Yao, W., Frydman, H., Larocque, D., & Simonoff, J. S. (2020). Ensemble Methods for Survival Data with Time-Varying Covariates. arXiv preprint, arXiv:2006.00567.
- Yao, W., Frydman, H., & Simonoff, J. S. (2021). An ensemble method for interval-censored time-to-event data. *Biostatistics*, 22, 198–213.
- Zhu, R. & Kosorok, M. R. (2012). Recursively imputed survival trees. *Journal of the American Statistical Association*, 107, 331–340.

---

Received 28 July 2018

Accepted 22 February 2021