

The joint survival super learner: A super learner for right-censored data

BY A. MUNCH

Section of Biostatistics, University of Copenhagen
a.munch@sund.ku.dk

AND T. A. GERDS

Section of Biostatistics, University of Copenhagen

SUMMARY

Risk prediction models are widely used to guide real-world decision-making in areas such as healthcare and economics, and they also play a key role in estimating nuisance parameters in semiparametric inference. The super learner is a machine learning framework that combines a library of prediction algorithms into a meta-learner using cross-validated loss. In the context of right-censored data, careful consideration must be given to both the choice of loss function and the estimation of expected loss. Moreover, estimators such as inverse probability of censoring weighting (IPCW) require accurate modeling and ~~estimates an estimator~~ of the censoring distribution. We propose a novel approach to super learning for survival analysis that jointly evaluates candidate learners for both the event-time distribution and the censoring distribution. Our method imposes no restrictions on the algorithms included in the library, accommodates competing risks, and does not rely on a single pre-specified estimator of the censoring distribution. We establish ~~theoretical guarantees for our proposed method, including a finite-sample oracle inequality. In a simulation study, our super learner was able to better account for different censoring mechanisms than existing methods~~ bound on the average price we pay for using cross-validation, and show that this price vanishes asymptotically, up to poly-logarithmic terms, provided that the size of the library does not grow faster than at a polynomial rate in the sample size. We demonstrate the practical utility of our method using prostate cancer data and illustrate its performance in small samples with simulated data.

Some key words: competing risks; cross-validation; loss based estimation; right-censored data; super learner

1. INTRODUCTION

Accurately predicting risk from time-to-event data is a central challenge in ~~diverse various~~ research fields, ~~including such as~~ epidemiology, economics, and weather forecasting, with applications in clinical decision making and policy interventions. For instance, in prostate cancer management, clinicians often need to estimate a patient's risk of disease progression and mortality over time to make informed decisions about treatment strategies such as active surveillance versus immediate intervention. Reliable ~~time-to-event risk predictions~~ risk prediction models can help tailor care to individual patients, avoid overtreatment, and allocate healthcare resources more effectively. Super learning (van der Laan et al., 2007), also known as ensemble learning or stacked regression (Wolpert, 1992; Breiman, 1996), provides a powerful approach to this problem

by combining multiple candidate prediction models to reduce the risk of bias incurred by a single potentially misspecified model. In survival analysis, a super learner may [for example](#) combine a stack of Cox regression models with a stack of random survival forests (Gerds & Kattan, 2021, Section 8.4). Such a strategy has recently produced KDpredict (<https://kdpredict.com/>) a model which jointly predicts the risks of kidney failure and all-cause mortality at multiple time horizons based on different sets of covariates (Liu et al., 2024). To evaluate the prediction performance of the learners, the super learner behind KDpredict uses inverse probability of censoring weighting (IPCW), where the censoring distribution is estimated under the restrictive assumption that it does not depend on the covariates. This is a potential source of bias which is difficult to overcome with the currently available methods.

In this [paper/article](#), we propose the *joint survival super learner*, a new super learner designed to handle the specific challenges of ensemble learning with right-censored data. The joint survival super learner simultaneously learns prediction models for the event-time and censoring distributions. The joint survival super learner is based on ~~the simple idea of an artificial~~ a competing risks model [for the observed data](#), in which censoring is included as a state of its own: ~~Candidate learners~~, such that at any time it is known in which state an individual is. We assume [conditionally independent censoring and exploit well-known relationships between the observed data distribution on the one side and the partly unobserved distributions of the event time and the censoring time on the other](#). Learners for the event-time and censoring hazard functions are then assessed based on how well they predict the state occupation probabilities of the ~~artificial~~ competing risks model over time, given baseline covariate information. Our estimation framework ~~allows for~~ [thus naturally incorporates](#) competing risks, avoids restrictive assumptions on the censoring distribution, and [produces an estimator for the censoring distribution](#). Our approach is also fully flexible with respect to the choice of learners. The latter is in contrast to other proposals which restrict the library of learners to specific model classes (Polley & van der Laan, 2011; Golmakani & Polley, 2020), [see as we discuss in more detail in](#) Section 3.

To analyse the theoretical properties of the joint survival super learner, we focus on the discrete super learner, which selects the model in the library with the best estimated performance (van der Laan et al., 2007). We provide theoretical guarantees for the performance of the joint survival super learner, and in particular show that the discrete joint survival super learner is consistent ~~under natural conditions and satisfies~~, [when the library of learners includes at least one consistent learner](#). We also derive a finite-sample oracle inequality [for the discrete joint survival super learner](#). We demonstrate how to construct a library ~~using common survival models and how to obtain risk predictions from the resulting ensemble of~~ learners using common methods for survival analysis and illustrate the use of the joint survival super learner using a prostate cancer data set.

The ~~rest of the paper/article~~ is organized as follows. We introduce our notation and framework in Section 2. Section 3 introduces loss-based super learning and ~~presents~~ [discusses other](#) existing super learners for right-censored data. In Section ~~??~~ [4.1](#) we define the joint survival super learner, while Section 5 provides theoretical guarantees. Section 6 reports the results of [a series of](#) numerical experiments, and Section 7 illustrates the method on prostate cancer data. We conclude with a discussion in Section 8. Proofs are collected in the Appendix. Code and an implementation of the joint survival super learner [in R \(R Core Team, 2024\)](#) are available at <https://github.com/ammudn/joint-survival-super-learner>.

2. NOTATION AND FRAMEWORK

In a competing risks framework (Andersen et al., 2012) with J competing risks, let T be a time to event variable, ~~$D \in \{1, 2\}$~~ $D \in \{1, 2, \dots, J\}$ the cause of the event, and $X \in \mathcal{X}$ a vector of baseline covariates taking values in a bounded subset $\mathcal{X} \subset \mathbb{R}^p$, $p \in \mathbb{N}$. Let $\tau < \infty$ be a fixed prediction horizon. We use ~~\mathcal{Q}~~ \mathcal{Q} to denote the collection of all probability measures on ~~$\{0, \tau\} \times \{1, 2\} \times \mathcal{X}$~~ $[0, \tau] \times \{1, 2, \dots, J\} \times \mathcal{X}$ such that $(T, D, X) \sim Q$ for some unknown $Q \in \mathcal{Q}$. For ~~$j \in \{1, 2\}$~~ $j \in \{1, 2, \dots, J\}$, the cause-specific conditional cumulative hazard functions $\Lambda_j: [0, \tau] \times \mathcal{X} \rightarrow \mathbb{R}_+$ are defined as

$$\Lambda_j(t | x) = \int_0^t \frac{Q(T \in ds, D = j | X = x)}{Q(T \geq s | X = x)}.$$

For ease of presentation we assume throughout from now on that $J = 2$ and that the map $t \mapsto \Lambda_j(t | x)$ is continuous for all x and j , however, all technical arguments extend naturally to the general case (Andersen et al., 2012). The event-free survival function conditional on covariates is given by

$$S(t | x) = \exp \{-\Lambda_1(t | x) - \Lambda_2(t | x)\}. \quad (1)$$

Let \mathcal{M}_τ denote the space of all conditional cumulative hazard functions on $[0, \tau] \times \mathcal{X}$. Any distribution $Q \in \mathcal{Q}$ can be characterized by

$$Q(dt, j, dx) = \{S(t- | x)\Lambda_1(dt | x)H(dx)\}^{\mathbb{1}_{\{j=1\}}} \{S(t- | x)\Lambda_2(dt | x)H(dx)\}^{\mathbb{1}_{\{j=2\}}},$$

where $\Lambda_j \in \mathcal{M}_\tau$ for $j = 1, 2$ and H is the marginal distribution of the covariates.

We consider the right-censored setting in which we observe $O = (\tilde{T}, \tilde{D}, X)$, where $\tilde{T} = \min(T, C)$ for a right-censoring time C , $\Delta = \mathbb{1}\{T \leq C\}$, and $\tilde{D} = \Delta D$. Let \mathcal{P} denote a set of probability measures on the sample space $O = [0, \tau] \times \{0, 1, 2\} \times \mathcal{X}$ such that $O \sim P$ for some unknown $P \in \mathcal{P}$. We assume that the event times and the censoring times are conditionally independent given covariates, $T \perp\!\!\!\perp C | X$. This implies that any distribution $P \in \mathcal{P}$ is characterized by a distribution $Q \in \mathcal{Q}$ and a conditional cumulative hazard function for C given X (c.f., Begun et al., 1983; Gill et al., 1997). We use $\Gamma \in \mathcal{M}_\tau$ to denote the cumulative hazard function of the conditional censoring distribution given covariates. For ease of presentation we assume that $t \mapsto \Gamma(t | x)$ is continuous for all x . We let $(t, x) \mapsto G(t | x) = \exp \{-\Gamma(t | x)\}$ denote the survival function of the conditional censoring distribution. The distribution P is characterized by

$$\begin{aligned} P(dt, j, dx) &= \{G(t- | x)S(t- | x)\Lambda_1(dt | x)H(dx)\}^{\mathbb{1}_{\{j=1\}}} \\ &\quad \{G(t- | x)S(t- | x)\Lambda_2(dt | x)H(dx)\}^{\mathbb{1}_{\{j=2\}}} \\ &\quad \{G(t- | x)S(t- | x)\Gamma(dt | x)H(dx)\}^{\mathbb{1}_{\{j=0\}}} \\ &= \{G(t- | x)Q(dt, j, dx)\}^{\mathbb{1}_{\{j \neq 0\}}} \\ &\quad \{G(t- | x)S(t- | x)\Gamma(dt | x)H(dx)\}^{\mathbb{1}_{\{j=0\}}}. \end{aligned} \quad (2)$$

Hence, we may write $\mathcal{P} = \{P_{Q, \Gamma} : Q \in \mathcal{Q}, \Gamma \in \mathcal{G}\}$ for some $\mathcal{G} \subset \mathcal{M}_\tau$. We also have H -almost everywhere

$$P(\tilde{T} > t | X = x) = S(t | x)G(t | x) = \exp \{-\Lambda_1(t | x) - \Lambda_2(t | x) - \Gamma(t | x)\}.$$

We assume that there exists $\kappa < \infty$ such that $\Lambda_j(\tau- | x) < \kappa$, for $j \in \{1, 2\}$, and $\Gamma(\tau- | x) < \kappa$ for almost all $x \in \mathcal{X}$. This implies that $G(\tau- | x)$ is bounded away from zero for almost all $x \in \mathcal{X}$.

Under these assumptions, the conditional cumulative hazard functions Λ_j and Γ can be identified from P by

$$\Lambda_j(t | x) = \int_0^t \frac{P(\tilde{T} \in ds, \tilde{D} = j | X = x)}{P(\tilde{T} \geq s | X = x)}, \quad (3)$$

$$\Gamma(t | x) = \int_0^t \frac{P(\tilde{T} \in ds, \tilde{D} = 0 | X = x)}{P(\tilde{T} \geq s | X = x)}. \quad (4)$$

Thus, we can consider Λ_j and Γ as operators which map from \mathcal{P} to \mathcal{M}_τ .

3. LOSS-BASED SUPER LEARNING

Loss-based super learning requires a library of ~~candidate models (or learners)~~ learners, a cross-validation algorithm, and a loss function for evaluating predictive performance on hold-out samples. Let $\mathcal{D}_n = \{O_i\}_{i=1}^n \in \mathcal{O}^n$ be a data set of i.i.d. observations from $P \in \mathcal{P}$, and \mathcal{A} a collection of candidate learners. Let Θ be the parameter space, which in our case is a class of functions representing different models. Each learner $a \in \mathcal{A}$ is a map $a: \mathcal{O}^n \rightarrow \Theta$ which takes a data set as input and returns an estimate $a(\mathcal{D}_n) \in \Theta$. Let $L: \Theta \times \mathcal{O} \rightarrow \mathbb{R}_+$ be a loss function, representing the performance of the model $\theta \in \Theta$ at the observation $O \in \mathcal{O}$, where lower values mean better performance.

The expected loss of a learner is estimated by splitting the data set \mathcal{D}_n into K disjoint approximately equally sized subsets $\mathcal{D}_n^1, \mathcal{D}_n^2, \dots, \mathcal{D}_n^K$ and then calculating the cross-validated loss

$$\hat{R}_n(a; L) = \frac{1}{K} \sum_{k=1}^K \frac{1}{|\mathcal{D}_n^k|} \sum_{O_i \in \mathcal{D}_n^k} L(a(\mathcal{D}_n^{-k}), O_i), \quad \text{with } \mathcal{D}_n^{-k} = \mathcal{D}_n \setminus \mathcal{D}_n^k. \quad (5)$$

The subset \mathcal{D}_n^{-k} is referred to as the k 'th training sample, while \mathcal{D}_n^k is referred to as the k 'th test or hold-out sample. The discrete super learner is defined as

$$\hat{a}_n = \underset{a \in \mathcal{A}}{\operatorname{argmin}} \hat{R}_n(a; L),$$

and depends on both the library of learners and the specific partitioning of the data into cross-validation folds $\mathcal{D}_n^1, \dots, \mathcal{D}_n^K$.

When designing a super learner for right-censored data, particular care must be taken in the choice of loss function and in the estimation of the expected loss. A commonly used loss function for right-censored data is the partial log-likelihood loss (e.g., Li et al., 2016; Yao et al., 2017; Lee et al., 2018; Katzman et al., 2018; Gensheimer & Narasimhan, 2019; Lee et al., 2021; Kvamme & Borgan, 2021). This loss function is also recommended for super learning with right-censored data by Polley & van der Laan (2011), under the assumption that data are observed in discrete time. However, the partial log-likelihood loss does not work well as a general purpose measure of performance in hold-out samples when data are observed in continuous time. The ~~is because~~ it reason is that the partial log-likelihood assigns an infinite value to any learner that predicts piecewise constant cumulative hazard functions, if the test set contains event times that ~~were~~ are not observed in the training set. This problem occurs with prominent survival learners including the Kaplan-Meier estimator, random survival forests, and semi-parametric Cox regression models, and these learners cannot be included in the library of the super learner proposed by Polley & van der Laan (2011). When a proportional hazards model is assumed, the baseline hazard function can be profiled out of the likelihood (Cox, 1972). The cross-validated partial log-likelihood loss (Verweij & van Houwelingen, 1993) has therefore been suggested as a loss function for super

learning by Golmakani & Polley (2020). ~~This~~ However, this choice of loss function restricts the library of learners to include only Cox proportional hazards models, and hence excludes many learners such as, e.g., random survival forests, additive hazards models, and accelerated failure time models. 150

Alternative approaches for super learning with right-censored data use an inverse probability of censoring weighted (IPCW) loss function (Graf et al., 1999; van der Laan & Dudoit, 2003; Molinaro et al., 2004; Keles et al., 2004; Hothorn et al., 2006; Gerds & Schumacher, 2006; Gonzalez Ginestet et al., 2021), censoring unbiased transformations (Fan & Gijbels, 1996; Steingrimsson et al., 2019), or pseudo-values (Andersen et al., 2003; Mogensen & Gerds, 2013; Sachs et al., 2019). All these methods rely on an estimator of the censoring distribution, and their drawback is that this estimator has to be pre-specified. Recent work by Han et al. (2021) and Westling et al. (2021) circumvents the need to pre-specify a censoring model by iterating between estimation of the outcome and censoring models. However, this iterative procedure is in general not guaranteed to converge to the true data-generating mechanism (Munch, 2024, Appendix A.4). 155
160

4. THE JOINT SURVIVAL SUPER LEARNER

4.1. *An artificial competing risks model*

The main idea of the joint survival super learner is to ~~jointly use learners for specify libraries of learners for the hazard functions~~ Λ_1 , Λ_2 , and Γ , and ~~to exploit the relations in equation (2) ; to learn a feature of the observed data distribution P . The discrete to define a joint loss function.~~ The joint survival super learner ranks thus evaluates a tuple of learners for $(\Lambda_1, \Lambda_2, \Gamma)$ based on how well they jointly ~~model-predict~~ the observed data and the discrete joint survival super learner chooses the best performing tuple. To formally introduce the joint survival super learner, we define the process 165

$$\eta(t) = \mathbb{1}\{\tilde{T} \leq t, \tilde{D} = 1\} + 2 \mathbb{1}\{\tilde{T} \leq t, \tilde{D} = 2\} - \mathbb{1}\{\tilde{T} \leq t, \tilde{D} = 0\}, \quad \text{for } t \in [0, \tau],$$

which takes values in $\{-1, 0, 1, 2\}$. The four values represent four mutually exclusive states. Specifically, value 0 represents the state where the individual is still event-free and uncensored, value 1 the state where the event of interest has occurred, value 2 the state where a competing risk has occurred, and value -1 the state where the observation is right-censored. The state occupation probabilities given baseline covariates X are given by the function 170

$$F(t, l, x) = P(\eta(t) = l \mid X = x), \quad (6)$$

for all $t \in [0, \tau]$, $l \in \{-1, 0, 1, 2\}$, and $x \in \mathcal{X}$. 175

The joint survival super learner is a super learner for the function-valued parameter $\theta(P) = F$ which is identified through equation (6). Under conditional independent censoring, each tuple $(\Lambda_1, \Lambda_2, \Gamma, H)$ characterizes a distribution $P \in \mathcal{P}$, c.f. equation (2), which in turn determines (F, H) . Hence, a learner for F can be constructed from learners for Λ_1 , Λ_2 , and Γ as follows:

$$\begin{aligned} F(t, 0, x) &= P(\tilde{T} > t \mid X = x) = \exp\{-\Lambda_1(t \mid x) - \Lambda_2(t \mid x) - \Gamma(t \mid x)\}, \\ F(t, 1, x) &= P(\tilde{T} \leq t, \tilde{D} = 1 \mid X = x) = \int_0^t F(s-, 0, x) \Lambda_1(ds \mid x), \\ F(t, 2, x) &= P(\tilde{T} \leq t, \tilde{D} = 2 \mid X = x) = \int_0^t F(s-, 0, x) \Lambda_2(ds \mid x), \\ F(t, -1, x) &= P(\tilde{T} \leq t, \tilde{D} = 0 \mid X = x) = \int_0^t F(s-, 0, x) \Gamma(ds \mid x). \end{aligned} \quad (7)$$

Equation (7) implies that a library for the joint survival super learner can be built from three libraries of learners: \mathcal{A}_1 , \mathcal{A}_2 , and \mathcal{B} , where \mathcal{A}_1 and \mathcal{A}_2 contain learners for the conditional cause-specific cumulative hazard functions Λ_1 and Λ_2 , respectively, and \mathcal{B} contains learners for the conditional cumulative hazard function of the censoring distribution. Taking the Cartesian product of these libraries, we obtain a library \mathcal{F} of learners for F :

$$\mathcal{F}(\mathcal{A}_1, \mathcal{A}_2, \mathcal{B}) = \{\varphi_{a_1, a_2, b} : a_1 \in \mathcal{A}_1, a_2 \in \mathcal{A}_2, b \in \mathcal{B}\}, \quad (8)$$

where in correspondence with the relations in equation (7),

$$\varphi_{a_1, a_2, b}(\mathcal{D}_n)(t, 0, x) = \exp \{-a_1(\mathcal{D}_n)(s | x) - a_2(\mathcal{D}_n)(s | x) - b(\mathcal{D}_n)(s | x)\},$$

$$\varphi_{a_1, a_2, b}(\mathcal{D}_n)(t, 1, x) = \int_0^t \varphi_{a_1, a_2, b}(\mathcal{D}_n)(s-, 0, x) a_1(\mathcal{D}_n)(ds | x),$$

$$\varphi_{a_1, a_2, b}(\mathcal{D}_n)(t, 2, x) = \int_0^t \varphi_{a_1, a_2, b}(\mathcal{D}_n)(s-, 0, x) a_2(\mathcal{D}_n)(ds | x),$$

$$\varphi_{a_1, a_2, b}(\mathcal{D}_n)(t, -1, x) = \int_0^t \varphi_{a_1, a_2, b}(\mathcal{D}_n)(s-, 0, x) b(\mathcal{D}_n)(ds | x).$$

$$\varphi_{a_1, a_2, b}(\mathcal{D}_n)(t, 0, x) = \exp \{-a_1(\mathcal{D}_n)(s | x) - a_2(\mathcal{D}_n)(s | x) - b(\mathcal{D}_n)(s | x)\},$$

$$\varphi_{a_1, a_2, b}(\mathcal{D}_n)(t, 1, x) = \int_0^t \varphi_{a_1, a_2, b}(\mathcal{D}_n)(s-, 0, x) a_1(\mathcal{D}_n)(ds | x),$$

$$\varphi_{a_1, a_2, b}(\mathcal{D}_n)(t, 2, x) = \int_0^t \varphi_{a_1, a_2, b}(\mathcal{D}_n)(s-, 0, x) a_2(\mathcal{D}_n)(ds | x),$$

$$\varphi_{a_1, a_2, b}(\mathcal{D}_n)(t, -1, x) = \int_0^t \varphi_{a_1, a_2, b}(\mathcal{D}_n)(s-, 0, x) b(\mathcal{D}_n)(ds | x).$$

Notably, the libraries \mathcal{A}_1 , \mathcal{A}_2 , and \mathcal{B} can be constructed using standard software for survival analysis. ~~In For example, in the R, for instance, we can construct Cox models as learners software we can specify various ways to include covariates in a Cox regression model and fit learners of the hazard functions~~ using the survival-package (Therneau, 2022), and we can ~~construct random survival forests as learners specify hyper parameters of a random survival forest and derive learners of the hazard functions~~ using the randomForestSRC-package (Ishwaran & Kogalur, 2025).

To evaluate how well a function F predicts the process η we use the integrated Brier score (Graf et al., 1999) $\bar{B}_\tau(F, O) = \int_0^\tau B_t(F, O) dt$, where B_t is the Brier score (Brier et al., 1950) at time $t \in [0, \tau]$,

$$B_t(F, O) = \sum_{l=1}^2 (F(t, l, X) - 1\{\eta(t) = l\})^2.$$

The Brier score is ~~here the the average~~ squared prediction error across ~~all of~~ the four states. Based on a split of a data set \mathcal{D}_n into K disjoint approximately equally sized subsets (c.f., Section 3),

each learner $\varphi_{a_1, a_2, b}$ in the library $\mathcal{F}(\mathcal{A}_1, \mathcal{A}_2, \mathcal{B})$ is evaluated using the cross-validated loss,

$$\hat{R}_n(\varphi_{a_1, a_2, b}; \bar{B}_\tau) = \frac{1}{K} \sum_{k=1}^K \frac{1}{|\mathcal{D}_n^k|} \sum_{O_i \in \mathcal{D}_n^k} \bar{B}_\tau(\varphi_{a_1, a_2, b}(\mathcal{D}_n^{-k}), O_i),$$

and the discrete joint survival super learner is

$$\hat{\varphi}_n = \underset{(a_1, a_2, b) \in \mathcal{A}_1 \times \mathcal{A}_2 \times \mathcal{B}}{\operatorname{argmin}} \hat{R}_n(\varphi_{a_1, a_2, b}; \bar{B}_\tau).$$

205

best performing tuple of hazard functions:

$$(\hat{\Lambda}_{1n}, \hat{\Lambda}_{2n}, \hat{\Gamma}_n) = \underset{(a_1, a_2, b) \in \mathcal{A}_1 \times \mathcal{A}_2 \times \mathcal{B}}{\operatorname{argmin}} \hat{R}_n(\varphi_{a_1, a_2, b}; \bar{B}_\tau). \quad (10)$$

4.1. Obtaining risk predictions

The joint survival super learner estimates the function F which depends on the censoring distribution and is therefore typically not of direct interest in itself. For instance, in the prostate cancer example we consider in Sections 6 and 7, the value $F(t, l, x)$ denotes the conditional probability that a patient with baseline characteristics $X = x$ will, before time point t : have had tumour recurrence before leaving the study ($l = 1$); have died without tumour recurrence before leaving the study ($l = 2$); have left the study (state $l = -1$); or be alive and part of the study without tumour recurrence ($l = 0$). The reference to still being part of the study is irrelevant from a new patient's perspective. We here demonstrate how clinically relevant risk predictions can instead be obtained from the joint survival super learner.

210

215

We recall that we work under Under the assumption of conditional independent censoring and positivity, as introduced stated in Section 2. Under these assumptions, it follows by from equations (3) and (4) and the definition of F that

$$\Lambda_j(t, x) = \int_0^t \frac{F(ds, j, x)}{F(s-, 0, x)}, \quad j \in \{1, 2\}. \quad (11)$$

Cause-specific risk predictions can be obtained from Λ_1 and Λ_2 using the formula (e.g., Benichou & Gail, 1990; Ozenne et al., 2017),

220

$$Q(T \leq t, D = j \mid X = x) = \int_0^t \exp \{-\Lambda_1(u \mid x) - \Lambda_2(u \mid x)\} \Lambda_j(du \mid x), \quad j \in \{1, 2\}.$$

Hence, given the joint survival super learner's estimate of F we can use equation to obtain estimates of the cause-specific cumulative hazard functions Λ_j , which can in turn be used to obtain estimates of the cause-specific risks through equation. For instance, in the prostate cancer example, this expression will provide the conditional probability that a patient with a certain set of baseline characteristics will, before time point t , have had tumour recurrence or have died without tumour recurrence.

225

We have suggested to implement (10) by substituting into the well-known formula (e.g., Benichou & Gail, 1990; Ozenne et al., 2017),

$$\hat{Q}_n(T \leq t, D = j \mid X = x) = \int_0^t \exp \{-\hat{\Lambda}_{1n}(u \mid x) - \hat{\Lambda}_{2n}(u \mid x)\} \hat{\Lambda}_{jn}(du \mid x), \quad j \in \{1, 2\}. \quad (12)$$

Furthermore, the joint survival super learner by building a library using learners of the cause-specific cumulative hazard functions, A_j , and the cumulative hazard function for censoring, Γ . With this implementation, we can directly input the highest ranked tuple of cause-specific hazard functions (A_1, A_2) provided by the joint survival super learner as input to equation provides an estimator of the censoring distribution:

$$\hat{G}_n(T \leq t \mid X = x) = \exp \{ -\hat{\Gamma}_n(t \mid x) \}.$$

5. THEORETICAL GUARANTEES

The use of cross-validation underlying the joint survival super learner is Cross-validation is the backbone of super learning and an intuitively reasonable procedure for fair model selection without overfitting. In this section, we follow the works adapt the work of van der Laan & Dudoit (2003) and van der Vaart et al. (2006) to and provide a theoretical justification for this practice the joint survival super learner in the form of a finite-sample oracle inequality. We begin by demonstrating that minimizing the integrated Brier score, as defined in Section 2.4.1, is statistically meaningful, in that perfect proper, in the sense that minimisation recovers the parameter of the data-generating distribution. Together with our finite-sample oracle inequality (Proposition 2 below), this implies that the joint survival super learner is consistent when it is based on a library that includes a at least one consistent learner. Another consequence of our finite-sample oracle inequality is that the joint survival super learner converges at (nearly) at the optimal rate achievable within the library of learners. This statement is made precise in Corollary 1 and the following discussion. Proofs are deferred to the Appendix.

A sensible loss function should attain the minimal expected value at the parameter corresponding to the data-generating distribution. Loss functions with this property are known as proper scoring rules, and as strictly proper scoring rules called proper, and strictly proper if the minimizer is unique (Gneiting & Raftery, 2007). Absence of properness makes it unclear why minimizing the (estimated) expected loss is interesting. Proposition 1 is a formal statement of the fact states that the integrated Brier score, as defined in our setting (e.f., Section 2.4.1, Section 4.1) is a strictly proper scoring rule. To state establish this result, recall that the function F implicitly depends on the data-generating probability measure $P \in \mathcal{P}$ but that this was so far suppressed in the notation. We now make this dependence explicit by writing F_P for the function determined by a given $P \in \mathcal{P}$ in accordance with equation (6). In the following we let $\mathcal{H}_{\mathcal{P}} = \{F_P : P \in \mathcal{P}\}$.

PROPOSITION 1. If $P \in \mathcal{P}$ then

$$F_P = \operatorname{argmin}_{F \in \mathcal{H}_{\mathcal{P}}} \mathbb{E}_P [\bar{B}_{\tau}(F, O)],$$

for all $l \in \{-1, 0, 1, 2\}$, almost all $t \in [0, \tau]$, and P -almost all $x \in \mathcal{X}$.

The discrete joint survival super learner defined in (10) provides an estimate of the function F

$$\hat{\varphi}_n = \varphi_{\hat{\Lambda}_{1n}, \hat{\Lambda}_{2n}, \hat{\Gamma}_n}$$

which is obtained by substituting $(\hat{\Lambda}_{1n}, \hat{\Lambda}_{2n}, \hat{\Gamma}_n)$ for (a_1, a_2, b) into the structural equations (9). To evaluate the performance of the joint survival super learner we might benchmark it against the data-generating distribution F_P , as this has which according to Proposition 1 has the smallest expected loss by Proposition 1. A more nuanced comparison is to benchmark it against

. Another useful theoretical benchmark is the so-called oracle learner which is the best learner available given the library and the training data. This is the so-called oracle learner, formally defined as included in the library of learners and formally defined by

$$\tilde{\varphi}_n = \operatorname{argmin}_{\varphi \in \mathcal{F}(\mathcal{A}_1, \mathcal{A}_2, \mathcal{B})} \tilde{R}_n(\varphi; \bar{B}_\tau), \quad \text{with} \quad \tilde{R}_n(\varphi; \bar{B}_\tau) = \frac{1}{K} \sum_{k=1}^K \mathbb{E}_P \left[\bar{B}_\tau(\varphi(\mathcal{D}_n^{-k}), O) \mid \mathcal{D}_n^{-k} \right],$$

where we use \mathbb{E}_P to denote the expectation under the distribution P for a new observation O independent of which is independent of the data \mathcal{D}_n^{-k} . Like the joint survival super learner, the oracle learner depends on the library of learners and on the specific datapartitions actual partition of the data, but unlike the joint survival super learner, it also depends on the unknown data-generating distribution. It is hence not available in practice and serves only as a theoretical benchmark.

In the following, we equip the space \mathcal{H}_P with the norm

$$\|F\|_P = \left\{ \sum_{l=1}^2 \int_0^\tau \mathbb{E}_P [F(t, l, X)^2] dt \right\}^{1/2}. \quad (13)$$

This norm is a natural performance which is equal to the excess risk, $\mathbb{E}_P [\bar{B}_\tau(F, O)] - \mathbb{E}_P [\bar{B}_\tau(F_P, O)]$ by, as shown in Lemma 1 in the Appendix and is a natural performance measure. For simplicity of presentation, we take n and the data partitions to be such that assume that all folds of the data partition have equal size, $|\mathcal{D}_n^{-k}| = n/K$ with K fixed. We will for a fixed number of folds K . We allow the number of learners to grow with n and write $\mathcal{F}_n = \mathcal{F}(\mathcal{A}_{1,n}, \mathcal{A}_{2,n}, \mathcal{B}_n)$ as short-hand notation emphasizing the dependence on the sample size. We now state a finite-sample inequality that bounds the performance of the joint survival super learner relative to that of the oracle learner.

PROPOSITION 2. For all $P \in \mathcal{P}$, $n \in \mathbb{N}$, $k \in \{1, \dots, K\}$, and $\delta > 0$,

$$\begin{aligned} \frac{1}{K} \sum_{k=1}^K \mathbb{E}_P [\|\hat{\varphi}_n(\mathcal{D}_n^{-k}) - F_P\|_P^2] &\leq (1 + 2\delta) \frac{1}{K} \sum_{k=1}^K \mathbb{E}_P [\|\tilde{\varphi}_n(\mathcal{D}_n^{-k}) - F_P\|_P^2] \\ &\quad + (1 + \delta) 16K\tau \left(13 + \frac{12}{\delta} \right) \frac{\log(1 + |\mathcal{F}_n|)}{n}. \end{aligned}$$

The expectations Note that the expectation in Proposition 2 reflect a mild abuse of notation, in that they are formally is taken with respect to the product measure P^n for the whole data set \mathcal{D}_n . This means that we are quantifying the average performance of the joint survival super learner across average training data all training data of size n . A corresponding quantity was called the expected true error rate in Efron & Tibshirani (1997). As with many finite-sample oracle inequalities, this result is of little direct practical utility because the right-hand side depends on data-dependent, unknown quantities. However, it does quantify how the number of folds, the time horizon, and the number of learners in the library can be expected to influence the performance. The result has the following asymptotic consequences.

COROLLARY 1. Assume that $|\mathcal{F}_n| = O(n^q)$, for some $q \in \mathbb{N}$ and that there exists a sequence $\varphi_n \in \mathcal{F}_n$, $n \in \mathbb{N}$, such that $\mathbb{E}_P [\|\varphi_n(\mathcal{D}_n^{-k}) - F_P\|_P^2] = C_P + O(n^{-\alpha})$, for some $\alpha \leq 1$ and $C_P \geq 0$.

- (a) If $\alpha = 1$, then $\frac{1}{K} \sum_{k=1}^K \mathbb{E}_P [\|\hat{\varphi}_n(\mathcal{D}_n^{-k}) - F_P\|_P^2] = C_P + O(\log(n)^{1+\varepsilon} n^{-1})$, $\forall \varepsilon > 0$.
- (b) If $\alpha < 1$, then $\frac{1}{K} \sum_{k=1}^K \mathbb{E}_P [\|\hat{\varphi}_n(\mathcal{D}_n^{-k}) - F_P\|_P^2] = C_P + O(n^{-\alpha})$.

Proposition 2 ~~provided a precise~~ thus provides a finite-sample bound on the average price we pay for using cross-validation, and Corollary 1 states that this price vanishes asymptotically, up to poly-logarithmic terms, provided that the size of the library does not grow faster than at a polynomial rate in the sample size. The situation $C_P = 0$ corresponds to a setting in which the library includes a consistent learner. Cases (a) and (b) correspond to situations where the oracle learner achieves a parametric or non-parametric asymptotic rate of convergence, respectively.

To illustrate the content of Corollary 1, consider first a situation where we use a library with an increasing number of semi-parametric Cox models with different interaction terms, as well as several Poisson regression models based on different discretisations of the time scale. Each of these models will independently achieve a parametric rate of convergence, and hence item (a) of Corollary 1 states that the joint survival super learner based on this library will achieve a near-parametric rate of convergence. The constant C_P can be taken equal to the distance to the least false model in the library, and so the joint survival super learner will approximate the least false model in the library at a near-parametric rate. Another situation appears if we add more flexible models to the library, such as Cox lasso or random survival forests. These models typically converge at non-parametric rates, with the fastest rate depending on the unknown data-generating distribution. Item (b) of Corollary 1 shows that the joint survival super learner achieves the same convergence rate as the best-performing algorithm in the library, without any knowledge of the data-generating distribution.

6. NUMERICAL EXPERIMENTS

In this section, we report results from a simulation study where we ~~consider estimation of~~ estimate the conditional survival function. ~~We restrict attention to the simpler setting without competing risks to enable~~ In order to enable a comparison with existing super learner methods ~~we restrict the attention to a setting without competing risks. We report the results of the numerical experiments in two parts.~~ In the first part, we compare the joint survival super learner ~~to~~ with two IPCW based discrete super learners that use either the Kaplan-Meier estimator or a Cox model to estimate the censoring probability (Gonzalez Ginestet et al., 2021). In the second part, we compare the joint survival super learner to the super learner proposed by Westling et al. (2021).

~~In both parts, we use the same data-generating mechanism. We generate data according to a distribution motivated from a real data set in which censoring depends on the baseline covariates. We simulate data based on~~ For our numerical experiments we have synthesized the data of the prostate cancer study of Kattan et al. (2000) based on structural equations where we specify Cox-Weibull regression models Bender et al. (2005) for the latent event and censoring times. The outcome of interest is the time to tumour recurrence, and five baseline covariates are used to predict outcome: prostate-specific antigen (PSA, ng/mL), Gleason score sum (GSS, values between 6 and 10), radiation dose (RD), hormone therapy (HT, yes/no) and clinical stage (CS, six values). The study was designed such that a patient’s radiation dose depended on when the patient entered the study (Gerds et al., 2013). This in turn implies that the time of censoring depends on the radiation dose. The data were re-analysed in (Gerds et al., 2013) where a sensitivity analysis was conducted based on simulated data. Here we use the same simulation setup, where event and censoring times are generated according to parametric Cox-Weibull models estimated from the original data, and the covariates are generated according to either marginal Gaussian normal or binomial distributions estimated from the original data (c.f., Gerds et al., 2013, Section 4.6). We refer to this simulation setting as ‘dependent censoring’. We also considered a simulation setting where data were generated in the same way, except that censoring was generated completely independently. We refer to this simulation setting as ‘independent censoring’.

For all super learners, we use a library consisting of three learners: The Kaplan-Meier estimator (Kaplan & Meier, 1958; Gerds, 2019), a Cox model with main effects (Cox, 1972; Therneau, 2022), and a random survival forest (Ishwaran et al., 2008; Ishwaran & Kogalur, 2025). We use the same library to learn the outcome distribution and the censoring distribution. The three learners in our library of learners can be used to learn the cumulative hazard functions of the outcome and the censoring distribution. The latter works by training the learner on the data set \mathcal{D}_n^c , where $\mathcal{D}_n^c = \{O_i^c\}_{i=1}^n$ with $O_i^c = (\tilde{T}_i, 1 - \Delta_i, X_i)$. When we say that we use a learner for the cumulative hazard function of the outcome to learn the cumulative hazard function of the censoring time, we mean that the learner is trained on \mathcal{D}_n^c .

We compare the joint survival super learner to two IPCW based super learners: The first super learner, called IPCW(Cox), uses a Cox model with main effects to estimate the censoring probabilities, while the second super learner, called IPCW(KM), uses the Kaplan-Meier estimator to estimate the censoring probabilities. The Cox model for the censoring distribution is correctly specified in both simulation settings while the Kaplan Meier estimator only estimates the censoring model correctly in the simulation setting where censoring is independent. Both IPCW super learners are fitted using the R-package `riskRegression` (Gerds et al., 2023). The IPCW super learners use the integrated Brier score up to a fixed time horizon (36 months). The marginal risk of the event before this time horizon is $\approx 24.6\%$. Under the ‘dependent censoring’ setting the marginal censoring probability before the time horizon is $\approx 61.9\%$. Under the ‘independent censoring’ setting the marginal censoring probability before this time horizon is $\approx 38.7\%$.

Each super learner provides a learner for the cumulative hazard function for the outcome of interest. From the cumulative hazard function, we obtain a risk prediction model as described in Section 4.1, with the special case of $\Lambda_2 = 0$. We measure the performance of each super learner by calculating the index of prediction accuracy (IPA) (Kattan & Gerds, 2018) at a fixed time horizon (36 months) for the risk prediction model provided by the super learner. The IPA is 1 minus the ratio between the model’s Brier score and the null model’s Brier score, where the null model is the model that does not use any covariate information. The IPA is approximated using a large ($n = 20,000$) independent data set of uncensored data. As a benchmark, we calculate the performance of the risk prediction model chosen by the oracle selector, which uses the large data set of uncensored event times to select the learner with the highest IPA.

The results are shown in Figure 1. We see that in the scenario where censoring depends on the covariates, using the Kaplan-Meier estimator to estimate the censoring probabilities provides a risk prediction model with an IPA that is lower than the risk prediction model provided by the joint survival super learner. The performance of the risk prediction model selected by the joint survival super learner is similar to the risk prediction model selected by the IPCW(Cox) super learner which a priori uses a correctly specified model for the censoring distribution. Both these risk prediction models are close to the performance of the oracle, except for small sample sizes.

We next compare the joint survival super learner to the super learner `survSL` (Westling et al., 2021). This is another super learner which like the joint survival super learner works without a pre-specified censoring model. Both the joint survival super learner and `survSL` provide estimates of the event-time and censoring distributions. Hence, we compare the performance of these methods with respect to both the outcome and the censoring distribution. Again we use the IPA to quantify the predictive performance.

The results are shown in Figures 2 and 3. We see that for most sample sizes, the joint survival super learner selects prediction models for both censoring and outcome which have similar or higher IPA compared to the prediction models selected by `survSL`.

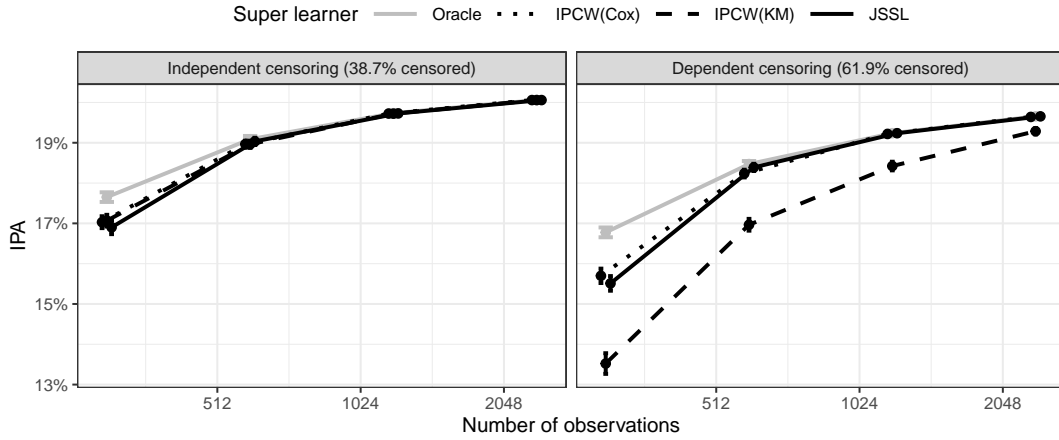


Fig. 1. For the risk prediction models provided by each of the super learners, the IPA is plotted against sample size. The results are averages across 1000 simulated data sets and the error bars are used to quantify the Monte Carlo uncertainty. JSSL denotes the joint survival super learner.

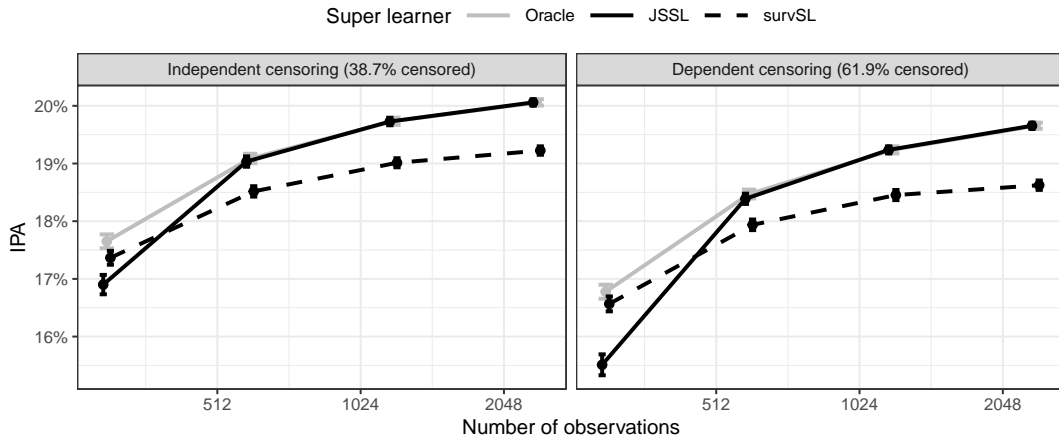


Fig. 2. For the risk prediction models of the outcome provided by each of the super learners, the IPA at the fixed time horizon is plotted against sample size. The results are averages across 1000 repetitions and the error bars are used to quantify the Monte Carlo uncertainty. JSSL denotes the joint survival super learner.

7. PROSTATE CANCER STUDY

We use the prostate cancer data of Kattan et al. (2000) to illustrate the use of the joint survival
 395 super learner in the presence of competing risks. We have introduced the data in Section 6. The
 data consists of 1,042 patients who are followed from start of followup until tumour recurrence,
 death without tumour recurrence, or end of followup (censored), whatever came first. We use the
 joint survival super learner to rank libraries of learners for the cause-specific cumulative hazard
 functions of tumour recurrence, death without tumour recurrence, and censoring. The libraries
 400 of learners each include five learners: the Nelson-Aalen estimator, three Cox regression models
 (unpenalized, lasso, elastic net) each including additive effects of the baseline covariates, and a
 random survival forest. We use the same set of learners to estimate the cumulative hazard function
 of tumour recurrence Λ_1 , the cumulative hazard function of death without tumour recurrence Λ_2 ,

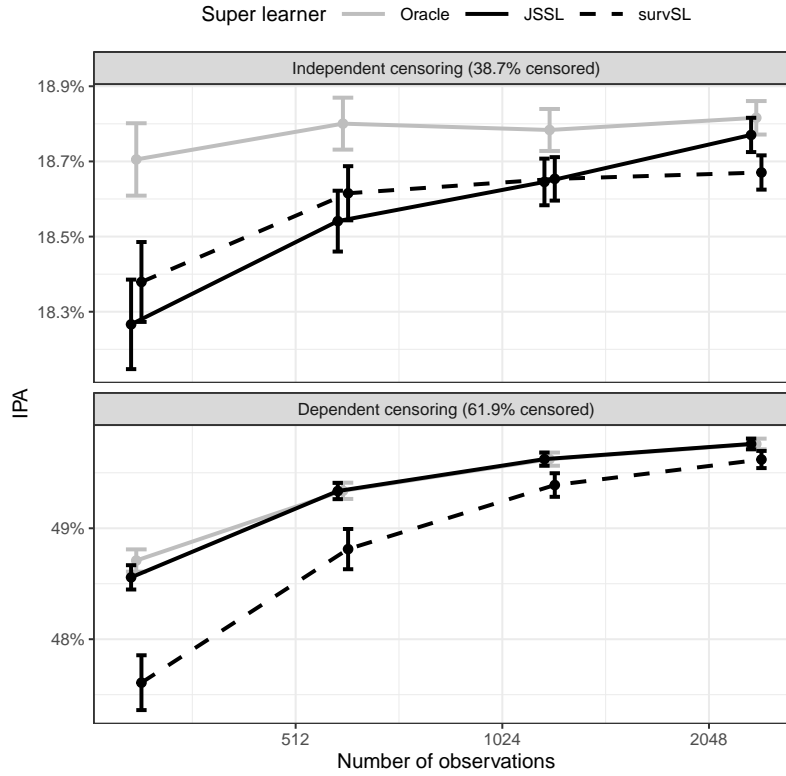


Fig. 3. For the risk prediction models of the censoring model provided by each of the super learners, the IPA at the fixed time horizon is plotted against sample size. The results are averages across 1000 repetitions and the error bars are used to quantify the Monte Carlo uncertainty. JSSL denotes the joint survival super learner.

and the cumulative hazard function of the conditional censoring distribution Γ . This gives a library consisting of $5^3 = 125$ learners for the conditional state occupation probability function F defined in equation (6). We use five folds for training and testing the models, and we repeat training and evaluation five times with different splits. The integrated Brier score (defined in Section 2.4.1) for all learners are shown in Figure 4. We see that the prediction performance is mostly affected by the choice of learner for the censoring distribution. Several combinations of learners give similar performance as measured by the integrated Brier score, provided that random survival forest is used to model the censoring distribution.

8. DISCUSSION

A major advantage of the joint survival super learner is that the performance of each combination of learners can be estimated without additional nuisance parameters. A potential drawback of our approach is that we are evaluating the loss of the learners on the level of the observed data distribution, while the target of the analysis is either the event-time distribution, or the censoring distribution, or both. Specifically, the finite-sample oracle inequality in Proposition 2 concerns the function F , which is a feature of $P \in \mathcal{P}$, while what we are typically interested in is Λ_j or S , which are features of $Q \in \mathcal{Q}$. We emphasize that, while the joint survival super learner provides estimates of Λ_j and Γ based on libraries \mathcal{A}_j and \mathcal{B} , the performance of these learners

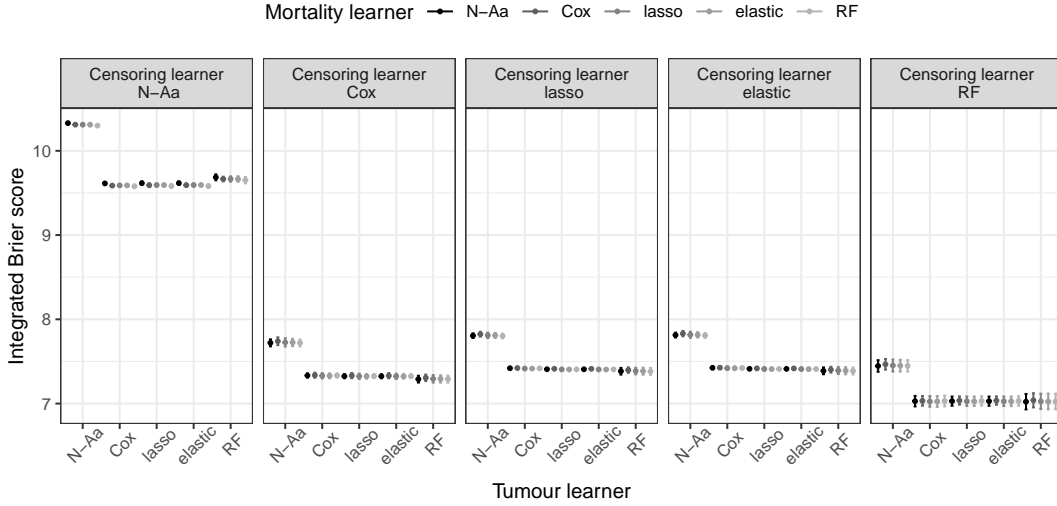


Fig. 4. The results of applying the 125 combinations of learners to the prostate cancer data set. The error bars are based on five repetitions using different splits. We refer to learners of Λ_1 , Λ_2 , and Γ as ‘Tumour learner’, ‘Mortality learner’, and ‘Censoring learner’, respectively.

is not assessed directly in terms of their respective target parameters, but only indirectly through the performance of F . For settings without competing risks, our numerical studies suggest that measuring the performance of F also leads to good performance for estimation of S .

Our proposed super learner can be implemented with a broad library of learners and using existing software. Furthermore, while the library $\mathcal{F}(\mathcal{A}_1, \mathcal{A}_2, \mathcal{B})$ consists of $|\mathcal{A}_1||\mathcal{A}_2||\mathcal{B}|$ many learners, we only need to fit $|\mathcal{A}_1| + |\mathcal{A}_2| + |\mathcal{B}|$ many learners in each fold. To evaluate the performance of each learner, we need to perform $|\mathcal{A}_1||\mathcal{A}_2||\mathcal{B}|$ many operations to calculate the integrated Brier score in each hold-out sample (one for each combination of the fitted models), but these operations are often negligible compared to fitting the models. Hence the joint survival super learner is essentially not more computationally demanding than any procedure that uses super learning to learn Λ_1 , Λ_2 , and Γ separately. While our proposal is based on constructing the library \mathcal{F} from libraries for learning Λ_1 , Λ_2 , and Γ , it could also be of interest to consider learners that estimate F directly.

In our numerical studies, we only considered learners of Λ_j and Γ that provide cumulative hazard functions which are piecewise constant in the time argument. This simplifies the calculation of F as the integrals in equation (7) reduce to sums. When Λ_j or Γ are absolutely continuous in the time argument, calculating F is more involved, but we expect that a good approximation can be achieved by discretisation.

Our original motivation for developing the joint survival super learner was for use within the framework of targeted or debiased machine learning – a general methodology that combines flexible, data-adaptive estimation with asymptotically valid inference for low-dimensional target parameters (van der Laan & Rose, 2011; Chernozhukov et al., 2018). In settings with right-censored competing risks, the relevant nuisance parameters often include the cause-specific and censoring cumulative hazard functions (e.g., van der Laan & Robins, 2003; Rytgaard & van der Laan, 2022). The joint survival super learner immediately provides estimates of these nuisance parameters and is hence particularly well suited for targeted and debiased machine learning. We

leave the study of the joint survival super learner in the context of targeted and debiased machine learning for a future [paperwork](#).

APPENDIX

Define $\bar{B}_{\tau,P}(F, o) = \bar{B}_{\tau}(F, o) - \bar{B}_{\tau}(F_P, o)$ and $R_P(F) = \mathbb{E}_P [\bar{B}_{\tau,P}(F, O)]$, where the integrated Brier score \bar{B}_{τ} was defined in Section [2.4.1](#). Recall the norm $\|\cdot\|_P$ defined in equation (13). 450

LEMMA 1. $R_P(F) = \|F - F_P\|_P^2$.

Proof. For any $t \in [0, \tau]$ and $l \in \{-1, 0, 1, 2\}$ we have

$$\begin{aligned} & \mathbb{E}_P [(F(t, l, X) - \mathbb{1}\{\eta(t) = l\})^2] \\ &= \mathbb{E}_P [(F(t, l, X) - F_P(t, l, X) + F_P(t, l, X) - \mathbb{1}\{\eta(t) = l\})^2] \\ &= \mathbb{E}_P [(F(t, l, X) - F_P(t, l, X))^2] + \mathbb{E}_P [(F_P(t, l, X) - \mathbb{1}\{\eta(t) = l\})^2] \\ &\quad + 2 \mathbb{E}_P [(F(t, l, X) - F_P(t, l, X))(F_P(t, l, X) - \mathbb{1}\{\eta(t) = l\})] \\ &= \mathbb{E}_P [(F(t, l, X) - F_P(t, l, X))^2] + \mathbb{E}_P [(F_P(t, l, X) - \mathbb{1}\{\eta(t) = l\})^2], \end{aligned} \quad 455$$

where the last equality follows from the tower property. Hence, using Fubini, we have

$$\mathbb{E}_P [\bar{B}_{\tau}(F, O)] = \|F - F_P\|_P^2 + \mathbb{E}_P [\bar{B}_{\tau}(F_P, O)].$$

Proof of Proposition 1. The result follows from Lemma 1. □ 460

Recall that we use \mathcal{F}_n to denote a library of learners for the function F , and that $\hat{\varphi}$ and $\tilde{\varphi}$ denotes, respectively, the discrete super learner and the oracle learner for the library \mathcal{F}_n , c.f., Section [2.4.1](#).

Proof of Proposition 2. Minimizing the loss \bar{B}_{τ} is equivalent to minimizing the loss $\bar{B}_{\tau,P}$, so the discrete super learner and oracle according to \bar{B}_{τ} and $\bar{B}_{\tau,P}$ are identical. By Lemma 1, $R_P(F) \geq 0$ for any $F \in \mathcal{H}_P$, and so using Theorem 2.3 from (van der Vaart et al., 2006) with $p = 1$, we have that for all $\delta > 0$, 465

$$\begin{aligned} & \frac{1}{K} \sum_{k=1}^K \mathbb{E}_P [R_P(\hat{\varphi}_n(\mathcal{D}_n^{-k}))] \\ & \leq (1 + 2\delta) \frac{1}{K} \sum_{k=1}^K \mathbb{E}_P [R_P(\tilde{\varphi}_n(\mathcal{D}_n^{-k}))] \\ & \quad + (1 + \delta) \frac{16K}{n} \log(1 + |\mathcal{F}_n|) \sup_{F \in \mathcal{H}_P} \left\{ M(F) + \frac{v(F)}{R_P(F)} \left(\frac{1}{\delta} + 1 \right) \right\} \end{aligned}$$

where for each $F \in \mathcal{H}_P$, $(M(F), v(F))$ is some Bernstein pair for the function $o \mapsto \bar{B}_{\tau,P}(F, o)$. As $\bar{B}_{\tau,P}(F, \cdot)$ is uniformly bounded by τ for any $F \in \mathcal{H}_P$, it follows from section 8.1 in (van der Vaart et al., 2006) that $(\tau, 1.5 \mathbb{E}_P [\bar{B}_{\tau,P}(F, O)^2])$ is a Bernstein pair for $\bar{B}_{\tau,P}(F, \cdot)$. Now, for any $a, b, c \in \mathbb{R}$ we have 470

$$\begin{aligned} (a - c)^2 - (b - c)^2 &= (a - b + b - c)^2 - (b - c)^2 \\ &= (a - b)^2 + (b - c)^2 + 2(b - c)(a - b) - (b - c)^2 \\ &= (a - b) \{(a - b) + 2(b - c)\} \\ &= (a - b) \{a + b - 2c\}, \end{aligned} \quad 475$$

so using this with $a = F(t, l, x)$, $b = F_P(t, l, x)$, and $c = \mathbb{1}\{\eta(t) = l\}$, we have by Jensen's inequality

$$\begin{aligned}
& \mathbb{E}_P [\bar{B}_{\tau, P}(F, O)^2] \\
& \leq 2\tau \mathbb{E}_P \left[\sum_{l=-1}^2 \int_0^\tau \left\{ (F(t, l, X) - \mathbb{1}\{\eta(t) = l\})^2 - (F_P(t, l, X) - \mathbb{1}\{\eta(t) = l\})^2 \right\}^2 dt \right] \\
& = 2\tau \mathbb{E}_P \left[\sum_{l=-1}^2 \int_0^\tau (F(t, l, X) - F_P(t, l, X))^2 \right. \\
& \quad \left. \times \{F(t, l, X) + F_P(t, l, X) - 2\mathbb{1}\{\eta(t) = l\}\}^2 dt \right] \\
& \leq 8\tau \mathbb{E}_P \left[\sum_{l=-1}^2 \int_0^\tau (F(t, l, X) - F_P(t, l, X))^2 dt \right] \\
& = 8\tau \|F - F_P\|_P^2.
\end{aligned}$$

Thus when $v(F) = 1.5 \mathbb{E}_P [\bar{B}_{\tau, P}(F, O)^2]$ we have by Lemma 1

$$\frac{v(F)}{R_P(F)} = 1.5 \frac{\mathbb{E}_P [\bar{B}_{\tau, P}(F, O)^2]}{\mathbb{E}_P [\bar{B}_{\tau, P}(F, O)]} \leq 12\tau,$$

and so using the Bernstein pairs $(\tau, 1.5 \mathbb{E}_P [\bar{B}_{\tau, P}(F, O)^2])$ we have

$$\sup_{F \in \mathcal{H}_P} \left\{ M(F) + \frac{v(F)}{R_P(F)} \left(\frac{1}{\delta} + 1 \right) \right\} \leq \tau \left(13 + \frac{12}{\delta} \right).$$

For all $\delta > 0$ we thus have

$$\begin{aligned}
\frac{1}{K} \sum_{k=1}^K \mathbb{E}_P [R_P(\hat{\varphi}_n(\mathcal{D}_n^{-k}))] & \leq (1 + 2\delta) \frac{1}{K} \sum_{k=1}^K \mathbb{E}_P [R_P(\tilde{\varphi}_n(\mathcal{D}_n^{-k}))] \\
& \quad + (1 + \delta) \log(1 + |\mathcal{F}_n|) \tau \frac{16K}{n} \left(13 + \frac{12}{\delta} \right),
\end{aligned}$$

and then the final result follows from Lemma 1. \square

Proof of Corollary 1. By definition of the oracle and Lemma 1,

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E}_P [\|\tilde{\varphi}_n(\mathcal{D}_n^{-k}) - F_P\|_P^2] \leq \frac{1}{K} \sum_{k=1}^K \mathbb{E}_P [\|\varphi_n(\mathcal{D}_n^{-k}) - F_P\|_P^2] = \mathbb{E}_P [\|\varphi_n(\mathcal{D}_n^{-k}) - F_P\|_P^2],$$

for all $n \in \mathbb{N}$, where the last equality follows because all the training sets \mathcal{D}_n^{-k} have the same distribution. The result then follows from Proposition 2 by letting δ grow to zero with n , for instance as $\delta_n = \log(n)^{-\varepsilon}$ for some $\varepsilon > 0$. \square

REFERENCES

- ANDERSEN, P. K., BORGAN, O., GILL, R. D. & KEIDING, N. (2012). *Statistical models based on counting processes*. Springer Science & Business Media.
- ANDERSEN, P. K., KLEIN, J. P. & ROSTHØJ, S. (2003). Generalised linear models for correlated pseudo-observations, with applications to multi-state models. *Biometrika*.
- BEGUN, J. M., HALL, W. J., HUANG, W.-M. & WELLNER, J. A. (1983). Information and asymptotic efficiency in parametric-nonparametric models. *The Annals of Statistics* **11**, 432–452.
- BENDER, R., AUGUSTIN, T. & BLETNER, M. (2005). Generating survival times to simulate Cox proportional hazards models. *Statistics in medicine* **24**, 1713–1723.
- BENICHOU, J. & GAIL, M. H. (1990). Estimates of absolute cause-specific risk in cohort studies. *Biometrics*, 813–826.

- BREIMAN, L. (1996). Stacked regressions. *Machine learning* **24**, 49–64.
- BRIER, G. W. et al. (1950). Verification of forecasts expressed in terms of probability. *Monthly weather review* **78**, 1–3. 505
- CHERNOZHUKOV, V., CHETVERIKOV, D., DEMIRER, M., DUFLO, E., HANSEN, C., NEWEY, W. & ROBINS, J. (2018). Double/debiased machine learning for treatment and structural parameters.
- COX, D. R. (1972). Regression models and life-tables. *Journal of the Royal Statistical Society: Series B (Methodological)* **34**, 187–202.
- EFRON, B. & TIBSHIRANI, R. (1997). Improvements on cross-validation: The .632+ bootstrap method. *Journal of the American Statistical Association* **92**, 548–560. 510
- FAN, J. & GIJBELS, I. (1996). *Local polynomial modelling and its applications*. Routledge.
- GENSHEIMER, M. F. & NARASIMHAN, B. (2019). A scalable discrete-time survival model for neural networks. *PeerJ* **7**, e6257.
- GERDS, T. A. (2019). *prodlim: Product-Limit Estimation for Censored Event History Analysis*. R package version 2019.11.13. 515
- GERDS, T. A. & KATTAN, M. W. (2021). *Medical risk prediction models: with ties to machine learning*. CRC Press.
- GERDS, T. A., KATTAN, M. W., SCHUMACHER, M. & YU, C. (2013). Estimating a time-dependent concordance index for survival prediction models with covariate dependent censoring. *Statistics in medicine* **32**, 2173–2184.
- GERDS, T. A., OHLENDORFF, J. S. & OZENNE, B. (2023). *riskRegression: Risk Regression Models and Prediction Scores for Survival Analysis with Competing Risks*. R package version 2023.03.22. 520
- GERDS, T. A. & SCHUMACHER, M. (2006). Consistent estimation of the expected Brier score in general survival models with right-censored event times. *Biometrical Journal* **48**, 1029–1040.
- GILL, R. D., VAN DER LAAN, M. J. & ROBINS, J. M. (1997). Coarsening at random: Characterizations, conjectures, counter-examples. In *Proceedings of the First Seattle Symposium in Biostatistics*. Springer. 525
- GNETING, T. & RAFTERY, A. E. (2007). Strictly proper scoring rules, prediction, and estimation. *Journal of the American statistical Association* **102**, 359–378.
- GOLMAKANI, M. K. & POLLEY, E. C. (2020). Super learner for survival data prediction. *The International Journal of Biostatistics* **16**, 20190065.
- GONZALEZ GINESTET, P., KOTALIK, A., VOCK, D. M., WOLFSON, J. & GABRIEL, E. E. (2021). Stacked inverse probability of censoring weighted bagging: A case study in the infcarehiv register. *Journal of the Royal Statistical Society Series C: Applied Statistics* **70**, 51–65. 530
- GRAF, E., SCHMOOR, C., SAUERBREI, W. & SCHUMACHER, M. (1999). Assessment and comparison of prognostic classification schemes for survival data. *Statistics in medicine* .
- HAN, X., GOLDSTEIN, M., PULL, A., WIES, T., PEROTTE, A. & RANGANATH, R. (2021). Inverse-weighted survival games. *Advances in Neural Information Processing Systems* **34**. 535
- HOTHORN, T., BÜHLMANN, P., DUDOIT, S., MOLINARO, A. & VAN DER LAAN, M. J. (2006). Survival ensembles. *Biostatistics* **7**, 355–373.
- ISHWARAN, H. & KOGALUR, U. (2025). *Fast Unified Random Forests for Survival, Regression, and Classification (RF-SRC)*. R package version 3.3.3. 540
- ISHWARAN, H., KOGALUR, U. B., BLACKSTONE, E. H. & LAUER, M. S. (2008). Random survival forests. *The annals of applied statistics* **2**, 841–860.
- KAPLAN, E. L. & MEIER, P. (1958). Nonparametric estimation from incomplete observations. *Journal of the American statistical association* **53**, 457–481.
- KATTAN, M. W. & GERDS, T. A. (2018). The index of prediction accuracy: an intuitive measure useful for evaluating risk prediction models. *Diagnostic and prognostic research* . 545
- KATTAN, M. W., ZELEFSKY, M. J., KUPELIAN, P. A., SCARDINO, P. T., FUKS, Z. & LEIBEL, S. A. (2000). Pretreatment nomogram for predicting the outcome of three-dimensional conformal radiotherapy in prostate cancer. *Journal of clinical oncology* **18**, 3352–3359.
- KATZMAN, J. L., SHAHAM, U., CLONINGER, A., BATES, J., JIANG, T. & KLUGER, Y. (2018). Deepsurv: personalized treatment recommender system using a Cox proportional hazards deep neural network. *BMC medical research methodology* **18**, 1–12. 550
- KELES, S., VAN DER LAAN, M. & DUDOIT, S. (2004). Asymptotically optimal model selection method with right censored outcomes. *Bernoulli* **10**, 1011–1037.
- KVAMME, H. & BORGAN, Ø. (2021). Continuous and discrete-time survival prediction with neural networks. *Lifetime Data Analysis* **27**, 710–736. 555
- LEE, C., ZAME, W., YOON, J. & VAN DER SCHAAR, M. (2018). Deephit: A deep learning approach to survival analysis with competing risks. In *Proceedings of the AAAI conference on artificial intelligence*, vol. 32.
- LEE, D. K., CHEN, N. & ISHWARAN, H. (2021). Boosted nonparametric hazards with time-dependent covariates. *Annals of Statistics* **49**, 2101. 560
- LI, Y., XU, K. S. & REDDY, C. K. (2016). Regularized parametric regression for high-dimensional survival analysis. In *Proceedings of the 2016 SIAM International Conference on Data Mining*. SIAM.
- LIU, P., SAWHNEY, S., HEIDE-JØRGENSEN, U., QUINN, R. R., JENSEN, S. K., MCLEAN, A., CHRISTIANSEN, C. F., GERDS, T. A. & RAVANI, P. (2024). Predicting the risks of kidney failure and death in adults with moderate to severe chronic

- 565 kidney disease: multinational, longitudinal, population based, cohort study. *British Medical Journal* **385**.
- MOGENSEN, U. B. & GERDS, T. A. (2013). A random forest approach for competing risks based on pseudo-values. *Statistics in medicine* **32**, 3102–3114.
- MOLINARO, A. M., DUDOIT, S. & VAN DER LAAN, M. J. (2004). Tree-based multivariate regression and density estimation with right-censored data. *Journal of Multivariate Analysis* **90**, 154–177.
- 570 MUNCH, A. (2024). *Targeted learning with right-censored data*. Phd thesis, University of Copenhagen.
- OZENNE, B., SØRENSEN, A. L., SCHEIKE, T., TORP-PEDERSEN, C. & GERDS, T. A. (2017). riskregression: Predicting the risk of an event using Cox regression models. *R Journal* **9**, 440–460.
- POLLEY, E. C. & VAN DER LAAN, M. J. (2011). Super learning for right-censored data. In *Targeted Learning: Causal Inference for Observational and Experimental Data*, M. J. van der Laan & S. Rose, eds. Springer, pp. 249–258.
- 575 R CORE TEAM (2024). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- RYTGAARD, H. C. & VAN DER LAAN, M. J. (2022). Targeted maximum likelihood estimation for causal inference in survival and competing risks analysis. *Lifetime Data Analysis*, 1–30.
- SACHS, M. C., DISCACCIATI, A., EVERHOV, Å. H., OLÉN, O. & GABRIEL, E. E. (2019). Ensemble prediction of time-to-event outcomes with competing risks: A case-study of surgical complications in Crohn’s disease. *Journal of the Royal Statistical Society Series C: Applied Statistics* **68**, 1431–1446.
- 580 STEINGRIMSSON, J. A., DIAO, L. & STRAWDERMAN, R. L. (2019). Censoring unbiased regression trees and ensembles. *Journal of the American Statistical Association*.
- THERNEAU, T. M. (2022). *A Package for Survival Analysis in R*. R package version 3.4-0.
- 585 VAN DER LAAN, M. J. & DUDOIT, S. (2003). Unified cross-validation methodology for selection among estimators and a general cross-validated adaptive epsilon-net estimator: Finite sample oracle inequalities and examples. Tech. rep., Division of Biostatistics, University of California.
- VAN DER LAAN, M. J., POLLEY, E. C. & HUBBARD, A. E. (2007). Super learner. *Statistical applications in genetics and molecular biology* **6**.
- 590 VAN DER LAAN, M. J. & ROBINS, J. M. (2003). *Unified methods for censored longitudinal data and causality*. Springer Science & Business Media.
- VAN DER LAAN, M. J. & ROSE, S. (2011). *Targeted learning: causal inference for observational and experimental data*. Springer Science & Business Media.
- VAN DER VAART, A. W., DUDOIT, S. & VAN DER LAAN, M. J. (2006). Oracle inequalities for multi-fold cross validation. *Statistics & Decisions* **24**, 351–371.
- 595 VERWEIJ, P. J. & VAN HOUWELINGEN, H. C. (1993). Cross-validation in survival analysis. *Statistics in medicine* **12**, 2305–2314.
- WESTLING, T., LUEDTKE, A., GILBERT, P. & CARONE, M. (2021). Inference for treatment-specific survival curves using machine learning. *arXiv preprint arXiv:2106.06602*.
- 600 WOLPERT, D. H. (1992). Stacked generalization. *Neural networks* **5**, 241–259.
- YAO, J., ZHU, X., ZHU, F. & HUANG, J. (2017). Deep correlational learning for survival prediction from multi-modality data. In *International conference on medical image computing and computer-assisted intervention*. Springer.

[Received on 2 January 2024. Editorial decision on 3 February 2025]