

## RESEARCH ARTICLE

# Case-Base Neural Networks: survival analysis with time-varying, higher-order interactions

Jesse Islam<sup>1</sup> | Maxime Turgeon<sup>2</sup> | Robert Sladek<sup>3</sup> | Sahir Bhatnagar<sup>4</sup><sup>1</sup>Department of Quantitative life sciences,  
McGill University<sup>2</sup>Department of Statistics, University of  
Manitoba<sup>3</sup>Department of Human Genetics, McGill  
University<sup>4</sup>Department of Biostatistics, McGill  
University**Correspondence**

Email: jesse.islam@mail.mcgill.ca

Neural network-based survival methods have been developed to model data-driven covariate interactions. While these methods have led to better predictive performance than regression-based approaches, they cannot model both time-varying interactions and complex baseline hazards. To address this, we propose Case-Base Neural Networks (CBNN) as a new approach that combines the case-base sampling framework with flexible architectures. Our method naturally accounts for censoring, does not require additional hyperparameters and uses binary cross-entropy as the loss function. Using a novel sampling scheme and data augmentation, we incorporate time directly into a feed-forward neural network. CBNN predicts the probability of an event occurring at a given moment and estimates the hazard function. We compare the performance of CBNN to survival methods based on regression and neural networks in two simulations and two real data applications. We report two time-dependent metrics for each model. In the simulations and real data applications, CBNN provides a more consistent predictive performance across time and outperforms the competing neural network approaches. For a simple simulation with an exponential hazard model, CBNN outperforms the other neural network methods. For a complex simulation, which highlights the ability of CBNN to model both a complex baseline hazard and time-varying interactions, CBNN outperforms all competitors. The first real data application shows CBNN outperforming all neural network competitors, while a second real data application shows competitive performance. This study highlights the benefit of combining case-base sampling with deep learning to provide a simple and flexible modeling framework for data-driven, time-varying interaction modeling of survival outcomes. An R package is available at <https://github.com/Jesse-Islam/cbnn>.

**KEYWORDS:**

survival analysis, machine learning, case-base, neural network

## 1 | INTRODUCTION

Smooth-in-time accelerated failure time (AFT) models can estimate absolute risks by modeling the hazard directly through a user-specified baseline hazard distribution<sup>1</sup>. Cox proportional hazards models are used more often than AFT models, causing analyses to be based on hazard ratios and relative risks rather than on survival curves and absolute risks<sup>2</sup>. The identification of an

appropriate distribution for the baseline hazard in an AFT model may be difficult for common diseases that have many interacting risk factors, or a Cox model where the disease pathogenesis may change with age, making it difficult to maintain the proportional hazards assumption. For example, previous studies of breast cancer incidence have discovered time-varying interactions with covariates of interest, such as tumor size<sup>3</sup>. One approach to provide flexibility in the baseline hazard involves using the basis of splines on time in our model<sup>4</sup>. However, regression-based models are limited in that they require prior knowledge about potential time-varying interactions and their quantitative effects.

Neural networks provide a data-driven approach to approximating interaction terms. For example, DeepSurv is a neural network-based proportional hazards model that implements the Cox partial log-likelihood as a custom loss function<sup>5</sup>, resulting in a stepwise absolute risk curve that cannot accommodate time-varying interactions. Compared to Cox regression, DeepSurv shows better performance on the Study to Understand Prognoses Preferences and Risk Treatments (SUPPORT) dataset<sup>6</sup>. To handle non-proportional hazards, a modification to the loss function was proposed<sup>7</sup>. To remove the need for tuning the number of layers and nodes, the concept of extreme learning machines has been applied to a Cox neural network model<sup>8</sup>.

As an alternative method that assumes a baseline hazard distribution, DeepHit specifies each survival time of interest in the model and directly estimates survival curves, rather than deriving a hazard function<sup>9</sup>. It assumes an inverse Gaussian distribution as the baseline hazard and it outperforms DeepSurv on the Molecular Taxonomy of Breast Cancer International Consortium (METABRIC) dataset<sup>10</sup>.

Providing further flexibility, Deep Survival Machines (DSM) is a parametric survival model using neural networks with a mixture of distributions as the baseline hazard<sup>11</sup>. On both the SUPPORT and METABRIC datasets, DSM outperforms DeepSurv and DeepHit<sup>11</sup>. However, like DeepHit and DeepSurv, DSM cannot model time-varying interactions.

We note that these alternative neural network approaches all require custom loss functions<sup>5,9,11</sup>. DeepHit introduces a hyperparameter weighing its two loss functions (negative log-likelihood and ranking losses) while DSM requires a two-phase learning process and user implementations for distributions beyond Log-Normal or Weibull<sup>9,11</sup>. Regression-based approaches require prior specification of all interaction terms, which makes it challenging to model covariate effects that change over time. The current neural network models provide flexibility at the cost of opacity, while regression models provide clarity at the cost of flexibility.

In this article, we propose Case-Base Neural Networks (CBNN) as a method that models time-varying interactions and a flexible baseline hazard using commonly available neural network components. Our approach to modeling the full hazard uses case-base sampling<sup>2</sup>. This sampling technique allows probabilistic models to predict survival outcomes. As part of the case-base framework, we use transformations of time as a feature (covariate) to specify different baseline hazards. For example, by including splines of time as covariates, we can approximate the Royston-Parmar flexible baseline hazard model<sup>4</sup>, however, this still requires explicit use of time-varying interactions. CBNN can model both without extra tuning parameters.

In Section 2, we describe how case-base sampling and neural networks are combined both conceptually and algebraically, along with our hyperparameter choices and software implementation. Section 3, describes our metrics and compares the performance of CBNN, DeepSurv, DeepHit, DSM, Cox regression and case-base using logistic regression (CBLR) on simulated data. Section 4 describes the real-data analysis, while Section 5 explores the implications of our results and contextualizes them within neural network survival analysis in a single event setting.

## 2 | CASE-BASE NEURAL NETWORK METHODOLOGY, METRICS AND SOFTWARE

In this section, we define case-base sampling, which converts the total survival time into discrete person-specific moments (person-moments). Then, we detail how neural networks can be used within this framework, explicitly incorporating time as a feature while adjusting for the sampling bias. Finally, we report on the software versions used. An R package is available for use at <https://github.com/Jesse-Islam/cbnn>. The entire code base to generate the analyses in this paper is available at <https://github.com/Jesse-Islam/cbnnManuscript>.

### 2.1 | Case-base sampling

Case-base sampling is an alternative framework for survival analysis<sup>2</sup>. In case-base sampling, we sample from the continuous survival time of each person in our dataset to create a *base series* of *person-moments*. This *base series* complements the *case series*, which contains all person-moments at which the event of interest occurs.

For each person-moment sampled, let  $X_i$  be the corresponding covariate profile  $(x_{i1}, x_{i2}, \dots, x_{ip})$ ,  $T_i$  be the time of the person-moment and  $Y_i$  be the indicator variable for whether the event of interest occurred at time  $T_i$ . We estimate the hazard function  $h(t | X_i)$  using the sampled person-moments. Recall that  $h(t | X_i)$  is the instantaneous potential of experiencing the event at time  $t$  for a given set of covariates  $X_i$ , assuming  $T_i \geq t$ .

Now, let  $b$  be the (user-defined) size of the *base series* and let  $B$  be the sum of all follow-up times for the individuals in the study. If we sample the *base series* uniformly across the study base, then the hazard function of the sampling process is equal to  $b/B$ . Therefore, we have the following equality<sup>a</sup>:

$$\frac{P(Y_i = 1 | X_i, T_i)}{P(Y_i = 0 | X_i, T_i)} = \frac{h(T_i | X_i)}{b/B}. \quad (1)$$

The odds of a person-moment being a part of the *case series* is the ratio of the hazard  $h(T_i | X_i)$  and the uniform rate  $b/B$ . Using (1), we can see how the log-hazard function can be estimated from the log-odds arising from case-base sampling:

$$\log(h(t | X_i)) = \log\left(\frac{P(Y_i = 1 | X_i, t)}{P(Y_i = 0 | X_i, t)}\right) + \log\left(\frac{b}{B}\right). \quad (2)$$

To estimate the correct hazard function, we adjusting for the bias introduced when sampling a fraction of the study base  $B$  by adding the offset term  $\log\left(\frac{b}{B}\right)$  as in (2). Next, we propose using neural networks to model the odds.

## 2.2 | Neural networks to model the hazard function

After case-base sampling, we pass all features, including time, into any user-defined feed-forward component, to which an offset term is added, then passed through a sigmoid activation function (Figure 1). We use the sigmoid function as its inverse is the odds, which we can use to calculate the hazard. The general form for the neural network using CBNN is:

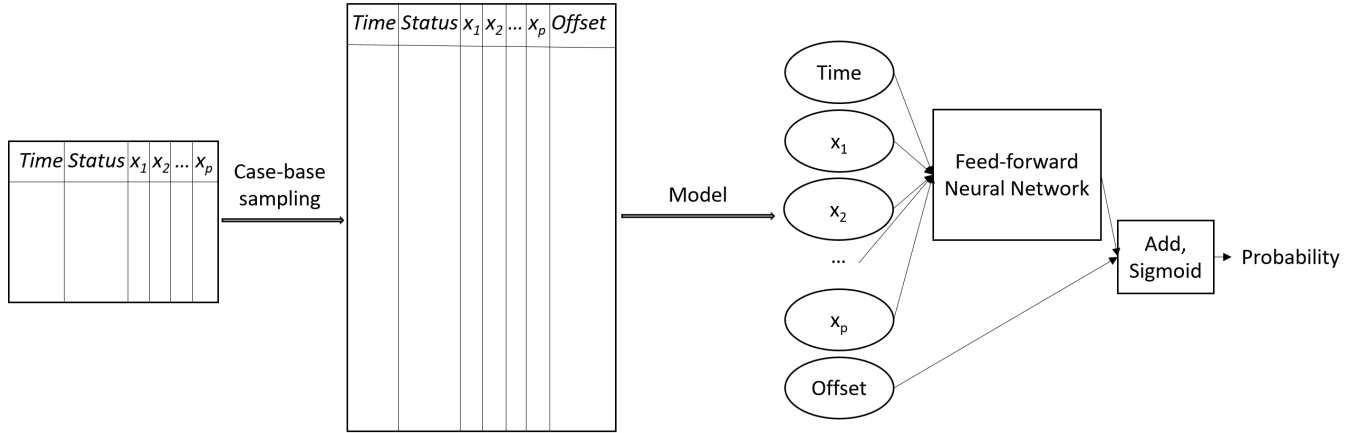
$$P(Y = 1 | X, T) = \text{sigmoid}\left(f_\theta(X, T) + \log\left(\frac{B}{b}\right)\right), \quad (3)$$

where  $T$  is a random variable representing the event time,  $X$  is the random variable for a covariate profile,  $f_\theta(X, T)$  represents any feed-forward neural network architecture,  $\log\left(\frac{B}{b}\right)$  is bias term set by case-base sampling,  $\theta$  is the set of parameters learned by the neural network and  $\text{sigmoid}(x) = \frac{1}{1+e^{-x}}$ . By approximating a higher-order polynomial of time using a neural network, the baseline hazard specification is now data-driven, where user-defined hyperparameters such as regularization, number of layers and nodes control the flexibility of the fukk hazard function. We provide a detailed description of the choices we made in the next sub-section.

The following derivation shows how our probability estimate is converted to odds:

$$\begin{aligned} \log(h(t | X)) &= \log\left(\frac{\text{sigmoid}\left(f_\theta(X, T) + \log\left(\frac{B}{b}\right)\right)}{1 - \text{sigmoid}\left(f_\theta(X, T) + \log\left(\frac{B}{b}\right)\right)}\right) + \log\left(\frac{b}{B}\right) \\ &= \log\left(\frac{\frac{\exp\left(f_\theta(X, T) + \log\left(\frac{B}{b}\right)\right)}{\exp\left(f_\theta(X, T) + \log\left(\frac{B}{b}\right)\right) + 1}}{1 - \frac{\exp\left(f_\theta(X, T) + \log\left(\frac{B}{b}\right)\right)}{\exp\left(f_\theta(X, T) + \log\left(\frac{B}{b}\right)\right) + 1}}\right) + \log\left(\frac{b}{B}\right) \\ &= \log\left(\exp\left(f_\theta(X, T) + \log\left(\frac{B}{b}\right)\right)\right) + \log\left(\frac{b}{B}\right) \\ &= f_\theta(X, T) + \log\left(\frac{B}{b}\right) + \log\left(\frac{b}{B}\right) \\ &= f_\theta(X, T). \end{aligned}$$

<sup>a</sup>We are abusing notation here, conflating hazards with probabilities. For a rigorous treatment, see Saarela & Hanley (2015) section 3<sup>12</sup>.



**FIGURE 1** Steps involved in CBNN from case-base sampling to the model framework we use for training. The first step is case-base sampling, completed before training begins. Next, we pass this sampled data through a feed-forward neural network. We add the offset and pass that through a sigmoid activation function, whose output is a probability. Once the neural network model completes its training, we can convert the probability output to a hazard, using it for our survival outcomes of interest.

We use binary cross-entropy as our loss function<sup>13</sup>:

$$L(\theta) = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(\hat{f}_\theta(x_i, t_i)) + (1 - y_i) \cdot \log(1 - \hat{f}_\theta(x_i, t_i)),$$

where  $\hat{f}_\theta(x_i, t_i)$  is our estimate for a given covariate profile and time,  $y_i$  is our target value specifying whether an event occurred and  $N$  represents the number of individuals in our training set.

Backpropagation with an appropriate minimization algorithm (e.g. Adam, RMSPropagation, stochastic gradient descent) is used to optimize the parameters in the model<sup>13</sup>. For our analysis, we use Adam as implemented in Keras<sup>13</sup>. Note that the size of the *case series* is fixed as the number of events, but we can make the *base series* as large as we want. A ratio of 100:1 *base series* to *case series* is sufficient<sup>2</sup>. We pass our feed-forward neural network through a sigmoid activation function (Figure 1). Finally, we can convert this model output to a hazard. When using our model for predictions, we manually set the offset term to 0 in the new data, as we account for the bias during the fitting process.

Since we are directly modeling the hazard, we can readily estimate the risk function ( $F$ ) at time  $t$  for a covariate profile  $X$ , viz.

$$F(t | X) = 1 - \exp\left(-\int_0^t h(u|X) du\right). \quad (4)$$

We use a finite Riemann sum<sup>14</sup> to approximate the integral in (4).

### 2.3 | Hyperparameter selection

Neural networks are flexible when defining the architecture and optimization parameters. These hyperparameter decisions can affect the estimated parameters and were chosen during a set of initial simulations to determine if CBNN can learn complex interactions in practice. We set a max number of epochs to be 2000, batch size as 512, learning rate as  $10^{-3}$ , decay as  $10^{-7}$ , patience as 10 epochs, {50, 50, 25, 25} nodes in each hidden layer with 50% dropout at each layer, a minimum delta loss on the validation set of  $10^{-7}$  over 10 epochs and Adam<sup>13</sup> as our optimizer. These choices may permit the model to approximate higher-order interactions while preventing over-fitting<sup>15</sup>. We fix a train-validation-test split that allows us to update the weights with a subset of the data (training), assess performance at each epoch (validation) and gauge performance of the final model (test) for each method. We select the best weights after training based on validation loss<sup>13</sup>.

## 2.4 | Software implementation

R<sup>16</sup> and python<sup>17</sup> are used to evaluate methods from both languages. We fit the Cox model using the **survival** package<sup>18</sup>, the CBLR model using the **casebase** package<sup>19</sup>, the DeepSurv model using the **survivalmodels** package<sup>20</sup>, the DeepHit model using **pycox**<sup>9</sup> and the DSM model using **DeepSurvivalMachines**<sup>11</sup>. We made the components of CBNB using the **casebase** package for the sampling step and the **keras**<sup>21</sup> package for our neural network architecture. The **simSurv** package<sup>22</sup> is used for our simulation studies, while **flexSurv**<sup>23</sup> is used to fit a flexible baseline hazard using splines for our complex simulation. We use the implementation of  $C_{IPCW}$  from the python package **sksurv**<sup>24</sup>. The **riskRegression** package<sup>25</sup> is used to get the Index of Prediction Accuracy (IPA metric). Both metrics are described in detail in the following section. We modify the **riskRegression** package to be used with any user supplied risk function  $F$ . To ensure that both R and Python-based models are running in unison on the same data through our simulations and bootstrap, we use the **reticulate** package<sup>26</sup>.

## 3 | SIMULATION STUDIES

In this section, we use simulated data to evaluate the performance of CBNB and compare our approach with existing regression-based (Cox, CBLR) and neural network-based (DeepHit, DeepSurv, DSM) methods. We specify a linear combination of each covariate as the linear predictor in regression-based approaches (Cox, CBLR), which contrasts with neural network approaches that allow for non-linear interactions. We simulate data under a simple exponential model and a complex baseline hazard with time-varying interactions, each with 10% random censoring. For both settings, we simulate three covariates:

$$z_1 \sim \text{Bernoulli}(0.5) \quad z_2 \sim \begin{cases} N(0, 0.5) & \text{if } z_1 = 0 \\ N(10, 0.5) & \text{if } z_1 = 1 \end{cases} \quad z_3 \sim \begin{cases} N(8, 0.5) & \text{if } z_1 = 0 \\ N(-3, 0.5) & \text{if } z_1 = 1. \end{cases}$$

The DeepHit-specific hyperparameter alpha is set to 0.5 (equal weight between its negative log-likelihood and ranking losses<sup>9</sup>). We modify the **DeepSurvivalMachines**<sup>11</sup> package to include dropout and a minimum delta loss during the fitting process. For DSM, we define a mixture model of six Weibull distributions for the baseline hazard. All other hyperparameters are held constant across all neural network methods in both the simulation studies and real data applications.

Besides the methods mentioned above, we include the Optimal model in our comparisons using CBLR. That is, we include the exact functional form of the covariates in a CBLR model (referred to as Optimal for simplicity). We calculate  $t$ -based 95% confidence intervals using 100 replications of the simulated data. For all analyses, we use 80% for training and 20% for the test set. 20% of the training set is kept for validation at each epoch. We predict risk functions  $F$  using (4) for individuals in the test set, which are used to calculate our  $C_{IPCW}$  and IPA scores.

### 3.1 | Performance metrics

We use two metrics to assess the performance of the different methods of interest: 1) (IPA)<sup>27</sup> and 2) inverse probability censoring weights-adjusted concordance index ( $C_{IPCW}$ )<sup>28</sup>, which we define below. Both these time-dependent metrics provide transparency as to when in follow-up time each model may perform better than the others.

#### 3.1.1 | Index of prediction accuracy (IPA)

The IPA is a function of the Brier score ( $BS(t)$ )<sup>29</sup>, which is defined as

$$BS(t) = \frac{1}{N} \sum_{i=1}^N \left( \frac{(1 - \hat{F}(t | X_i))^2 \cdot I(T_i \leq t, \delta_i = 1)}{\hat{G}(T_i)} + \frac{\hat{F}(t | X_i)^2 \cdot I(T_i > t)}{\hat{G}(t)} \right), \quad (5)$$

where  $\delta_i = 1$  shows individuals who have experienced the event,  $N$  represents the number of samples in our dataset over which we calculate  $BS(t)$ ,  $\hat{G}(t) = P[c > t]$  is a non-parametric estimate of the censoring distribution,  $c$  is censoring time and  $T_i$  is an individual's survival or censoring time. The Brier score provides a score that accounts for the information loss because of censoring. There are three categories of individuals that may appear within the dataset once we fix our  $t$  of interest. Individuals who experienced the event before  $t$  are present in the first term of the equation. The second term of the equation includes individuals who experience the event or are censored after  $t$ . Those censored before  $t$  are the third category of people.

The inverse probability censoring weights (IPCW) adjustment ( $G(\cdot)$ ) is to account for these category three individuals whose information is missing. The IPA score as a function of time is given by

$$\text{IPA}(t) = 1 - \frac{BS_{\text{model}}(t)}{BS_{\text{null}}(t)},$$

where  $BS_{\text{model}}(t)$  represents the Brier score over time  $t$  for the model of interest and  $BS_{\text{null}}(t)$  represents the Brier score if we use an unadjusted Kaplan-Meier (KM) curve as the prediction for all observations<sup>27</sup>. Note that IPA has an upper bound of one, where positive values show an increase in performance over the null model and negative values show that the null model performs better. These scores show how performance changes over follow-up time.

A potential artifact of IPA is that the score is unstable at earlier and later survival times. This is because of a near equivalent Brier score among each model and the null model. At small values, a difference of 0.1 creates a more significant fold change than at larger values. As the Brier score is potentially small at the start and at the end of their respective curves, The IPA score may be unstable at the same locations.

### 3.1.2 | Inverse probability censoring weights-adjusted concordance index

The  $C_{IPCW}$  is a non-proper, rank-based metric that does not depend on the censoring times in the test data<sup>28</sup>. The  $C_{IPCW}$  is given by

$$C_{IPCW}(t) = \frac{\sum_{i=1}^N \sum_{j=1}^N \delta_i \left\{ \hat{G}(T_i) \right\}^{-2} I(T_i < T_j, T_i < t) I\left(\hat{F}(t|X_i) > \hat{F}(t|X_j)\right)}{\sum_{i=1}^N \sum_{j=1}^N \delta_i \left\{ \hat{G}(T_i) \right\}^{-2} I(T_i < T_j, T_i < t)}. \quad (6)$$

where the follow-up period of interest is  $(0, t)$ ,  $I(\cdot)$  is an indicator function,  $\hat{F}(X_i, t)$  is the risk function estimated for everyone in the study at time  $t$  and  $C_{IPCW}$  can compare the performance of different models, where a higher score is better. Note that the  $C_{IPCW}$  may produce misleading performance, as it ranks based on survival times, not event status<sup>30</sup>. This metric is considered an unbiased population concordance measure because of the IPCW adjustment<sup>28</sup>.

## 3.2 | Simple simulation: constant baseline hazard

We simulate data from a simple model that primarily depends on a constant baseline hazard:

$$\log h(t | X_i) = \beta_1 z_1 + \beta_2 z_2 + \beta_3 z_3,$$

The covariate effects are given by  $\beta_1 = 0.1, \beta_2 = 0.1, \beta_3 = 0.1$ . Once we simulate survival times, we introduce 10% random censoring.

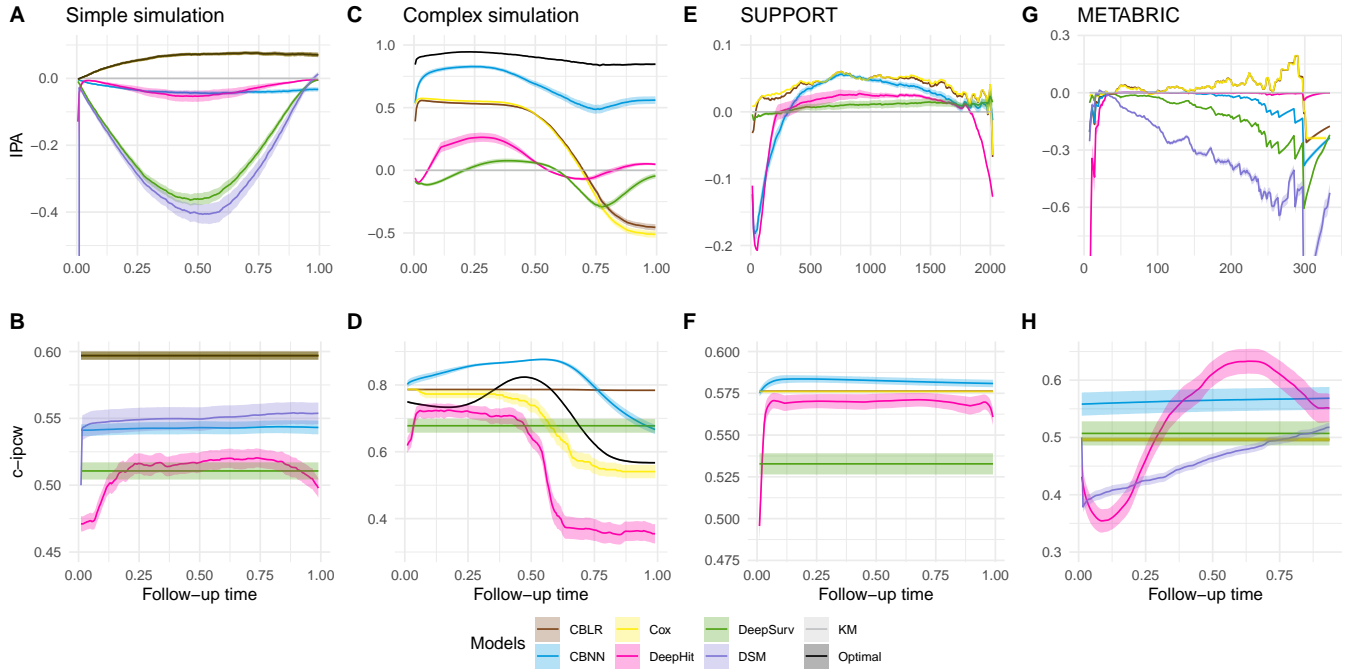
### 3.2.1 | Performance comparison in simple simulation

Figure 2 A, B and Table 1 A show the results for the simple simulation. The regression-based methods (CBLR, Cox, Optimal) outperform the neural network ones in the simple simulation setting. Among the neural network approaches, CBNN outperforms all other methods in terms of both IPA and  $C_{IPCW}$  (Figure 2 A, B). Specifically, we see CBNN is consistent across time with smaller confidence intervals compared to DeepHit, DeepSurv and DSM.

In this simple setting, the regression models are much closer to the Optimal model, while the neural network models perform worse than the KM null model. The wide confidence bands in  $C_{IPCW}$  suggest the neural network models may be over-parameterized.

## 3.3 | Complex simulation: flexible baseline hazard, time-varying interactions

This simulation demonstrates performance with the presence of a complex baseline hazard and a time-varying interaction. Originally used to show the spline-based hazard model proposed by Royston and Parmar<sup>4</sup>, the breast cancer dataset provides



**FIGURE 2** Summarizes the simple simulation (A, B), complex simulation (C, D), SUPPORT case study (E, F) and METABRIC case study (G, H) results. The first row shows the IPA score for each model in each study over follow-up time. Negative values mean our model performs worse than the null model and positive values mean the model performs better. The second row demonstrates the  $C_{IPCW}$  score for each model in each study over follow-up time. A score of 1 is the maximum performance for either metric. Each model-specific metric in each study shows a 95-percent confidence interval over 100 iterations. The models of interest are case-base with logistic regression (CBLR), Case-Base Neural Networks (CBNN), Cox, DeepHit, DeepSurv, Deep Survival Machines (DSM), Optimal (a CBLR model with the exact interaction terms and baseline hazard specified) and Kaplan-Meier (to serve as a baseline, predicting the average for all individuals).

a complex hazard from which we simulate, available in the **flexsurv** R package<sup>23</sup>. To increase the complexity of our data-generating mechanism for this simulation, we design the model as follows:

$$\log h(t | X_i) = \sum_{i=1}^5 (\gamma_i \cdot \psi_i) + \beta_1(z_1) + \beta_2(z_2) + \beta_3(z_3) + \tau_1(z_1 \cdot z_2 \cdot t) + \tau_2(z_1 \cdot z_3) + \tau_3(z_2 \cdot z_3),$$

**TABLE 1** Four tables representing performance at certain follow-up times for the simple simulation, complex simulation, SUPPORT and METABRIC. Each table shows performance for each method in each study at 25%, 50%, 75% and 100% of follow-up time. The bold elements show the best model for each study, at each follow-up time of interest. These tables are included to provide exact measures at certain intervals. The models of interest are: Cox, case-base with logistic regression (CBLR), DeepSurv, DeepHit, Case-Base Neural Network (CBNN), Optimal and Deep Survival Machines (DSM).

<b>A: Simple</b>								
Method	IPA				c-ipcw			
	25	50	75	100	25	50	75	100
Cox	<b>0.05 (0.05,0.06)</b>	<b>0.07 (0.07,0.08)</b>	<b>0.08 (0.07,0.08)</b>	<b>0.07 (0.06,0.08)</b>	<b>0.6 (0.59,0.6)</b>	<b>0.6 (0.59,0.6)</b>	<b>0.6 (0.59,0.6)</b>	<b>0.6 (0.59,0.6)</b>
CBLR	<b>0.05 (0.05,0.06)</b>	<b>0.07 (0.07,0.08)</b>	<b>0.08 (0.07,0.08)</b>	<b>0.07 (0.06,0.08)</b>	<b>0.6 (0.59,0.6)</b>	<b>0.6 (0.59,0.6)</b>	<b>0.6 (0.59,0.6)</b>	<b>0.6 (0.59,0.6)</b>
DeepSurv	-0.26 (-0.27,-0.25)	-0.36 (-0.38,-0.34)	-0.22 (-0.24,-0.21)	0 (-0.01,0)	0.51 (0.5,0.52)	0.51 (0.5,0.52)	0.51 (0.5,0.52)	0.51 (0.5,0.52)
DeepHit	-0.03 (-0.04,-0.02)	-0.05 (-0.07,-0.04)	-0.03 (-0.04,-0.02)	-0.01 (-0.01,0)	0.52 (0.51,0.52)	0.52 (0.51,0.52)	0.52 (0.51,0.53)	0.5 (0.49,0.5)
CBNN	-0.03 (-0.04,-0.03)	-0.04 (-0.05,-0.04)	-0.04 (-0.05,-0.03)	-0.03 (-0.04,-0.03)	0.54 (0.54,0.55)	0.54 (0.54,0.55)	0.54 (0.54,0.55)	0.54 (0.54,0.55)
Optimal	<b>0.05 (0.05,0.06)</b>	<b>0.07 (0.07,0.08)</b>	<b>0.08 (0.07,0.08)</b>	<b>0.07 (0.06,0.08)</b>	<b>0.6 (0.59,0.6)</b>	<b>0.6 (0.59,0.6)</b>	<b>0.6 (0.59,0.6)</b>	<b>0.6 (0.59,0.6)</b>
DSM	-0.28 (-0.3,-0.27)	-0.4 (-0.43,-0.37)	-0.28 (-0.3,-0.25)	0.01 (0.01,0.02)	0.55 (0.54,0.56)	0.55 (0.54,0.56)	0.55 (0.54,0.56)	0.55 (0.55,0.56)

<b>B: Complex</b>								
Method	IPA				c-ipcw			
	25	50	75	100	25	50	75	100
Cox	0.55 (0.54,0.56)	0.45 (0.44,0.46)	-0.19 (-0.21,-0.17)	-0.51 (-0.54,-0.48)	0.77 (0.76,0.79)	0.74 (0.72,0.76)	0.56 (0.54,0.58)	0.54 (0.52,0.56)
CBLR	0.53 (0.52,0.54)	0.45 (0.44,0.46)	-0.15 (-0.16,-0.13)	-0.46 (-0.48,-0.43)	0.79 (0.78,0.79)	0.79 (0.78,0.79)	0.79 (0.78,0.79)	<b>0.78 (0.78,0.79)</b>
DeepSurv	0.03 (0.01,0.05)	0.06 (0.05,0.07)	-0.27 (-0.29,-0.25)	-0.04 (-0.07,-0.02)	0.68 (0.66,0.7)	0.68 (0.66,0.7)	0.68 (0.66,0.7)	0.68 (0.66,0.7)
DeepHit	0.26 (0.22,0.3)	0.07 (0.06,0.08)	-0.06 (-0.08,-0.04)	0.05 (0.04,0.06)	0.72 (0.7,0.74)	0.65 (0.62,0.67)	0.37 (0.33,0.4)	0.35 (0.33,0.38)
CBNN	0.83 (0.81,0.84)	0.69 (0.67,0.71)	0.49 (0.45,0.52)	0.56 (0.53,0.59)	<b>0.86 (0.85,0.86)</b>	<b>0.87 (0.87,0.88)</b>	<b>0.8 (0.79,0.81)</b>	0.67 (0.65,0.68)
Optimal	<b>0.94 (0.94,0.95)</b>	<b>0.9 (0.9,0.91)</b>	<b>0.84 (0.84,0.85)</b>	<b>0.85 (0.84,0.85)</b>	0.74 (0.74,0.75)	0.82 (0.82,0.82)	0.63 (0.63,0.63)	0.57 (0.56,0.57)

<b>C: SUPPORT</b>								
Method	IPA				c-ipcw			
	25	50	75	100	25	50	75	100
Cox	<b>0.04 (0.04,0.04)</b>	<b>0.05 (0.05,0.05)</b>	<b>0.04 (0.04,0.04)</b>	-0.06 (-0.07,-0.06)	<b>0.58 (0.58,0.58)</b>	<b>0.58 (0.58,0.58)</b>	<b>0.58 (0.58,0.58)</b>	<b>0.58 (0.58,0.58)</b>
CBLR	<b>0.04 (0.04,0.04)</b>	<b>0.05 (0.05,0.05)</b>	<b>0.04 (0.04,0.04)</b>	-0.07 (-0.07,-0.06)	<b>0.58 (0.58,0.58)</b>	<b>0.58 (0.58,0.58)</b>	<b>0.58 (0.58,0.58)</b>	<b>0.58 (0.58,0.58)</b>
DeepSurv	0.01 (0,0.01)	0.01 (0.01,0.02)	0.01 (0.01,0.02)	<b>0.01 (0.01,0.02)</b>	0.53 (0.53,0.54)	0.53 (0.53,0.54)	0.53 (0.53,0.54)	0.53 (0.53,0.54)
DeepHit	0.02 (0.01,0.03)	0.02 (0.02,0.03)	0.02 (0.02,0.02)	-0.13 (-0.13,-0.12)	0.57 (0.57,0.57)	0.57 (0.57,0.57)	0.57 (0.57,0.57)	0.56 (0.56,0.57)
CBNN	0.03 (0.03,0.04)	<b>0.05 (0.04,0.05)</b>	0.03 (0.02,0.03)	-0.01 (-0.02,0)	<b>0.58 (0.58,0.59)</b>	<b>0.58 (0.58,0.58)</b>	<b>0.58 (0.58,0.58)</b>	<b>0.58 (0.58,0.58)</b>

<b>D: METABRIC</b>								
Method	IPA				c-ipcw			
	25	50	75	100	25	50	75	100
Cox	0.00 (-0.01,0.00)	<b>0.03 (0.03,0.04)</b>	<b>0.06 (0.05,0.06)</b>	-0.23 (-0.25,-0.23)	0.5 (0.49,0.5)	0.5 (0.49,0.5)	0.5 (0.49,0.5)	0.5 (0.49,0.5)
CBLR	<b>0.00 (0.00,0.00)</b>	<b>0.03 (0.03,0.03)</b>	<b>0.06 (0.05,0.06)</b>	-0.18 (-0.18,-0.17)	0.5 (0.49,0.5)	0.5 (0.49,0.5)	0.5 (0.49,0.5)	0.5 (0.49,0.5)
DeepSurv	-0.02 (-0.02,-0.01)	-0.08 (-0.08,-0.08)	-0.26 (-0.27,-0.25)	-0.22 (-0.23,-0.21)	0.51 (0.49,0.53)	0.51 (0.49,0.53)	0.51 (0.49,0.53)	0.51 (0.49,0.53)
DeepHit	<b>0.00 (0.00,0.00)</b>	0.00 (0.00,0.00)	-0.01 (-0.01,0.00)	<b>0.00 (0.00,0.00)</b>	0.45 (0.43,0.47)	<b>0.6 (0.57,0.62)</b>	<b>0.63 (0.6,0.65)</b>	0.55 (0.53,0.58)
CBNN	<b>0.00 (0.00,0.00)</b>	-0.01 (-0.01,0.00)	-0.09 (-0.1,-0.09)	-0.24 (-0.25,-0.23)	<b>0.56 (0.54,0.58)</b>	0.56 (0.54,0.58)	0.57 (0.55,0.59)	<b>0.57 (0.55,0.59)</b>
DSM	-0.14 (-0.14,-0.13)	-0.31 (-0.32,-0.29)	-0.52 (-0.55,-0.49)	-0.52 (-0.57,-0.48)	0.42 (0.42,0.43)	0.46 (0.45,0.47)	0.49 (0.48,0.5)	0.52 (0.51,0.53)



where  $\gamma_1 = 3.9, \gamma_2 = 3, \gamma_3 = -0.43, \gamma_4 = 1.33, \gamma_5 = -0.86, \beta_1 = 1, \beta_2 = 1, \beta_3 = 1, \tau_1 = 10, \tau_2 = 2, \tau_3 = 2$  and  $\psi$  are basis splines. To generate the *gamma* coefficients, an initial cubic splines model with three knots is fit using the *flexsurvspline* function with an intercept-only model on the data, which returns the coefficients ( $\gamma$ ) we use in our model. Note that we fix these values for the analysis. The  $\beta$  coefficients represent direct effects,  $\tau_2$  and  $\tau_3$  represent interactions and  $\tau_1$  is a time-varying interaction.

### 3.3.1 | Performance comparison in complex simulation

Figure 2 C, D and Table 1 B show the performance over time on a test set in the complex simulation. Apart from the Optimal regression model, CBNN outperforms the competing models when examining IPA and up to the 75-percentile of follow-up time for  $C_{IPCW}$ . We expected the Optimal model to perform best in both metrics. However, this was not the case for  $C_{IPCW}$  and may be due to an artifact of concordance-based metrics, where a misspecified model may perform better than a correctly specified one<sup>30</sup>. We attribute the performance of CBNN to its flexibility in modeling time-varying interactions and baseline hazard, flexibility the other neural network models do not have.

## 4 | APPLICATION TO SUPPORT AND METABRIC DATA

Our complex simulation demonstrates the superior performance of CBNN in ideal conditions with clean data. To obtain a more realistic performance assessment, we compared models using two real datasets with a time-to-event outcome. The first case study examines the SUPPORT dataset<sup>6</sup>. The second case study examines the METABRIC dataset<sup>10</sup>. We use the same hyperparameters as in the simulation studies. As we do not know the true model for the real data, we exclude the Optimal model. We split the datasets keeping 20% of the observations as a test set. 20% of the training set is kept aside for validation at each epoch. We predict risk functions for everyone in the test set, which is used to calculate our metrics. We conduct 100 bootstrap re-samples for the real data applications to obtain confidence intervals.

### 4.1 | Performance evaluation using the SUPPORT dataset

The SUPPORT dataset tracks the time until death for seriously ill patients at five American hospitals<sup>6</sup>. We use a pre-processed version of the dataset made available in the DeepSurv package<sup>5</sup>. This dataset contains 9104 samples and 14 covariates (age, sex, race, number of comorbidities, presence of diabetes, presence of dementia, presence of cancer, mean arterial blood pressure, heart rate, respiration rate, temperature, white blood cell count, serum's sodium and serum's creatinine)<sup>5</sup>. Patients with missing features were excluded and 68.10% of the patients died during the 5.56-year study period<sup>5</sup>.

Figure 2 E, F and Table 1 C demonstrates the performance over time on a test set. The regression models (CBLR, Cox) perform best considering IPA, followed by CBNN from the 25<sup>th</sup> to 100<sup>th</sup> percentile of follow-up time. We note a drop in performance for CBNN before the 25<sup>th</sup> percentile of follow-up. For  $C_{IPCW}$ , CBNN outperforms the competing models consistently over follow-up time. Note that performance is similar for all models, aside from DeepSurv whose  $C_{IPCW}$  is lower than the rest (Figure 2 E, F and Table 1 C).

### 4.2 | Performance evaluation using the METABRIC dataset

METABRIC is a 30-yearlong study aiming to discover the molecular drivers of breast tumors, following 2000 individuals with breast cancer until death<sup>10</sup>. They described these growths as primary invasive breast carcinomas, with the goal of discovering both genetic and clinical risk factors for breast cancer survival<sup>10</sup>. We used the processed dataset made available through DeepSurv<sup>5</sup>, which includes 1980 patients of which 57.72% die due to breast cancer within a median 10 years of follow-up<sup>5</sup>. There are 11 covariates in total: 4 RNA-Seq gene expressions (MKI67, EGFR, PGR and ERBB2) and 5 clinical features (hormone treatment indicator, radiotherapy indicator, chemotherapy indicator, ER-positive indicator and age at diagnosis)<sup>5</sup>.

Figure 2 G, H and Table 1 D shows the performance on a test set over time. The IPA scores suggest that regression models outperform competing models on this dataset, as all the neural network models are equal to or perform worse than KM over follow-up time. Our CBNN model is comparable to KM until around the 50-percentile of follow-up time, after which CBNN, DeepSurv and DSM drop in performance. The  $C_{IPCW}$  produces a different ranking. Our CBNN model outperforms the other models up to the 25<sup>th</sup> percentile of follow-up time, whereas DeepHit performs best from the 50<sup>th</sup> to 75<sup>th</sup> percentile of follow-up. With this metric, CBNN and DeepHit outperform the regression models. This disagreement between IPA and  $C_{IPCW}$  may be

due to the misspecification issue of concordance-based metrics<sup>30</sup>. The neural network models may be over-parameterized, as shown by the wide confidence bands in  $C_{IPCW}$ .

## 5 | DISCUSSION

CBNN models survival outcomes by using neural networks on case-base sampled data. We incorporate follow-up time as a feature, providing a data-driven estimate of a flexible baseline hazard and time-varying interactions in our hazard function. The three competing neural network models we evaluated cannot model time-varying interactions by design<sup>11,59</sup>. DSM requires a mixture component distributions for a flexible baseline hazard to be fit<sup>11</sup>. As our goal is to fix the design and compare performance, we did not change any shared hyperparameters. With our choice of shared hyperparameters and model design, DSM did not converge in the complex simulation and SUPPORT case study. Despite this limitation, we include this method when it did function as it can fit flexible baseline hazards. Compared to CBNN, the remaining two models also have limitations. DeepSurv is a proportional hazards model and does not estimate the baseline hazard<sup>5</sup>. DeepHit requires an alpha hyperparameter, is restricted to a single distribution for the baseline hazard and models the survival function directly<sup>9</sup>. The alternative neural network methods match on time, while CBNN models time directly.

To assess performance among these models, we use both IPA and  $C_{IPCW}$  metrics. Concordance-based measures are commonly used in survival analysis to compare models and we opt to keep them in our analyses. However,  $C_{IPCW}$  is a non-proper metric and may cause misspecified models to appear better than they should<sup>30</sup>. Therefore, we contextualize our  $C_{IPCW}$  results in relation to IPA, a proper scoring rule. The model rankings between  $C_{IPCW}$  and IPA differed for both the complex simulation and METABRIC application. Wider  $C_{IPCW}$  confidence intervals for DeepHit and DeepSurv show potential misspecification of the models (Figure 2 D, H).

With this interpretation of  $C_{IPCW}$  in mind, we assess two simulations with minimal noise and two case studies in real scenarios. The simple simulation demonstrates potential pitfalls associated with neural network models, particularly overparameterization. All neural network-based approaches performed worse than null KM model (IPA score), while the regression-based performed better. We attribute this to potential overparameterization in the neural network models, as the wide confidence intervals suggest over-fitting, even with dropout. Both DeepSurv and DSM are affected, while DeepHit and CBNN are less so. This is a limitation to our strategy for evaluating the methods, which uses a fixed study design across all assessments.

From this baseline benchmark, we move on to a complex simulation that requires a method that can learn both time-varying interactions and have a flexible baseline hazard. Here, CBNN demonstrates a distinct advantage over all other methods. The regression models show improved performance over the null KM model, while the competing neural network models performed worse. Based on our complex simulation results (Figure 2 C, D and Table 1 B), CBNN outperforms the competitors when time-varying interactions and a complex baseline hazard are present. This simulation shows how CBNN can perform under ideal conditions, while the following two analyses on real data serve to assess its performance in realistic conditions.

In the SUPPORT and METABRIC case studies, flexibility in both interaction modeling and baseline hazard did not improve CBNNs relative performance, suggesting that this flexibility does not aid prediction in either case study. From a biological perspective, the 30-year follow-up time in the METABRIC study may contain competing causes of death. The causes at the start of the study may not match the causes towards the end, potentially explaining the drop in performance as we reach later survival times with fewer individuals. The baseline hazard is unlikely to cause this drop as DSM is the most flexible competing model in our comparison. Over-fitting is also unlikely given the tight confidence intervals. Further research is required to determine the cause of the drop in performance seen in the METABRIC case study. In both case studies, CBNN outperforms the competing neural network methods.

Taken together, CBNN outperforms all competitors in the complex simulation, demonstrating its value in survival settings that may involve time-varying interactions and a complex baseline hazard. Once we perform case-base sampling and adjust for the sampling bias, we can use a sigmoid activation function to predict our hazard function. Our approach simplifies the incorporation of censored individuals, allowing survival outcomes to be treated as binary ones. Forgoing the requirement of custom loss functions, CBNN only requires the use of standard components in machine learning libraries (specifically, the add layer to adjust for sampling bias and the sigmoid activation function). Due to the simplicity in its implementation and by extension user experience, CBNN is both a user-friendly approach to data-driven survival analysis and is easily extendable to any feed-forward neural network framework.

## References

1. Kleinbaum DG, Klein M. *Survival analysis: a self-learning text*. 3. Springer . 2012.
2. Hanley JA, Miettinen OS. Fitting smooth-in-time prognostic risk functions via logistic regression. *The International Journal of Biostatistics* 2009; 5(1).
3. Coradini D, Daidone MG, Boracchi P, et al. Time-dependent relevance of steroid receptors in breast cancer. *Journal of clinical oncology* 2000; 18(14): 2702–2709.
4. Royston P, Parmar MK. Flexible parametric proportional-hazards and proportional-odds models for censored survival data, with application to prognostic modelling and estimation of treatment effects. *Statistics in medicine* 2002; 21(15): 2175–2197.
5. Katzman JL, Shaham U, Cloninger A, Bates J, Jiang T, Kluger Y. DeepSurv: personalized treatment recommender system using a Cox proportional hazards deep neural network. *BMC Medical Research Methodology* 2018; 18(1): 24.
6. Knaus WA, Harrell FE, Lynn J, et al. The SUPPORT prognostic model: Objective estimates of survival for seriously ill hospitalized adults. *Annals of Internal Medicine* 1995; 122(3): 191–203.
7. Faraggi D, Simon R. A neural network model for survival data. *Statistics in medicine* 1995; 14(1): 73–82.
8. Wang H, Li G. Extreme learning machine Cox model for high-dimensional survival analysis. *Statistics in medicine* 2019; 38(12): 2139–2156.
9. Lee C, Zame WR, Yoon J, Schaar Mvd. DeepHit: A Deep Learning Approach to Survival Analysis With Competing Risks.. *AAAI Conference on Artificial Intelligence* 2018: 2314–2321. Software from pycox version 0.2.2.
10. Curtis C, Shah SP, Chin SF, et al. The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups. *Nature* 2012; 486(7403): 346–352.
11. Nagpal C, Li XR, Dubrawski A. Deep survival machines: Fully parametric survival regression and representation learning for censored data with competing risks. *IEEE Journal of Biomedical and Health Informatics* 2021. Software downloaded March 10, 2021.
12. Saarela O, Hanley JA. Case-base methods for studying vaccination safety. *Biometrics* 2015; 71(1): 42–52.
13. Gulli A, Pal S. *Deep learning with Keras*. Packt Publishing Ltd . 2017.
14. Hughes-Hallett D, Gleason AM, McCallum WG. *Calculus: Single and multivariable*. John Wiley & Sons . 2020.
15. Srivastava N, Hinton G, Krizhevsky A, Sutskever I, Salakhutdinov R. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research* 2014; 15(1): 1929–1958.
16. R Core Team . *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing; Vienna, Austria: 2021.
17. Van Rossum G, Drake FL. *Python 3 Reference Manual*. Scotts Valley, CA: CreateSpace . 2009. Version 3.6.13.
18. Terry M. Therneau , Patricia M. Grambsch . *Modeling Survival Data: Extending the Cox Model*. New York: Springer . 2000. R package version 3.2-11.
19. Bhatnagar SR, Turgeon M, Islam J, Hanley JA, Saarela O. casebase: An Alternative Framework For Survival Analysis and Comparison of Event Rates. *arXiv preprint arXiv:2009.10264* 2020.
20. Sonabend R. *survivalmodels: Models for Survival Analysis*. ; : 2021. R package version 0.1.9.
21. Allaire J, Chollet F. *keras: R Interface to 'Keras'*. ; : 2021. R package version 2.7.0.
22. Brilleman SL, Wolfe R, Moreno-Betancur M, Crowther MJ. Simulating Survival Data Using the simsurv R Package. *Journal of Statistical Software* 2020; 97(3): 1–27. doi: 10.18637/jss.v097.i03

23. Jackson C. flexsurv: A Platform for Parametric Survival Modeling in R. *Journal of Statistical Software* 2016; 70(8): 1–33. R package version 2.0doi: 10.18637/jss.v070.i08
24. Pölsterl S. scikit-survival: A Library for Time-to-Event Analysis Built on Top of scikit-learn. *Journal of Machine Learning Research* 2020; 21(212): 1-6. Version 0.14.0.
25. Gerds TA, Kattan MW. *Medical Risk Prediction Models: With Ties to Machine Learning (1st ed.)*. Chapman and Hall/CRC . 2021. R package version 2021.10.10.
26. Ushey K, Allaire J, Tang Y. *reticulate: Interface to 'Python'.* ; : 2021. R package version 1.22.
27. Kattan MW, Gerds TA. The index of prediction accuracy: an intuitive measure useful for evaluating risk prediction models. *Diagnostic and prognostic research* 2018; 2(1): 1–7.
28. Uno H, Cai T, Pencina MJ, D’Agostino RB, Wei LJ. On the C-statistics for evaluating overall adequacy of risk prediction procedures with censored survival data. *Statistics in medicine* 2011; 30(10): 1105–1117.
29. Graf E, Schmoor C, Sauerbrei W, Schumacher M. Assessment and comparison of prognostic classification schemes for survival data. *Statistics in medicine* 1999; 18(17-18): 2529–2545.
30. Blanche P, Kattan MW, Gerds TA. The c-index is not proper for the evaluation of-year predicted risks. *Biostatistics* 2019; 20(2): 347–357.

