#### RESEARCH ARTICLE



Check for updates

# IDNetwork: A deep illness-death network based on multi-state event history process for disease prognostication

Aziliz Cottin<sup>1</sup> | Nicolas Pecuchet<sup>1</sup> | Marine Zulian<sup>1</sup> | Agathe Guilloux<sup>2,3</sup> | Sandrine Katsahian<sup>4,5,6,7</sup>

<sup>1</sup>Healthcare and Life Sciences Research, Dassault Systemes, Velizy-Villacoublay, France

<sup>2</sup>CNRS, Université Paris-Saclay, Paris, France

<sup>3</sup>Laboratoire de Mathématiques et Modélisation d'Evry, Université d'Evry, Evry-Courcouronnes, France

<sup>4</sup>AP-HP, Hôpital Européen Georges Pompidou, Unité de Recherche Clinique, APHP Centre, Paris, France

<sup>5</sup>Inserm, Centre d'Investigation Clinique 1418 (CIC1418) Epidémiologie Clinique, Paris, France

<sup>6</sup>Inserm, Centre de recherche des Cordeliers, Sorbonne Université, Université de Paris, Paris, France

<sup>7</sup>HeKA, INRIA PARIS, Paris, France

#### Correspondence

Aziliz Cottin, Healthcare and Life Sciences Research, Dassault Systemes, 78140 Velizy-Villacoublay, France. Email: aziliz.cottin@3ds.com Multi-state models can capture the different patterns of disease evolution. In particular, the illness-death model is used to follow disease progression from a healthy state to an intermediate state of the disease and to a death-related final state. We aim to use those models in order to adapt treatment decisions according to the evolution of the disease. In state-of-the art methods, the risks of transition between the states are modeled via (semi-) Markov processes and transition-specific Cox proportional hazard (P.H.) models. The Cox P.H. model assumes that each variable makes a linear contribution to the model, but the relationship between covariates and risks can be more complex in clinical situations. To address this challenge, we propose a neural network architecture called illness-death network (IDNetwork) that relaxes the linear Cox P.H. assumption within an illness-death process. IDNetwork employs a multi-task architecture and uses a set of fully connected subnetworks in order to learn the probabilities of transition. Through simulations, we explore different configurations of the architecture and demonstrate the added value of our model. IDNetwork significantly improves the predictive performance compared to state-of-the-art methods on a simulated data set, on two clinical trials for patients with colon cancer and on a real-world data set in breast cancer.

#### KEYWORDS

deep learning, illness-death process, neural networks, stratified medicine, survival analysis

# 1 | INTRODUCTION

Disease prognosis is of major importance for physicians when making medical decisions and requires specialized algorithms to estimate the risks of a patient. In this line of work and within event history analysis, we propose a novel algorithm for individual prognostication in a three-state illness-death model.

Event history analysis, also known as survival analysis, aims at predicting the time until the occurrence of a future event(s) of interest and is used in multiple areas including healthcare, economy, finance, and engineering. In particular, survival analysis is widely used in healthcare to model patient survival outcome in order to understand disease progression. In clinical practice, clinicians may be more interested in the complete evolution of a disease and not only in a unique or composite event. The multi-state<sup>1</sup> approach has been developed as a generalization of survival analysis when multiple events can occur successively over time. In the present work, we focus on the illness-death model which is a multi-state model composed of three states: "healthy"; "relapsed" or "diseased"; "dead." Illness-death model is the most frequent

structure used to follow the evolution of cancer patients through an intermediate non-fatal relapse state and a death state as in ovarian cancer<sup>2</sup> or in chronic myeloid leukemia.<sup>3</sup> Other applications of the illness-death model include Alzheimer's disease<sup>4</sup> and cardiovascular disease.<sup>5</sup>

There are two main literature streams of event history analysis. The first is based on traditional statistic theories including three approaches. (i) The nonparametric approaches include the Kaplan Meier<sup>6</sup> and the Nelson-Aalen<sup>7</sup> estimators and do not enable individual prediction. (ii) The parametric models impose a precise form of the hazard rate and the influence of covariates on it but allow individual prediction. (iii) Finally, semi-parametric models can be viewed as a compromise between nonparametric and parametric approaches. They include the widely used Cox<sup>8</sup> proportional hazard (P.H.) model and have been extended to multi-state analysis. However, these models style assume strong assumptions on the relation between the covariates and the event times distribution. These assumptions have shown limitations in many real-world applications and have been found to be violated especially in clinical areas.

To address these challenges, a second variety of literature proposes new machine learning algorithms. In particular, neural networks have been developed to extend the Cox P.H. model in a statistical assumption-free framework. One of the advantages of neural networks is that they can fit highly nonlinear patterns in the data by using multiple layers and nonlinear activation functions. In addition, the use of fully connected networks allows to take into consideration possible interactions between covariates with no prior assumption. Hence, traditional artificial neural networks have been successfully introduced for survival analysis by Faraggi and Simon. More recently, deep neural networks have been extended by Luck et al, tatzman et al, Fotso, Kvamme et al, among others. He model within the gradients' survival as compared to the Cox P.H. model. Regardless, their approaches are still limited to the case of a unique clinical event. Lee et al extended deep neural networks to handle competing events. To the best of our knowledge, no nonlinear methods, especially deep neural networks, have been explicitly extended for multi-state analysis and in particular for an illness-death process. Thus, addressing the linear limitation of the Cox P.H. model within the illness-death modeling framework is one of the focuses of this work.

While the Faraggi and Simon<sup>10</sup> approach uses neural networks to parameterize the Cox's linear regression function, most of the recent methods directly predict a discrete-time distribution of the event times as an output of the neural network. As an approximation for continuous-time survival data, they all perform a division of the continuous time scale into discrete-time intervals. Alternatively to discrete-time approach, piecewise survival models<sup>20</sup> perform a discretization of the time scale but each subject's duration of exposure during the interval is taken into account. Hence, the approximation error that arises when a discrete-time method is used can be reduced with piecewise approximations. See the work of Kvamme and Borgan<sup>21</sup> for more details. Thus, developing a well approximated continuous-time method instead of a discrete-time method is the second objective of this article.

There are few recommendations for dividing intervals and selecting the interval cutpoints in piecewise survival models. This has a significant impact on the model performance and can cause either over-fitting (for a larger number of intervals) or under-fitting. However, Kvamme and Borgan<sup>21</sup> conduct a simulation study by varying the interval cutpoints determination method in a discrete time approach. Inspired by their work, we discuss some methods for determining the optimal interval cutpoints. In addition, we propose a regularization method in order to minimize the risk of over-fitting related to the number of intervals.

In the present work, we propose a deep learning architecture, illness-death network (IDNetwork), for illness-death model that encompasses a multi-task neural network including one subnetwork shared for all the transitions and three transition-specific subnetworks. After a presentation of the classical illness-death model, we (i) derive a new form of the log-likelihood of a piecewise constant illness-death process, (ii) build the network architecture IDNetwork, (iii) implement in Python the pipeline of our method including performances criteria evaluation. We finally conduct experiments on a simulated nonlinear data set and on real data sets of patients with colon cancer and with breast cancer.

# 2 | THE ILLNESS-DEATH MODEL

In this section, we introduce the traditional illness-death model, see the work of Andersen et  $al^{22}$  for a complete presentation.

## 2.1 | Notation

Multi-state models are a generalization of survival models when multiple events of interest can occur over time. A multi-state process  $(E(t), 0 < t \le +\infty)$  is a continuous-time stochastic process that describes the states occupied by a patient over time. In this article, we consider an illness-death multi-state process with three states: state 0 is the initial state "Event-free," state 1 is an intermediate state "Relapse," state 2 is an absorbing state "Death." The illness-death process is characterized by three irreversible transitions: from 0 to 1  $(0 \to 1)$ , from 0 to 2  $(0 \to 2)$ , from 1 to 2  $(1 \to 2)$ , where transitions from state 0 are competing, transitions  $0 \to 1$  and  $1 \to 2$  are successive. It assumes that all subjects are in state 0 at time t = 0 (ie,  $\mathbb{P}(E(0) = 0) = 1$ ).

The evolution of an illness-death process can be characterized by three random variables (r.v.)  $T_{ql}$ , for  $(q, l) \in \{(0, 1), (0, 2), (1, 2)\}$ , associated with each of the three transitions. They represent the transition times from state q to state l ( $q \neq l$ ). A subject leaving state 0 will enter either state 1 at time  $T_{01}$  or state 2 at time  $T_{02}$ . A subject having state 1 at  $T_{01}$  will enter in state 2 at time  $T_{01} + T_{12}$ . This can be summarized as follows.  $T_{01}$ , the exit time from state 0 is

$$T_0 = \inf_{t>0} \{E(t) \neq 0\} = \min(T_{01}, T_{02})$$

and is recorded together with  $D_0 \in \{1, 2\}$  which indicates the entered state.  $T_2$  the entry time to state 2

$$T_2 = \inf_{t>0} \{E(t) = 2\} = T_0 + \mathbb{1} \{D_0 = 1\} T_{12}$$

and characterizes the total survival time.

In clinical settings, the true transition times are commonly partially observed because of right-censoring. To model this phenomenon, we introduce C a non-negative censoring r.v. that precludes its observation. Let  $\tilde{T}_0 = \min{(T_0, C)}$  and  $\tilde{T}_2 = \min{(T_2, C)}$  be the observed event times. Together with these event times, we observe a vector of covariates X of dimension P and we assume that  $C \perp \perp (T_0, T_2) \mid X$ . We also observe the binary labels  $\delta_{0l} = \mathbb{1}$   $\{D_0 = l, T_0 \leq C\}$   $\{l = 1, 2\}$ ,  $\delta_{12} = \delta_{01} \mathbb{1}$   $\{T_2 \leq C\}$  that indicate the status of the transitions, where  $\delta_{ql} = 1$  indicates an entry in state l from q and  $\delta_{ql} = 0$  indicates a censored transition.

## 2.2 Transition intensities and the Cox P.H. model

Illness-death processes are traditionally described with counting processes.<sup>23</sup> In an illness-death model, the observation of process E is equivalent to the observation of the three-variate process  $t \ge 0 \mapsto N(t) = (N_{01}(t), N_{02}(t), N_{12}(t))$ . where

$$N_{ql}(t) = \text{card} \{0 \le s \le t : E(s-) = q, E(s) = l\}, \text{ for } (q, l) \in \{(0, 1), (0, 2), (1, 2)\},\$$

the transition from q to l happened before t. To take into account the presence of censoring, we define in addition two previsible processes  $Y_0$  and  $Y_1$  as

$$Y_q(t) = 1 \{ E(t-) = q \}, \text{ for } q = 0, 1, t > 0,$$

they indicate if the patient is in state q before time t.

The three-variate counting process N is conventionally associated with a set of transition intensities, or transition-specific hazard functions. We define the three processes by making specific assumptions on cancer evolution over time.<sup>24</sup> For the transitions  $0 \to 1$  and  $0 \to 2$ , we consider time nonhomogeneous Markovian processes. For transition  $1 \to 2$ , we perform a time transformation, following Anderson et al,<sup>22</sup> and we consider a time homogeneous semi-Markovian process (the probability of transiting from state 1 to state 2 at time t depends only on the duration  $t - T_0$  already spent in 1). Wherever convenient, we use the duration variable t and t instead of the time variable t for the transition t and t in the transition t and t instead of the time variable t for the transition t and t instead of the time variable t for the transition t and t in the transition t and t in the transition t

$$\begin{split} \alpha_{0l}^*(t|X) &= \lim_{h \to 0} \frac{1}{h} \mathbb{P}\left( E(t+h) = l \mid E(t-) = 1, X \right), \; \text{ for } \; l = 1, 2, \\ \alpha_{12}^*(d|X) &= \lim_{h \to 0} \frac{1}{h} \mathbb{P}\left( E(d+h) = l \mid E(d-) = 1, X \right), \end{split}$$

on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons License

that express the instantaneous risk to transit from a state to another at a specific time for the covariates X.

Transition specific Cox P.H. models have been proposed to model multi-state processes in order to evaluate covariates effects on each transition. Covariates effects are introduced in the transition intensities by means of log linear risk function as follows:

$$\alpha_{0l}^*(t|X) = \alpha_{0l}^{*,0}(t) \exp(X\beta_{0l}^*), \text{ for } l = 1, 2,$$
 (1)

$$\alpha_{12}^*(d|X) = \alpha_{12}^{*,0}(d) \exp\left(X\beta_{12}^*\right),\tag{2}$$

where, for  $(q, l) \in \{(0, 1), (0, 2), (1, 2)\}$ ,  $\alpha_{ql}^{*,0}$  is the unknown baseline transition intensity related to the transition  $q \to l$  (ie, the underlying hazard when all the covariates are equal to zero) and  $\beta_{ql}^*$  is the vector of regression coefficients to be estimated.

In general, the process associated with each of the three transitions may depend on the time t of arrival in the state (Markov process) or on the time d since the entry to the state (semi-Markov process). The choice of the time scale<sup>3</sup> is an important step in disease modeling because it will induct how the disease will evolve over time. We discuss our choice in Section 3.1.

# 2.3 | The log-likelihood

The log-likelihood associated with the observation of the three-dimensional censored counting process  $t \mapsto N(t \wedge C) = (N_{01}(t \wedge C), N_{02}(t \wedge C), N_{12}(t \wedge C))$  on  $[0, \tau]$ , where  $\tau$  is the horizon time, is given by the log product of the three transition-specific likelihoods:

$$\log \mathcal{L} = \log \left( \mathcal{L}^{0 \to 1} \times \mathcal{L}^{0 \to 2} \times \mathcal{L}^{1 \to 2} \right),$$

where  $\mathcal{L}^{0\to 1}$ ,  $\mathcal{L}^{0\to 2}$ ,  $\mathcal{L}^{1\to 2}$  are the likelihoods associated with each transition, see Andersen et al,<sup>22</sup> and the covariates X. They are given, for the proposals  $\alpha_{0l}$  and  $\alpha_{12}$ , for l=1,2, by the following equations:

$$\log \mathcal{L}^{0 \to l} = \int_0^{\tau} \log \left( \alpha_{0l}(t|X) \right) Y_0(t) \mathbb{1} \{ C \ge t \} dN_{0l}(t) - \int_0^{\tau} \alpha_{0l}(t|X) Y_0(t) \mathbb{1} \{ C \ge t \} dt, \text{ for } l = 1, 2,$$
 (3)

$$\log \mathcal{L}^{1\to 2} = \int_0^{\tau} \log \left(\alpha_{12}(d|X)\right) Y_1(t) \mathbb{1}\{C \ge d\} dN_{12}(t) - \int_0^{\tau} \alpha_{12}(d|X) Y_1(t) \mathbb{1}\{C \ge d\} dd. \tag{4}$$

The log-likelihood of Equation (3) is traditionally used to estimate the coefficients  $\beta_{ql}^*$  while the unknown functions  $\alpha_{ql}^{*,0}$  can be estimated with Nelson-Aalen<sup>9</sup> or spline estimators.<sup>25</sup> Other models have been proposed for the transition intensities,<sup>26</sup> the simplest assuming constant transition intensities  $\alpha_{ql}^*(t) = \alpha_{ql}^*$ . Finally, models with time varying coefficients, with proposals defined in Equations (1) and (2), have been studied, see Martinussen and Scheike<sup>27</sup> for a complete review. Among these proposals, Murphy and Sen<sup>28</sup> considered piecewise constant estimators that allow for time-dependent covariate effects.

Under the assumption of a piecewise constant model, the optimization of the log-likelihood is facilitated. Indeed, they allow to perform a continuous model that would have necessitated the use of the Cox partial log-likelihood<sup>8</sup> and significantly impact on the computational cost when using classical stochastic gradient descent algorithm because of the presence of two cumulative sums. <sup>14,29</sup> Whereas a discrete-time model approximates the full log-likelihood by dividing the time axis into discrete time intervals. Piecewise constant approaches are a compromise to reduce the approximation error that arise when a discrete-time method is used by computing the cumulative functions by taking into account the subject's duration of exposure in the intervals, with no supplementary computational cost, see the Supplementary materials (Section 3) for details.

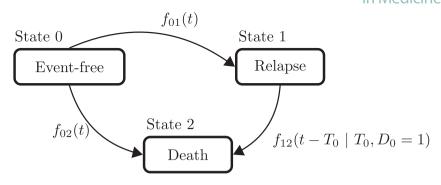


FIGURE 1 Illustration of our illness-death process

## 3 | OUR METHODOLOGY

In this section, we present our approach to model an illness-death process.

## 3.1 | Functions of interest and data

Instead of considering the transition intensities, we will write the log-likelihood in terms of transition probabilities  $f_{ql}^*(.|X)$ , that we define in this paragraph. See an illustration in Figure 1.

Let  $f_{01}^*$  and  $f_{02}^*$  be the infinitesimal probabilities of experiencing, respectively, transitions  $0 \to 1$  and  $0 \to 2$ , defined as

$$f_{0l}^*(t|X) = \lim_{h \to 0} \frac{1}{h} \mathbb{P}\left(E(t+h) = l, E(t-) = 0|X\right) = \lim_{h \to 0} \frac{1}{h} \mathbb{P}\left(t \le T_0 \le t + h, D_0 = l|X\right),$$

and  $F_{01}^*$  and  $F_{02}^*$  their cumulative counterparts,

$$F_{0l}^*(t|X) = \mathbb{P}\left(E(t) \neq l, D_0 = l|X\right) = \mathbb{P}\left(T_0 \leq t, D_0 = l|X\right) = \int_0^t f_{0l}^*(s|X)ds, \text{ for } l = 1, 2 \text{ and } t \geq 0$$

expresses the probability that a transition  $0 \to l$  occurs on or before time t. With these definitions, we define

$$f_0^*(t|X) = f_{01}^*(t|X) + f_{02}^*(t|X)$$

as the infinitesimal probability of exiting state 0 at time t and

$$F_0^*(t|X) = F_{01}^*(t|X) + F_{02}^*(t|X)$$

as the probability of having exited state 0 before time *t*.

For the transition  $1 \to 2$ , the functions of interest are defined conditionally to  $T_0$ ,  $D_0 = 1$ . To simplify the notations, we drop this conditioning in the definitions. We define  $f_{12}^*$  as the infinitesimal probability of experiencing transition  $1 \to 2$  such that,

$$f_{12}^*(d|X) := f_{12}^*(d|T_0, D_0 = 1, X) = \lim_{h \to 0} \frac{1}{h} \mathbb{P}(E(d+h) = 2, E(d-) = 1|X) = \lim_{h \to 0} \frac{1}{h} \mathbb{P}(d \le T_2 - T_0 \le d + h|X),$$

and  $F_{12}^*$  as its cumulative counterpart such that,

$$F_{12}^*(d|X) := F_{12}^*(d|T_0, D_0 = 1, X) = \mathbb{P}(E(d) = 2|X) = \mathbb{P}(T_2 - T_0 \le d|X) = \int_0^d f_{12}^*(s|X)ds, \text{ for } d \le 0.$$

 $F_{01}^*$ ,  $F_{02}^*$ , and  $F_{12}^*$  are commonly referred to cumulative incidence functions (CIFs) in the literature<sup>30</sup> and are our main functions of interest. Indeed, for a patient with covariates X, if he/she still is state, the physician will need to estimate precisely the values  $F_{01}^*(t_0|X)$  and  $F_{02}^*(t_0|X)$  where  $t_0$  is a certain horizon time, chosen according to the pathology. Even if they are formula linking the transition intensities<sup>23</sup> to these probabilities, see Supplementary materials (Section 1), we chose to work with transition probabilities.

Our data consist in the observation of n independent and identically distributed r.v. in  $\mathbb{R}^P \times \mathbb{R}_+ \times \{0,1\} \times \{0,1\} \times \mathbb{R}_+ \times \{0,1\}$  such that

$$\mathcal{D}_i = \{X_i, (\tilde{T}_0^i, \delta_{01}^i, \delta_{02}^i), (\tilde{T}_2^i, \delta_{12}^i)\}_{1 \le i \le n},$$

where  $X_i = (X_{i1}, ..., X_{iP})^T$  is a vector of P covariates observed at baseline. From these observations and for each subject i, we aim to estimate the true transition specific density probabilities conditionally to the clinical features  $X_i$  in order to predict the individual CIFs defined in Equations (9) and (10).

# 3.2 Writing of the log-likelihood in terms of the functions of interest

We show that the conventional illness-death log-likelihood defined in Equation (3) can be rewritten in terms of the density probability functions introduced in Section 3.1, see the Supplementary materials (Section 2) for a formal proof.

We define the log-likelihood  $\ell_n$  by dividing the contributions in three distinct parts:

$$\ell_n = \log \mathcal{L}_n = \log \mathcal{L}_n^{0 \to 1} + \log \mathcal{L}_n^{0 \to 2} + \log \mathcal{L}_n^{1 \to 2}$$
(5)

$$= \frac{1}{n} \sum_{i=1}^{n} \left[ \mathcal{E}_{i}^{0 \to 1} + \mathcal{E}_{i}^{0 \to 2} + \mathcal{E}_{i}^{1 \to 2} \right], \tag{6}$$

where  $\ell_i^{0\to 1} + \ell_i^{0\to 2}$  is the log contribution of patient i from state 0. From state 0, patient i with an event a  $\tilde{T}_0^i$  can contribute in three ways. He (i) can experience a transition  $0\to 1$ ; (ii) can experience a transition  $0\to 2$ ; (iii) can be censored at  $\tilde{T}_0^i$ . Thus  $\ell_0^i$  is given, for proposals  $f_{01}$ ,  $f_{02}$ ,  $F_0$ , and binary labels  $\delta_{01}$ ,  $\delta_{02}$  that indicate an entry in another state or a censored transition, by

$$\mathcal{\ell}_{i}^{0 \to 1} + \mathcal{\ell}_{i}^{0 \to 2} = \sum_{l=1,2} \{ \delta_{0l}^{i} \log(f_{0l}(\tilde{T}_{0}^{i} | X_{i})) \} + (1 - (\delta_{01}^{i} + \delta_{02}^{i})) \log(1 - F_{0}(\tilde{T}_{0}^{i} | X_{i})). \tag{7}$$

On the other hand,  $\mathcal{C}_i^{1\to 2}$  is the log contribution of patient i from the time he has entered state 1 (only for i such that  $\delta_{01}^i=1$ ). Following the previous reasoning, patient i can contribute in two ways at time  $\tilde{T}_2^i-\tilde{T}_0^i$ . He (i) can experience a transition  $1\to 2$ ; (ii) can be censored. Thus  $\mathcal{C}_i^{1\to 2}$  is given, for proposals  $f_{12}$ ,  $F_{12}$ , and the binary label  $\delta_{12}$  that indicate an entry in state 2 or a censored transition, by

$$\mathcal{E}_{i}^{1 \to 2} = \delta_{01}^{i} \delta_{12}^{i} \log(f_{12}(\tilde{T}_{2}^{i} - \tilde{T}_{0}^{i} | X_{i})) + \delta_{01}^{i} (1 - \delta_{12}^{i}) \log(1 - F_{12}(\tilde{T}_{2}^{i} - \tilde{T}_{0}^{i} | X_{i})). \tag{8}$$

The log-likelihood in Equation (5) has been derived under a time nonhomogeneous Markovian assumption for transitions  $0 \to 1$  and  $0 \to 2$ , and a time homogeneous semi-Markovian assumption for transition  $1 \to 2$ . Depending on the application, time assumptions can be reformulated, it would change the log likelihood above.

# 3.3 Rewriting of the log-likelihood for piecewise constant proposals

We propose in the present article to consider the class of piecewise constant transition-specific probabilities following the ideas of Murphy and Sen.<sup>28</sup> More specifically candidate estimators will be assumed to be constant on K time intervals. We define  $\tau$  as the maximum horizon time window and we divide the time axis into K disjoint time intervals:

 $v_1 = [a_0, a_1), \dots, v_K = [a_{K-1}, a_K)$ , with  $a_0 = 0$  and  $a_K = \tau$ . For any time  $s \in [0, \tau[$ , we denote by  $v_{k(s)}$  the time interval to which s belongs. so that :

$$f_{0l}(t|X) = f_{0l}(v_{k(t)}|X)$$
, for  $l = 1, 2$ ,  
 $f_{12}(t|X) = f_{0l}(v_{k(d)}|X)$ .

Piecewise constant models allow to keep constant the probabilities within observed time intervals. In this model, effects of the covariates are time-dependent.<sup>31</sup>

We rewrite the log-likelihood of an illness-death process for piecewise constant proposals. It has to be understood as an approximation and it is well-established that the approximation error can be bounded when the true functions are smooth.<sup>32</sup>

Assuming a piecewise constant model allows to consider a nonhomogeneous model (see Section 3.1) while retaining the hypothesis of homogeneity within the same time interval, following the lines of Kvamme and Borgan,<sup>21</sup> Friedman et al.<sup>20</sup> Its means that the density probabilities are constant within each interval but covariates effects are parameterized to be different in each interval (see Section 2.2). Hence, we can express  $f_{ql}$  (q = 0, 1 and l = 1, 2) as step functions such that  $f_{ql}(\cdot, |X|) = f_{ql}(v_{k(\cdot)}, |X|)$  that depend on the covariates and on time. Equivalently, since the density probabilities are assumed to be piecewise constant, their corresponding cumulative counterparts  $F_{0l}$  (l = 1, 2) and  $F_{12}$  are piecewise linear:

$$F_{0l}(t \mid X) = \sum_{k=1}^{k(t)-1} f_{0l}(v_k \mid X) |v_k| + (t - a_{k(t)-1}) f_{0l}(v_{k(t)} \mid X), \tag{9}$$

$$F_{12}(d \mid X) = \sum_{k=1}^{k(d)-1} f_{12}(\nu_k \mid X) |\nu_k| + (d - a_{k(d)-1}) f_{12}(\nu_{k(d)} \mid X), \tag{10}$$

where  $|v_k|$  is the length of interval  $v_k$ . In Equations (9) and (10), we see that each duration of exposure during the intervals is taken into account. In comparison, only whether an event occurred or not in a given time interval is taken into account in discrete-time models, disregarding the duration of exposure in the given time interval.

In real clinical data, the r.v.  $T_0$ ,  $T_2$  can take values after  $\tau$  (a patient can leave state 0 or state 1 after  $\tau$ ). Hence, under the piecewise constant assumption, the following equations are satisfied:

$$\sum_{k=1}^{K} f_0(\nu_k|X)|\nu_k| + 1 - F_0(\tau|X) = 1,$$

$$\sum_{k=1}^{K} f_{12}(\nu_k|X)|\nu_k| + 1 - F_{12}(\tau|X) = 1,$$
(11)

where  $|v_k|$  is the length of interval  $v_k$ . See Lee et al<sup>18</sup> and Kvamme and Borgan<sup>21</sup> for similar remarks.

Regarding the constraint of Equation (11), we consider a supplementary interval  $v_{K+1} = [\tau, +\infty)$ . Consequently, we fulfill the constraint by defining  $1 - F_0(\tau|X) = f_0(v_{K+1}|X)$ ,  $1 - F_{12}(\tau|X) = f_{12}(v_{K+1}|X)$ .

Under the assumption of a piecewise constant model, we rewrite the log-likelihood in Equation (5). We rewrite the sum of the first two contributions in Equation (7) as follows:

$$\mathcal{\ell}_{i}^{0 \to 1} + \mathcal{\ell}_{i}^{0 \to 2} = \sum_{l=1,2} \left\{ \delta_{0l}^{i} \log \left( f_{0l} \left( \nu_{k(\tilde{T}_{0}^{i})} \mid X_{i} \right) \right) \right\} + (1 - (\delta_{01}^{i} + \delta_{02}^{i})) \log (1 - F_{0}(\tilde{T}_{0}^{i} \mid X_{i})).$$

We rewrite the third contribution in Equation (8) as follows:

$$\mathcal{E}_i^{1 \to 2} = \delta_{01}^i \delta_{12}^i \log(f_{12}(v_{k(\tilde{T}_2^i - \tilde{T}_0^i)} \mid X_i)) + \delta_{01}^i (1 - \delta_{12}^i) \log(1 - F_{12}(\tilde{T}_2^i - \tilde{T}_0^i \mid X_i)).$$

Estimators of  $f_{0l}$  and  $f_{12}$  will be precisely defined in Section 4.1.

FIGURE 2 IDNetwork architecture. "FC layer" refers to fully connected layer

## 4 | DESCRIPTION OF IDNETWORK

In this section, we describe how we use a new deep learning approach to parameterize the step probability functions  $f_{01}$ ,  $f_{02}$ ,  $f_{12}$  over the interval  $[0, \tau]$ . Our deep learning architecture, called IDNetwork, model the relationship between covariates and these transition probabilities with no linear assumption. Unlike classical methods,  $^{9,33}$  IDNetwork is divided into transition-specific tasks and uses nonlinear activation functions to capture nonlinearity between covariates and transition probabilities. We propose a loss function that encompasses the negative log-likelihood of Equation (5) and that is tuned to automatically choose a good time division K (ie, the number of time intervals) in order to minimize the risk of over-fitting. In addition, we propose two methods to select the interval cutpoints.

# 4.1 | Network architecture

Inspired by the work of Lee et al<sup>18</sup> and Fotso,<sup>13</sup> we develop an architecture (see an illustration in Figure 2) with three task-specific subnetworks that are related to the three transitions of an illness-death process. Multi-task learning is done with hard parameter sharing<sup>34</sup> in order to extract common and specific patterns from the patient's characteristics (ie, the baseline covariates). It is composed of a first subnetwork shared between the three transitions and of three transition-specific subnetworks. Two different softmax output layers are used to transform the transition-specific subnetworks outputs into time-dependent probabilities. One softmax layer is related to the exit from state 0 (ie, the transitions  $0 \to 1$  and  $0 \to 2$ ), the other to the exit from state 1 (ie, the transition  $1 \to 2$ ).

# 4.1.1 | Input layer

The input layer is composed of the matrix X of P baseline covariates for the n individuals.

# 4.1.2 | Covariates shared subnetwork

The shared subnetwork takes as input the *input layer* and contains L fully connected (FC) hidden layers with l units. Its output is a vector  $\mathbf{z} = g^{\text{input}}(X)$  in  $\mathbb{R}^l$  that captures shared patterns between the three transitions ( $g^{\text{input}}$  is a nonlinear activation function).

# 4.1.3 | Transition-specific subnetworks

Each transition-specific subnetwork takes  $\mathbf{z}$  as input and contains  $L^{ql}$  fully connected hidden layers with  $l^{ql}$  units. Its output is a vector  $\mathbf{y}_{ql} = g^{ql}(\mathbf{z})$ , that is a transition-specific transformation of the shared features ( $g^{ql}$  is a nonlinear activation function). Given the unbalanced number of observations for the three transitions that may exist in real data, the range of model complexity is different for each of the three transitions. If the architecture is too complex for a transition, the model will poorly be generalist on new data (over-fitting). If the architecture is not complex enough for a transition, the model will not capture all the information in the data (under-fitting). To find the best configuration that will result in the best model performance for each of the three transitions, we will set the structure of each subnetwork independently. Hence, a transition with a high number of observations (oftenly the transition  $0 \to 1$ ) can support more hidden layers than a transition with less observations (oftenly the transition  $0 \to 2$ ).

# 4.1.4 | Probabilistic output layers

The output of the network is composed of two probabilistic layers that map the transition-specific outcomes  $\mathbf{y}_{ql}$  into time-dependent probabilities. Each output of the network is a fully connected layer. The first output layer is related to transitions  $0 \to 1$  and  $0 \to 2$ , the second output layer is related to transition  $1 \to 2$ . Each output layer is built in two steps:

- 1. Each uses primarily a fully connected layer with a linear activation function, noted  $g^{linear}$ , to transform the transition-specific outputs  $\mathbf{y}_{ql}$  into vectors of lengths K+1 (ie, the number of time intervals). The first output layer (resp. the second output layer) transform the transition-specific outputs  $(\mathbf{y}_{01}, \mathbf{y}_{02})^T$  (resp. the transition-specific  $\mathbf{y}_{12}$ ) into vectors  $(\phi_{01}, \phi_{02})^T$  (resp.  $\phi_{12}$ ), each of length K+1.
- 2. Subsequently, each output layer uses a fully connected layer with a weighted softmax activation function. The first output layer (resp. the second output layer) uses a weighted softmax activation function  $\sigma_0$  (resp.  $\sigma_2$ ) to transform  $(\phi_{01}, \phi_{02})^T$  (resp.  $\phi_{12}$ ) into probabilities and provide an estimation of  $f_{01}^*, f_{02}^*$  (resp.  $f_{12}^*$ ). The use of softmax activation functions ensure the fulfillment of the model constraint defined in Equation (11). Each softmax function is weighted by the length of the time intervals to provide an estimation of the density probabilities under the assumption of a piecewise constant model.
- 3. Finally the output layers are characterized by the vectors:

$$\begin{split} \hat{\mathbf{f}}_0 &= (\hat{\mathbf{f}}_{01}, \hat{\mathbf{f}}_{02})^T = \sigma_0 \left( \mathbf{g}^{\text{linear}} \left( (\mathbf{y}_{01}, \mathbf{y}_{02})^T \right) \right) = \sigma_0 \left( (\phi_{01}, \ \phi_{02})^T \right), \\ \hat{\mathbf{f}}_{12} &= \sigma_2 \left( \mathbf{g}^{\text{linear}} \left( \mathbf{y}_{12} \right) \right) = \sigma_2 \left( \phi_{12} \right), \end{split}$$

where

$$\hat{\mathbf{f}}_{ql} = (\hat{f}_{ql}(\nu_k \mid X))_{0 < k \leq K+1}, \text{ for } (q, l) \in \{(0, 1), (0, 2), (1, 2)\},$$

and

$$\begin{split} \hat{f}_{0l}(v_k \mid X) &= \frac{\exp\left[\phi_{0l}^k(X)\right]}{\sum_{j=1}^{K+1} \left(\exp\left[\phi_{01}^j(X)\right] + \exp\left[\phi_{02}^j(X)\right]\right) \left|v_j\right|}, \text{ for } l = 1, 2, \\ \hat{f}_{12}(v_k \mid X) &= \frac{\exp\left[\phi_{12}^k(X)\right]}{\sum_{j=1}^{K+1} \exp\left[\phi_{12}^j(X)\right] \left|v_j\right|}. \end{split}$$

# 4.2 | Loss function and mitigation of the number of time intervals effect via penalization

To learn IDNetwork parameters, we minimize a total loss function,

$$\ell_{\text{total}} = -\ell^{K+1} + P_{\lambda},\tag{12}$$

that sums the negative log-likelihood and a penalization term.

The first term  $-\ell^{K+1}$  is a revising of the negative log-likelihood  $-\ell$  defined in Equation (5) considering a supplementary interval (K+1) in accordance with the constraint described in Equation (11).

The second term  $P_{\lambda}$  is a penalization term related to  $\mathcal{E}^{K+1}$  allowing to smooth the effect of a non-optimal number of time intervals (ie, a non-optimal value for K). The choice of K has a significant impact on the performance: the number of nodes grows with K, which might cause over-fitting (for large value of K) or under-fitting. Lee et al $^{35(\text{Section4})}$  suggest to choose a large K but prevents over-fitting by using L1 regularization over weights in the output layer. Kvamme and Borgan $^{21(\text{Section4.1})}$  suggest to fix a small value for K in order to reduce the size of the output layers as much as possible. However, the automatic selection of K can be fixed by applying a temporal smoothing technique. Following Möst $^{36}$  and Tibshirani et al, $^{37}$  we apply a temporal smoothness constraint by penalizing, in the weight matrices (resp. the bias vectors) of the output layers, the first order differences of the weights (resp. the bias) associated with two adjacent time intervals. Let's consider  $W = (W^1, W^{12})^T$ ,  $B = (B^1, B^{12})^T$  the weight and bias parameters associated with the two output layers, with  $W^1 = (W^{01}, W^{02})^T \in \mathbb{R}^{(l^0+l^{02})\times 2(K+1)}$ ,  $W^{12} \in \mathbb{R}^{(l^{12}\times (K+1))}$  and  $B^1 = (B^{01}, B^{02})^T \in \mathbb{R}^{2(K+1)}$ ,  $B^{12} \in \mathbb{R}^{K+1}$ . For  $k = 1, \ldots, K$ , we compute

$$\Delta_{w_{j,k}^{ql}} = w_{j,k+1}^{ql} - w_{j,k}^{ql}, \ \Delta_{b_{k}^{ql}} = b_{k+1}^{ql} - b_{k}^{ql},$$

the weight and bias differences associated with the transition  $q \to l$ , neuron j and adjacent time intervals  $v_k$ ,  $v_{k+1}$ . Then the penalty term of our loss function in Equation (12) has the form

$$P_{\lambda}(B, W) = \sum_{ql} \left( \lambda_w^{ql} \sum_{j=1}^{l^{ql}} \sum_{k=1}^{K} \left| \Delta_{w_{j,k}^{ql}} \right| + \lambda_b^{ql} \sum_{k=1}^{K} \left| \Delta_{b_k^{ql}} \right| \right),$$

where  $\lambda_w^{ql}$  and  $\lambda_b^{ql}$  are transition-specific positive constants determining the amount of smoothing to be applied for each transition. For  $\lambda_w^{ql} \to +\infty$  (respectively  $\lambda_b^{ql} \to +\infty$ ), all differences will be set to zero resulting in constant weights (respectively constant bias). This penalization term allows to minimize the risk of over-fitting, for the three transitions independently, for larger values of K.

# 4.3 | Selection of the interval cutpoints $a_k$ 's

Under the piecewise constant assumption (see Section 3.3), the definition of the density probabilities requires time to be on the form  $0 = a_0 < a_1 < \cdots < a_K = \tau$ . Hence, we need to perform a division of the time scale and the selection of the interval cutpoints  $a_k$ 's (k = 1, ..., K). On the one hand, we would like to select sufficiently wide interval cutpoints to retain enough information in each interval (ie, keeping a sufficient number of observed transitions within each interval). On the other hand, we would like to select sufficiently narrow cutpoints to ensure that significant temporal changes in the density probabilities can be identified.

To select the cutpoints, the most obvious way would be to choose K equidistant cutpoints in  $[0, \tau]$  (ie, uniform intervals each of length  $K/\tau$ ). An alternative is to select the cutpoints based on the distribution of the transition times. This approach inducts different cutpoints  $a_k^{*,1}$ ,  $a_k^{*,2}$  for the two output layers of IDNetwork. In that case, for the first output layer related to the transitions  $0 \to l$  (l = 1, 2), we can select the cutpoints by estimating the K quantiles  $q_k^{*,1}$  ( $k = 1, \ldots, K$ ) of the marginal distribution of the sojourn duration in state 0 (ie,  $1 - F_0^*(.)$ ). In the same way, for the second output layer related to the transition  $1 \to 2$ , we can select the cutpoints based on the K quantiles  $q_k^{*,2}$  ( $k = 1, \ldots, K$ ) of the distribution of the sojourn duration in state 1 among patients at risk (ie,  $1 - F_{12}^*(.)$ ). The quantiles  $q_k^{*,1}$ ,  $q_k^{*,2}$  (for  $k = 1, \ldots, K$ ) verify:

$$\begin{split} 1 &= 1 - F_0^*(0) = q_0^{*,1} > q_1^{*,1} > \dots > q_K^{*,1} = 1 - F_0^*(\tau), \\ 1 &= 1 - F_{12}^*(0) = q_0^{*,2} > q_1^{*,2} > \dots > q_K^{*,2} = 1 - F_{12}^*(\tau), \end{split}$$

and can be found by estimating the duration distributions  $F_0^*$  and  $F_{12}^*$  via the nonparametric Aalen-Johansen<sup>38</sup> estimator. This approach ensures that each time interval has the same decrease in the duration distribution estimates, such that

$$\begin{aligned} q_{k+1}^{*,1} - q_k^{*,1} &= \left(1 - F_0^*(\tau)\right) / K, \\ q_{k+1}^{*,2} - q_k^{*,2} &= \left(1 - F_{12}^*(\tau)\right) / K. \end{aligned}$$

Finally, by denoting  $\hat{F}_0^{\text{AJ}}$ ,  $\hat{F}_{12}^{\text{AJ}}$  the Aalen-Johansen estimators of  $F_0^*$ ,  $F_{12}^*$ , the interval cutpoints are found by solving

$$1 - \hat{F}_0^{AJ}(a_k^1) = q_k^1,$$
  
$$1 - \hat{F}_{12}^{AJ}(a_k^2) = q_k^2.$$

See the work of Kvamme and Borgan<sup>21</sup> for a similar approach in the case of a unique event.

# 5 | PREDICTION TASK AND BENCHMARK

# 5.1 | Prediction of the individual CIFs

In this subsection, we define the predictions of interest according to the time scales defined below. From the output of our network (ie, the step functions  $\hat{f}_{0l}(\cdot \mid X)$  (l=1,2),  $\hat{f}_{12}(\cdot \mid X)$ ), we can derive the estimation of the CIFs.

For a new patient j with the baseline covariates  $X_j$ , we note the estimated CIFs, derived from Equations (9) and (10), as  $\hat{F}_{0l}(.|X_j)$  (l=1,2),  $\hat{F}_{12}(.|X_j)$ , such that:

$$\hat{F}_{0l}(t|X_j) = \sum_{k=1}^{k(t)-1} \hat{f}_{0l}(v_k|X_j)|v_k| + (t - a_{k(t)-1})\hat{f}_{0l}(v_{k(t)}|X_j) \text{ for } l = 1, 2,$$

$$\hat{F}_{12}(d|X_j) = \sum_{k=1}^{k(d)-1} \hat{f}_{12}(v_k|X_j)|v_k| + (d - a_{k(d)-1})\hat{f}_{12}(v_{k(d)}|X_j).$$

We will use the estimated CIFs to assess the predictive performance of IDNetwork.

## 5.2 | Predictive evaluation criteria

In event history analysis, commonly used performance measures are the time-dependent AUC (for discrimination) and the time-dependent Brier score (BS) (for calibration). On the basis of the transition-specific time properties defined in Section 3.1, we adapt the definitions of the time-dependent AUC<sup>39</sup> and the time-dependent Brier score.<sup>40</sup>

To evaluate predictive performances related to transitions  $0 \to 1$  and  $0 \to 2$ , all the patients are considered. For the transition  $1 \to 2$ , predictions of interest are formulated conditionally to be in state 1 at time  $T_0$  and from the duration  $t - T_0$ . Hence only the patients at risk for experiencing the transition are considered (ie, only the patients who have already experienced a transition  $0 \to 1$ ).

For two patients i and j, the transition-specific time-dependent AUC measures the probability that a patient i who experienced the transition ql before time t has greater probability of occurrence of the transition than a patient j who has survived to the transition. For, on the one hand the transitions  $0 \to 1$ ,  $0 \to 2$ , and on the other hand the transition  $1 \to 2$ , we define

$$\begin{split} & \text{AUC}^{0l}(t) = \mathbb{P}(F_{0l}^*\left(t \mid X_i\right) > F_{0l}^*\left(t \mid X_j\right) \mid T_0^i \leq t, \ T_0^j > t, \ D_0^i = l), \text{for } l = 1, 2, \\ & \text{AUC}^{12}(d) = \mathbb{P}(F_{12}^*\left(d \mid X_i\right) > F_{12}^*\left(d \mid X_j\right) \mid T_2^i - T_0^i \leq d, \ T_2^j - T_0^j > d, D_0^i = 1, \ D_0^j = 1). \end{split}$$

Commonly, the AUC is defined as the integration of the ROC curve opposing specificity (Sp) and sensitivity (Se).

The transition-specific time-dependent Brier score measures the mean difference between the predicted probability of occurrence of the transition at time *t* and the observed status of the transition, such that:

$$BS^{0l}(t) = \frac{1}{n} \sum_{i=1}^{n} \left[ \mathbb{1} \left\{ T_0^i > t \right\} - F_{0l}^*(t \mid X_i) \right]^2 \text{ for } l = 1, 2,$$

$$\mathrm{BS}^{12}(d) = \frac{1}{n_{12}} \sum_{i: D_0^i = 1} \left[ \mathbb{1} \left\{ T_2^i - T_0^i > d \right\} - F_{12}^*(d \mid X_i) \right]^2,$$

where  $n_{12}$  is the number of patients at risk for transition  $1 \rightarrow 2$ .

Several estimators have been developed taking into account the loss of information due to censoring by using the inverse probability of censoring weight (IPCW) methods. To evaluate predictive performances related to transitions  $0 \rightarrow 1$  and  $0 \rightarrow 2$ , we use classical IPCW estimators already developed in the literature.<sup>39,40</sup> For the transition  $1 \rightarrow 2$ , we rewrite the classical estimators under the semi-Markovian assumption considering only the patients at risk for experiencing the transition. Exact expressions of the estimators that we have implemented are given in the Supplementary materials (Section 4).

The time-dependent AUC and the time-dependent BS can be extended to the interval  $]0, \tau]$  by computing, respectively, the integrated AUC (iAUC) and the integrated Brier score (iBS) as follows:

$$iAUC^{ql} = \frac{1}{\tau} \int_0^{\tau} AUC^{ql}(t)dt, iBS^{ql} = \frac{1}{\tau} \int_0^{\tau} BS^{ql}(t)dt.$$

# 5.3 | Softwares and benchmark

Predictive performances of IDNetwork in predicting the CIFs are compared in terms of discrimination (with the iAUC) and calibration (with the iBS) with two state-of-the-art statistical methods: the multi-state Cox P.H model (msCox), that is defined in Section 2.2, from the R library mstate<sup>9</sup> and a spline-based version of the Cox multi-state model (msSplineCox) from the R library flexsurv.<sup>25</sup> We also compare IDNetwork with a simplified linear version of IDNetwork (LinearID-Network), see details in the Supplementary materials (Section 5). We implement IDNetwork and LinearIDnetwork in Python within a Tensorflow environment.\*

# 5.4 | Validation

We perform two sets of experiments on (1) simulated data sets and on (2) three real clinical data sets. For the two sets of experiments, we score predictive performances of the methods through internal validation.<sup>41</sup> We employ Monte Carlo Simulations (MCS) to validate experiments on simulations by generating M data sets. We employ Monte Carlo Cross Validation (MCCV) to validate model performance and estimate model variance on the real data sets by randomly splitting M times each data set. For the two sets of experiments, we set M = 20 due to a high computational time. We validate performance and estimate variance of IDNetwork as follows:

- 1. For each iteration m (m = 1, ..., M) (either the data set m for simulations, or the data set from split m for the real data sets):
  - We split the data set  $\mathcal{D}_m$  into  $\mathcal{D}_m^{\text{train}}/\mathcal{D}_m^{\text{test}}/\mathcal{D}_m^{\text{validation}}$  (70% for training, 10% for early stopping and hyper-parameters tuning, 20% for validation).
  - IDNetwork hyper-parameters are tuned by performing B = 60 random searches on  $\mathcal{D}_m^{\text{train}}$ . Each random set of hyper-parameters is evaluated on the set  $\mathcal{D}_m^{\text{test}}$ .
  - We choose the set of hyper-parameters maximizing the iAUCs (averaged across the three transitions) on the set  $\mathcal{D}_m^{\text{test}}$ .
  - With the optimal set of hyper-parameters, we estimate model parameters on  $\mathcal{D}_m^{\text{train}}$ .

- We compute model performance on the external validation set  $\mathcal{D}_m^{ ext{val}}$ .
- 2. We estimate model performance by computing the median ( $\pm$  standard deviation (SD)) iAUC (higher the better) and iBS (lower the better) on the validation sets  $\mathcal{D}_m^{\text{val}}$  (m = 1, ..., M).
- 3. We estimate predictive performance of the other methods in the same way (excluding the hyper-parameters tuning).
- 4. We statistically compare performances of IDNetwork over the other methods using a bilateral Wilcoxon<sup>42</sup> signed rank test. In the results,  $\cdot$  indicates a *P*-value less than 0.1, † less than 0.05, ‡ less than 0.01, \* less than 0.001.
- 5. In the results, bold values indicate the best model performance.

Experimental details on IDNetwork's hyper-parameters tuning are given in the Supplementary materials (Section 6). The pseudocode for the validation process of IDNetwork is given in the Supplementary materials (Section 7). On real clinical data sets, an external validation<sup>41</sup> of the predictive performance of the models can be conducted by computing for each iteration m (m = 1, ..., M) the criteria on an independent data set as well.

#### 6 | EXPERIMENTS ON SIMULATED DATA SETS

Through simulations, we aim to first illustrate the effects of IDNetwork parameterization on the predictive performance and to secondly compare IDNetwork performance with other methods cited above.

# 6.1 | Data simulation

For the first set of experiments, we conduct Monte Carlo simulations by generating M=20 data sets with the same parameters. We generate continuous-time illness-death data and we fix the horizon-time window at  $\tau=100$ . We generate three sets of data sets:  $\mathcal{D}_{\text{nonlin}}^{2000}$ ,  $\mathcal{D}_{\text{nonlin}}^{5000}$ ,  $\mathcal{D}_{\text{nonlin}}^{20\ 000}$  varying, respectively, the training sample size in  $n_{\text{tr}}=2000$ , 5000, 20000 (80% the data sets) and the validation sample size in  $n_{\text{val}}=500$ , 1250, 5000 (20% the data sets). For each observation i ( $1 \le i \le n$ ) we generate four 2-dimensional baseline variables, i each drawn from a multivariate Gaussian distribution with mean 0 and a matrix of variance covariance  $\Sigma_p$ :

$$X_i = \left(X_i^{(1)}, X_i^{(2)}, X_i^{(3)}, X_i^{(4)}\right)^T \text{ with } X_i^{(p)} \in \mathbb{R}^2 \sim \mathcal{N}\left(0, \Sigma_p\right), \ 1 \leq p \leq 4,$$

where the entries of the matrix  $\Sigma_p^{1/2}$  are simulated from i.i.d. uniform variables on [0, 1].

We aim to generate the processes  $(T_0, D_0)$  and  $T_2$ , such that the illness-death times  $T_{ql}$ , for  $(q, l) \in \{(01), (02), (12)\}$ , are simulated through Cox transition-specific hazard functions:<sup>43</sup>

$$T_{ql}^i \sim \alpha_{ql}^*(t|X_i) = \alpha_{ql}^{*,0}(t) \exp\left(g_{ql}\left(X_i, \beta_{ql}^*\right)\right)$$

where  $g_{ql}(.)$  is a transition-specific risk function,  $\beta_{ql}^* = \left(\beta_{ql}^{*,(1)},\beta_{ql}^{*,(2)},\beta_{ql}^{*,(3)},\beta_{ql}^{*,(4)}\right)^T$  with  $\beta_{ql}^{*,(p)} \in \mathbb{R}^2$  for  $1 \le p \le 4$  are fixed effect coefficients, and  $\alpha_{ql}^{*,0}(.)$  is the baseline hazard function. We generate the three baseline hazard functions as follows:  $\alpha_{ql}^{*,0}(.) \sim \text{Weibull}$  (scale = 0.01, shape = 1.2).

We set the transition-specific risk functions to be nonlinear using quadratic functions, in the spirit of Lee et al, 18 as

$$g_{01}\left(X_{i},\beta_{01}^{*}\right)=\left(X_{i}^{(1)}\beta_{01}^{*,(1)}+X_{i}^{(2)}\beta_{01}^{*,(2)}\right)^{2},\ g_{02}\left(X_{i},\beta_{02}^{*}\right)=\left(X_{i}^{(2)}\beta_{02}^{*,(2)}+X_{i}^{(3)}\beta_{02}^{*,(3)}\right)^{2},\ g_{12}\left(X_{i},\beta_{12}^{*}\right)=\left(X_{i}^{(3)}\beta_{12}^{*,(3)}+X_{i}^{(4)}\beta_{12}^{*,(4)}\right)^{2}.\tag{13}$$

We fix arbitrary values for the fixed effects coefficients. Hence, in this simulation scheme, the Cox's linear assumption does not hold anymore.

From the simulated  $T_{ql}$ , the simulation of the processes  $(T_0, D_0)$  and  $T_2$  has to respect constraints of the model (ie,  $T_{01}$  and  $T_{02}$  are competing,  $T_{01}$  and  $T_{12}$  are recurrent) in order to generate an identifiable Cox model . We fix r = 30% such that 30% of patients from state 0 are censored, and 30% of patients at risk for transition  $1 \rightarrow 2$  are censored from state 1, see Table 1. We refer the reader to the Supplementary materials (Section 8) for more details on the simulation.

TABLE 1 Descriptive statistics on the number of (No.) observations in the simulated data sets

No. observations (%)						
Data set	$0 \rightarrow 1$	$0 \rightarrow 2$	$0 \rightarrow cens.$	$1 \rightarrow 2$	$1 \rightarrow cens.$	Total
$\mathcal{D}_{ m nonlin.}^{2000}$	720 (36%)	680 (34%)	600 (30%)	504 (70% <sup>a</sup> )	216 (30% <sup>a</sup> )	2000
$\mathcal{D}_{ m nonlin.}^{ m 5000}$	1817 (36%)	1683 (34%)	1500 (30%)	1272 (70% <sup>a</sup> )	545 (30% <sup>a</sup> )	5000
$\mathcal{D}_{ m nonlin.}^{ m 20~000}$	7234 (36%)	6766 (34%)	6000 (30%)	5064 (70% <sup>a</sup> )	2170 (30% <sup>a</sup> )	20 000

<sup>&</sup>lt;sup>a</sup>Among patients at risk.

# 6.2 | Simulation study

# 6.2.1 Understanding the effect of K and n

To get a better understanding of the methodologies discussed in Sections 4.2 and 4.3, we perform a simulation study where we vary the size n of the data sets, the number K of time intervals used for the piecewise approximations and the selection methods of the interval cutpoints. Gensheimer and Narasimhan<sup>17</sup> performed a similar study in a discrete-time approach by varying the value K with the conclusion that there is no difference in the predictive performance. Kvamme and Borgan<sup>21</sup> performed a similar study as well in both discrete-time and piecewise approaches with the conclusion that there is no difference in predictive performances between the methods to select the interval cutpoints. However, they concluded that smaller values for K are better for smaller values of n.

For evaluation, we use the internal validation process described in Section 5.3. We use the iAUC and iBS measures, in addition to the transition-specific integrated mean absolute error (iMAE<sub>ql</sub>) between the estimated CIFs  $\hat{F}_{ql}(.|X_i)$  and the true CIFs  $F_{ql}^*(.|X_i)$ :

$$\mathrm{iMAE}_{ql} = \frac{1}{\tau} \int_0^\tau \frac{1}{n_{\mathrm{val}}} \sum_{i=1}^{n_{\mathrm{val}}} \left| \hat{F}_{ql}(t|X_i) - F_{ql}^*(t|X_i) \right|, \ \ \mathrm{for} \ \ (q,l) \in \{(0,1),(0,2),(1,2)\},$$

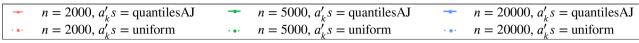
where  $n_{\rm val}$  is the number of subjects in the validation sets.

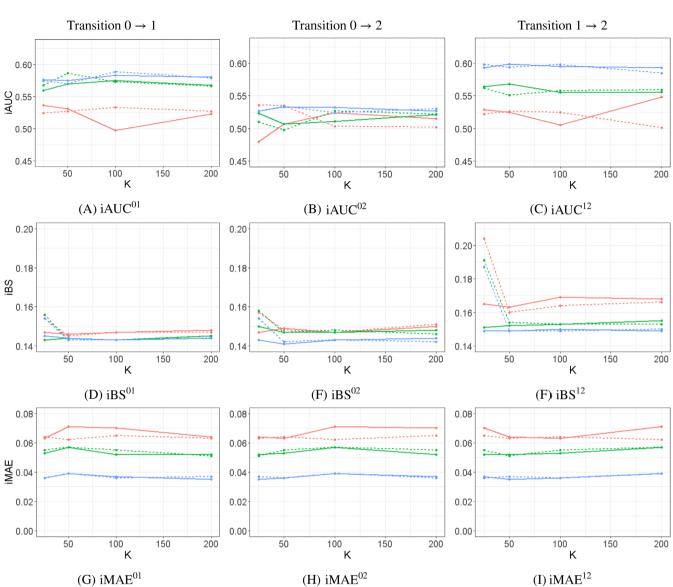
For the discretization of the time scale in the three data sets, we vary the number K of time intervals in K = 25, 50, 100, 200. We applied two methods to select the interval cutpoints with either equidistant cutpoints (uniform) or cutpoints obtained with the Aalen-Johansen quantiles (quantiles AJ).

In Figure 3, we plot the transition-specific validated criteria (iAUCs, iBSs, iMAEs) of the cutpoints selection methods versus the values of K and n. We integrate the AUC and BS measures at all the 4 equidistant time points in  $[0,\tau]$  (ie, at times  $t=4,8,12,\ldots,96$ , for computational cost reasons). We integrate the MAE measure at all the 100 discrete time points in  $[0,\tau]$  (ie, at times  $t=1,2,3\ldots,100$ ). We can see that the selection methods of the interval cutpoints give similar performances in terms of iAUC and iMAE. However, in term of iBS, the uniform selection method give slightly poorer performance for a smaller value of K (for K=25) than the quantilesAJ selection method. In terms of iAUC, iBS and iMAE, it is evident to see that a larger value of n increases the performances for all the values of K. However, for transition  $0 \to 2$ , we can see that the performance differences in terms of iAUC are very unstable when varying simultaneously the value of n, the value of K and the methods for selecting the interval cutpoints.

# 6.2.2 Understanding the role of $P_{\lambda}$

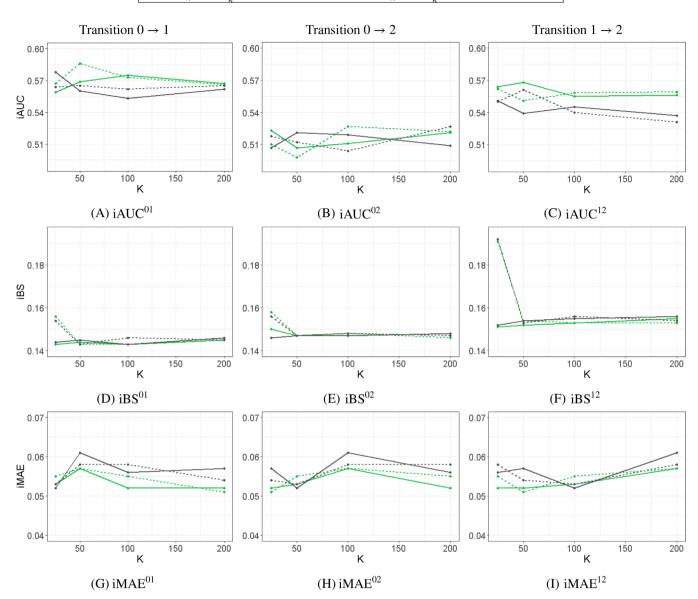
To get a better understanding of the role of the penalization term  $P_{\lambda}$  in the loss function in Equation (12), we perform a second study on set of simulated data sets  $\mathcal{D}_{\text{nonlin}}^{5000}$  where we vary the effect of  $P_{\lambda}$ , the number K of time intervals (K = 25, 50, 100, 200) and the selection methods of the interval cutpoints (uniform, quantilesAJ).





**FIGURE 3** Median iAUCs, iBSs, and iMAEs per transition, for each specified values of *K* and *n* in the simulation study. The full lines represent the method quantilesAJ, while the dotted lines represent the method uniform

In Figure 4, we plot the transition-specific validated criteria (iAUCs, iBSs, iMAEs) versus the values of K. The full lines represent the method quantilesAJ, while the dotted lines represent the method uniform. The green color represents the performance of IDNetwork with  $P_{\lambda} \neq 0$ , while the brown color represents the performance of IDNetwork with  $P_{\lambda} = 0$ . In terms of iAUCs, the penalized model (ie, with  $P_{\lambda} \neq 0$ ) outperforms the unpenalized model (ie, with  $P_{\lambda} = 0$ ) for all the values of K for the transitions  $0 \to 1$  and  $1 \to 2$ . In terms of iBSs, the penalized and the unpenalized model gives slightly better performance for transition  $1 \to 2$ . In terms of iMAEs, the penalized model outperforms the unpenalized model for all the values of K for the three transitions. Subsequently, we can see an interaction between the value of K and the effect of the penalization regarding the iAUCs and the iMAEs. For the smaller value of K (ie, K = 25), the penalized and the unpenalized models are equivalent. For the larger value of K (ie, K = 200), the penalized model outperforms the unpenalized model. Thus, it is evident that when the value of K increases, the performance of the unpenalized model decreases (in particular for the method quantilesAJ) while the performance of the penalized model



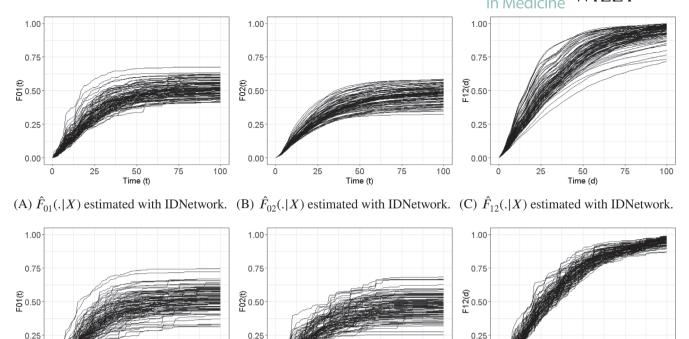
**FIGURE 4** Illustration of the effect of the penalization term  $P_{\lambda}$  in the loss function in Equation (12) on the set of simulated data sets  $\mathcal{D}_{\text{nonlin.}}^{5000}$ : median iAUCs, iBSs, and iMAEs per transition, for each specified values of K. The full lines represent the method quantilesAJ, while the dotted lines represent the method uniform. The red color represents the performance of IDNetwork with  $P_{\lambda} \neq 0$ , while the brown color represents the performance of IDNetwork with  $P_{\lambda} = 0$ 

is stable. Hence, the penalization term  $P_{\lambda}$  allows to mitigate the effect of a too large value of K on the IDNetwork performance.

In addition, the role of the penalization term  $P_{\lambda}$  is to provide smooth estimates of the density probability functions. Indeed, it may be reasonable to assume that large temporal variations in the density functions over successive time intervals should be smoothed. For a too large value of K, it will result in the emergence of high jumps when computing the CIFs. This problem can be fixed by applying a temporal smoothing method as done by  $P_{\lambda}$ . To illustrate our point, we use the simulated data set  $\mathcal{D}_{\text{nonlin.}}^{5000}$  and generate 100 additional observations. We estimate the CIFs of the additional observations with IDNetwork in the case of  $P_{\lambda} = 0$  (unpenalized model) and in the case of  $P_{\lambda} \neq 0$  (penalized model). Results per transition are shown in Figure 5. For the transitions  $0 \to 1$  and  $0 \to 2$ , the penalization term smoothes the estimates when high jumps are estimated between adjacent time intervals. For the transition  $1 \to 2$ , estimates of IDNetwork with

0.00

(D)



**FIGURE 5** Estimated  $\hat{F}_{ql}(.|X)$  in  $[0,\tau]$  on the set of simulated data sets  $\mathcal{D}^{5000}_{\text{nonlin.}}$  for 100 additional simulated patients  $((q,l) \in \{(0,1),(0,2),(1,2)\})$ 

(E)  $\hat{F}_{02}(.|X)$ 

IDNetwork( $P_{\lambda} = 0$ ).

0.00

100

with

75

estimated

 $P_{\lambda} = 0$  result in a loss of the variability that is corrected by setting  $P_{\lambda} \neq 0$ . We refer the reader to the true functions in the Supplementary materials (Section 8.2); it is evident that the penalization term smooth very well the jumps that are nonexistent in the true curves.

0.00

25

(F)  $\hat{F}_{12}(.|X)$ 

IDNetwork( $P_{\lambda} = 0$ ).

100

with

75

estimated

100

with

75

Time (t)

estimated

# 6.3 | Benchmark

 $\hat{F}_{01}(.|X)$ 

IDNetwork( $P_{\lambda} = 0$ ).

We use the set of nonlinear simulated data sets  $\mathcal{D}_{\text{nonlin.}}^{5000}$  to compare the predictive performance of IDNetwork with the state-of-the-art methods (see Section 5.3). To benchmark IDNetwork, we set K = 100 and divide the time scale into intervals of uniform length.

The integrated predictive performances are shown in Table 2. Detailed results per evaluation times are displayed in the Supplementary materials (Section 9.1). In these simulations, the Cox's linear assumption does not hold anymore. Consequently, as expected, IDNetwork significantly outperforms msCox and msSplineCox with a P-value less than 0.001 for the three transitions in terms of iAUC and iBS. (except for the transition  $1 \rightarrow 2$  where msCox and msSplineCox outperform IDNetwork but with no statistical difference). IDNetwork significantly outperforms the linear version LinearIDNetwork with a P-value less than 0.001 as well. Moreover, we note that LinearIDNetwork outperforms msCox and msSplineCox. This illustrates the effect of a deep learning approach as compared with a statistical approach. We also evaluate the predictive performances of IDNetwork on a linear simulated data set. Results are shown in the Supplementary materials (Section 9.2).

# 7 | APPLICATION ON REAL CLINICAL DATA SETS

We conduct experiments on real illness-death data from two clinical trials in colon cancer and one clinical trial in breast cancer.

10970258, 2022, 9, Downloaded from https://onlinelibrary.wiley

on Wiley Online Library for rules of use; OA articles are governed by the applicable Creative Commons License

**TABLE 2** Predictive performance (median  $\pm$  SD) on the validation sets (internal validation with Monte Carlo simulations) for the set of nonlinear simulated data sets  $\mathcal{D}_{\text{nonlin}}^{5000}$ , with K = 100 and a uniform subdivision of the time scale

		Transition			
Criteria	Algorithm	0  o 1	0 → 2	1 → 2	Average
iAUC	msCox	$0.508^* \pm 0.04$	$0.489^{\circ} \pm 0.05$	$0.499^* \pm 0.04$	$0.499^* \pm 0.02$
	msSplineCox	$0.509^* \pm 0.05$	$0.488^{\circ} \pm 0.05$	$0.501^* \pm 0.04$	$0.498^* \pm 0.01$
	LinearIDNetwork	$0.517^* \pm 0.05$	$0.499 \pm 0.06$	$0.509^* \pm 0.03$	$0.502^* \pm 0.02$
	IDNetwork	$\textbf{0.573} \pm \textbf{0.03}$	$\textbf{0.527} \pm \textbf{0.03}$	$\textbf{0.558} \pm \textbf{0.04}$	$\textbf{0.545} \pm \textbf{0.01}$
iBS	msCox	$0.229^* \pm 0.01$	$0.241^* \pm 0.01$	$\textbf{0.152} \pm \textbf{0.01}$	$0.206^* \pm 0.01$
	msSplineCox	$0.238^* \pm 0.01$	$0.249^* \pm 0.01$	$\textbf{0.152} \pm \textbf{0.01}$	$0.211^* \pm 0.01$
	LinearIDNetwork	$0.161^* \pm 0.01$	$0.160^* \pm 0.01$	$0.158 \pm 0.01$	$0.160^* \pm 0.01$
	IDNetwork	$\textbf{0.143} \pm \textbf{0.01}$	$\textbf{0.148} \pm \textbf{0.01}$	$0.153 \pm 0.01$	$\textbf{0.148} \pm \textbf{0.01}$

TABLE 3 Descriptive statistics on the number of (No.) observations in the colon cancer data sets

	No. observation	No. observations (%)						
Data set	$0 \rightarrow 1$	$0 \rightarrow 2$	$0 \rightarrow cens.$	$1 \rightarrow 2$	$1 \rightarrow cens.$	Total		
NCT00079274	623 (29%)	81 (4%)	1427 (67%)	276 (44% <sup>a</sup> )	347 (56% <sup>a</sup> )	2121		
NCT00275210	279 (25%)	14 (1%)	829 (74%)	132 (47% <sup>a</sup> )	147 (53% <sup>a</sup> )	1122		

<sup>&</sup>lt;sup>a</sup>Among patients at risk.

TABLE 4 Descriptive statistics on the transition times (in months) in the colon cancer data sets

	Transi	tion							
	$0 \rightarrow 1$			<b>0</b> → 2			$1 \rightarrow 2$		
Data set	Q1	Median	Q3	Q1	Median	Q3	Q1	Median	Q3
NCT00079274	11	15	23	6	14	32	17	25	33
NCT00275210	10	14	22	4	7	23	16	22	30

#### 7.1 Data on colon cancer

# 7.1.1 Description of the data sets

We use two data sets from Phase III clinical trials evaluating endpoints relapse-free survival (RFS) and overall-survival (OS) in non-metastatic colon cancer. (1) The study NCT00079274<sup>†</sup> contains 2121 observed patients followed for 60 months (5 years) for RFS and for 96 months (8 years) for OS.<sup>44</sup> It presents 67% of censoring from state 0 and 56% of censoring from state 1 among patients at risk. (2) The study NCT00275210<sup>‡</sup> contains 1122 patients followed over 60 months for RFS and OS.<sup>45</sup> The data set presents 74% of censoring from state 0 and 53% from state 1 among patients at risk. Descriptive statistics of the colon data sets are shown in Tables 3 and 4.

The preprocessing of this two data sets requires a preliminary evaluation of the compatible features (same covariates and same distributions of the covariates) and to adjust the length of follow-up between both data sets. Thus, we finally restrict our attention to 9 baseline clinical covariates (8 categorical and 1 numerical) including the following features: BMI, sex, race, age, tumor histological type, number of positive lymph nodes, cancer stage, ECOG performance status, presence of bowel obstruction/perforation. In the study NCT00079274, outcome RFS has been right-censored at 5 years and outcome OS at 8 years. Whereas in the study NCT00275210, both outcomes have been right-censored at 5 years. We adjust the length of follow-up of both studies choosing a value for  $\tau$  compatible with both. Then for the experimentation,

**TABLE 5** Predictive performance (median  $\pm$  SD) for the data sets NCT00079274, NCT00275210 on colon cancer on (1) the validation sets (internal validation), (2) the external NCT00275210 test set (external validation), with a uniform subdivision of the time scale in K = 48 (months)

			Transition			
Evaluation	Criteria	Algorithm	$0 \rightarrow 1$	$0 \rightarrow 2$	$1 \rightarrow 2$	Average
(1) Internal	iAUC	msCox	$\boldsymbol{0.677 \pm 0.03}$	$0.658 \pm 0.08$	$0.682 \pm 0.04$	$\boldsymbol{0.663 \pm 0.03}$
		msSplineCox	$0.672 \pm 0.03$	$0.600 \pm 0.09$	$0.692^\dagger \pm 0.04$	$0.659 \pm 0.03$
		LinearIDNetwork	$0.666 \pm 0.03$	$0.617 \pm 0.09$	$0.665^{\circ} \pm 0.05$	$0.637 \pm 0.04$
		IDNetwork	$0.669 \pm 0.03$	$\boldsymbol{0.660 \pm 0.09}$	$0.662 \pm 0.06$	$0.651 \pm 0.03$
	iBS	msCox	$\textbf{0.153} \pm \textbf{0.01}$	$0.032^{\ddagger} \pm 0.00$	$0.201^{\circ} \pm 0.02$	$0.130 \pm 0.01$
		msSplineCox	$0.156 \pm 0.01$	$0.032^{\circ} \pm 0.01$	$0.192^\dagger \pm 0.03$	$\textbf{0.126} \pm \textbf{0.01}$
		LinearIDNetwork	$0.154 \pm 0.01$	$0.028 \pm 0.00$	$0.207 \pm 0.03$	$0.130 \pm 0.01$
		IDNetwork	$\textbf{0.153} \pm \textbf{0.01}$	$\boldsymbol{0.027 \pm 0.00}$	$0.212 \pm 0.03$	$0.129\pm0.01$
(2) External	iAUC	msCox	$0.669^{\dagger} \pm 0.00$	$0.601^* \pm 0.02$	$0.559 \pm 0.02$	$0.610^* \pm 0.01$
		msSplineCox	$0.668^{\dagger} \pm 0.00$	$0.598^* \pm 0.03$	$\boldsymbol{0.566 \pm 0.01}$	$0.613^* \pm 0.00$
		LinearIDNetwork	$0.670 \pm 0.01$	$0.694 \pm 0.03$	$0.533^\dagger \pm 0.02$	$0.642^\dagger \pm 0.01$
		IDNetwork	$\boldsymbol{0.673 \pm 0.01}$	$\textbf{0.713} \pm \textbf{0.05}$	$0.562 \pm 0.03$	$\boldsymbol{0.651 \pm 0.02}$
	iBS	msCox	$0.157^\dagger \pm 0.00$	$0.013^* \pm 0.00$	$0.184^* \pm 0.01$	$0.118^* \pm 0.00$
		msSplineCox	$0.159^\dagger \pm 0.00$	$0.013^* \pm 0.00$	$0.176^* \pm 0.01$	$0.116^* \pm 0.00$
		LinearIDNetwork	$\textbf{0.154} \pm \textbf{0.00}$	$\textbf{0.011} \pm \textbf{0.00}$	$\textbf{0.146} \pm \textbf{0.01}$	$\boldsymbol{0.103 \pm 0.00}$
		IDNetwork	$\textbf{0.154} \pm \textbf{0.01}$	$\textbf{0.011} \pm \textbf{0.00}$	$0.150 \pm 0.01$	$0.105 \pm 0.00$

Note: We integrate the AUC and BS measures at all the 30 equidistant time points in [60,  $\tau$ ] (ie, at every month from 2 months).

the study NCT00079274 will be used for internal validation (training and validation) and the study NCT00275210 for external validation.<sup>41</sup>

For the two data sets, missing values were imputed by the median value for numerical features and by the mode for categorical features. We apply one-hot encoding on categorical features and standardize numerical features with the Z-score. For each data set, even though IDNetwork is a continuous-time method, we still need to subdivide the time axis into time intervals. We fix a uniform length for the time intervals to 1 month such that the time interval for month j,  $v_j = [j-1,j)$ , includes all the events that occurred on the daily time interval  $[(j-1) \times 30.5, j \times 30.5)$ . We set K = 48 (ie,  $\tau = 48$  months = 4 years) and subdivide the time axis into monthly intervals between 0 and 48 months and set event times after 48 months in a last interval  $v_{49}$ .

# 7.1.2 | Benchmark

We conduct (1) internal and (2) external validation. The integrated predictive performances are shown in Table 5. Results per evaluation times are detailed in the Supplementary materials (Section 10). On the validation splits, neither method outperforms the other for the three transitions independently. In addition, all the models are equivalent (ie, no significant differences) on average. However, IDNetwork shows significant better performance on the external NCT00275210 validation study for all the transitions (excluding in terms of iAUC for transition  $1 \rightarrow 2$  where msSplineCox outperforms IDNetwork but with no statistical significance, in terms if iBS for transition  $1 \rightarrow 2$  where LinearIDNetwork display better performance but with no statistical significance). On average, IDNetwork displays significant better iAUC and iBS than msCox and msSpline Cox with a P-value less than 0.001, a significant better iAUC than LinearIDNetwork with a P-value less than 0.05.

Mahidol University Library & Knowledge Center, Wiley Online Library on [20/03/2025]. See the Terms

TABLE 6 Descriptive statistics on the number of (No.) observations in the METABRIC data set

	No. observations (%)							
Data set	$0 \rightarrow 1$	0  o 2	$0 \rightarrow cens.$	$1 \rightarrow 2$	$1 \rightarrow cens.$	Total		
METABRIC	677 (36%)	509 (27%)	717 (38%)	593 (88% <sup>a</sup> )	84 (12% <sup>a</sup> )	1903		

<sup>&</sup>lt;sup>a</sup>Among patients at risk.

TABLE 7 Descriptive statistics on the transition times (in months) in the METABRIC data set

	Transi	tion							
0 → 1			<b>0</b> → 2	_ 0 → 2			1 → 2		
Data set	Q1	Median	Q3	Q1	Median	Q3	Q1	Median	Q3
METABRIC	20	39	81	65	113	121	36	61	110

# 7.2 | Data on breast cancer

# 7.2.1 | Description of the data set

We use the data from the Molecular Taxonomy of Breast Cancer International Consortium (METABRIC) cohort. <sup>46</sup> For the METABRIC§ data set, we include 1903 patients followed for 360 months (30 years) for RFS and OS, with 38% of censoring from state 0 and 13% of censoring from state 1 among patients at risk. Descriptive statistics of the data set are shown in Tables 6 and 7.

This data set contains clinical, histo-pathological, gene copy number and gene expression features used to determine breast cancer subgroups. Based on the literature, we select 17 baseline clinical, histo-pathological and gene copy number features (12 categorical, 5 numerical). It includes the following features: age, inference on the menopausal status, Nottingham Prognostic Index (NPI), immunohistochemical oestrogen-receptor (ER) status, number of positive lymph nodes, cancer grade, tumor size, tumor histological type, cellularity, Her2 copy number by SNP6, Her2 expression, ER Expression, progesterone (PR) expression, type of breast surgery, cancer molecular subtype (pam50 subgroup, integrative cluster), chemotherapy regimen, hormonal regimen, radiotherapy regimen.

The data set contains around 24 000 gene expression features. Before applying methods on this data set, we conduct a preprocessing step by using statistical methods, in order to select a subset of relevant genes to be included in IDNetwork. Several approaches have been reported to integrate gene expression data into survival models. These approaches are based either on dimension reduction, on genes or metagenes selection<sup>47</sup> or, more recently, on the use of a large number of gene expression values (>1000) with the development of deep learning methods.  $^{48,49}$  We choose a selection approach in order to extract a ranked list of cancer-related genes based on their P-values from independent transition-specific Cox P.H. models. The genes are selected using an independent data set in order to control for selection bias; we use the Breast Invasive Carcinoma TCGA PanCancer data set. We describe our selection process in two steps as follows:

## 1. For each gene expression feature:

- We fit independent transition-specific Cox P.H. models including the clinical covariates and each gene expression feature
- For each of the three transitions, we compute the *P*-value from a Wald test on the estimated coefficient related to the specific gene.

# 2. For each transition:

- We adjust the *P*-values with a Benjamini-Hochberg multi-test correction for each transition independently. See the Supplementary materials (Section 11) for details on the *P*-values per transition.
- We rank the *P*-values for each transition.

TABLE 8 Number of (No.) genes selected per transition on the Breast Invasive Carcinoma TCGA PanCancer data set for each threshold

	No. genes sele	No. genes selected per transition				
Threshold ( $\alpha$ )	$0 \rightarrow 1$	0  o 2	$1 \rightarrow 2$	Total <sup>a</sup>		
0.001	8	0	1	9		
0.005	17	1	1	19		
0.01	26	2	1	29		
0.05	76	6	3	84		
0.1	110	36	11	156		

<sup>&</sup>lt;sup>a</sup>Total number of unique genes selected.

• We use different thresholds, noted α, of the *P*-values for gene selection. We present the number of gene selected for each threshold in Table 8. For the transition 0 → 1, several genes are selected (up to 110 genes for the threshold 0.1). While for the transitions 0 → 2 and 1 → 2, very few genes are selected. In order to assess the biological meaning of the selected genes, we perform a gene set enrichment analysis<sup>50</sup> for each transition-specific list of genes against the Hallmark<sup>51</sup> collection of the Molecular Signatures Database (MSigDB). For the transitions 0 → 2 and 1 → 2, there is no significant enriched terms. For the transition 0 → 1, there are 9 significant enriched terms. We refer the reader to the Supplementary materials (Section 12) for detailed results on the analysis. In particular, the most significant enriched term "Epithelial Mesenchymal Transition" is related to a core biological component of the metastatic process in breast cancer<sup>52</sup> and therefore well related to a primary cancer progression through a state relapse. Hence, the gene selected on the transition 0 → 1 display a coherent functional profile associated to cancer relapse. The addition of the selected genes in an illness-death model is thus expected to improve its performance for the transition 0 → 1 more than the others.

For all the features, missing values were imputed by the median value for numerical features and by the mode for categorical features. We apply one-hot encoding on categorical features and standardize numerical features with the Z-score. We fix a uniform length for the time intervals to 1 month such that the time interval for month j,  $v_j = [j-1,j)$ , includes all the events that occurred on the daily time interval  $[(j-1) \times 30.5, j \times 30.5)$ . We fix K=120 (ie,  $\tau=120$  months = 10 years) and subdivide the time axis into monthly intervals between 0 and 120 months and set event times after 120 months in a last interval  $v_{121}$ .

# 7.2.2 | Benchmark

With the Metabric data set we aim to compare the scalability performance of IDNetwork with other methods when including a large number of features and varying the amount of features. The integrated predictive performances are shown in Tables 9 and 10. Detailed results per evaluation times are displayed in the Supplementary materials (Section 13). We compare different models with different lists of genes varying the threshold  $\alpha$  of the P-values. For each of the four compared algorithms, the models show different performances in terms of iAUC. The algorithms msCox and msSplineCox display similar iAUC for transition  $0 \to 1$ , better iAUC for transition  $0 \to 2$  and poorer iAUC for transition  $1 \to 2$  when the value of  $\alpha$  increases. The algorithms IDNetwork and LinearIDNetwork improve their iAUC for transition  $0 \to 1$  and  $1 \to 2$ when the value of  $\alpha$  increases. As expected, IDNetwork and LinearIDNetwork display similar iAUC for the transitions 0  $\rightarrow$  2 and 1  $\rightarrow$  2 when the value of  $\alpha$  increases, as most of the genes included are informative on transition 0  $\rightarrow$  1. Hence, the two deep learning methods have a better iAUC for the transition  $0 \to 1$  thanks to the addition of the selected gene features in the modelization. This result is consistent with the gene set enrichment analysis realized in Section 7.2.1 in which we have highlighted than the list of genes selected for the transition  $0 \to 1$  is particularly related to the metastatic process whereas there are no enriched terms in the genes selected for the transitions  $0 \to 2$  and  $1 \to 2$ . In terms of iBS, we can see a decreasing of the msCox and msSpline performances when the value of  $\alpha$  increases. Whereas, LinearIDNetwork and IDNetwork performances remain stable. Hence it seems that the deep learning methods are less sensitive in terms of iBS to the increase in the number of features. Finally, we chose the model  $M_{0.1}^{BH}$  as the best model to predict transitions

of use; OA articles are governed by the applicable Creative Commons License

**TABLE 9** Integrated AUCs (median  $\pm$  SD) on the validation sets (internal validation) for the METABRIC data set, with a uniform subdivision of the time scale in K = 120 (months)

			Transition					
Criteria	Model	Algorithm	0 → 1	0 → 2	1 → 2	Average		
AUC	$M_0$	msCox	$\boldsymbol{0.686 \pm 0.02}$	$0.702^\dagger \pm 0.05$	$0.725 \pm 0.04$	$0.711^\dagger \pm 0$		
		msSplineCox	$\boldsymbol{0.686 \pm 0.02}$	$0.702^\dagger \pm 0.05$	$0.729 \pm 0.04$	$0.710^{\dagger} \pm 0$		
		LinearIDNetwork	$0.649^{\ddagger} \pm 0.04$	$0.609^* \pm 0.05$	$0.705 \pm 0.04$	$0.657^* \pm 0$		
		IDNetwork	$0.681 \pm 0.03$	$0.671 \pm 0.05$	$\boldsymbol{0.732 \pm 0.05}$	$0.689 \pm 0.$		
	$\mathrm{M}_{0.001}^{\mathrm{BH}}$	msCox	$\boldsymbol{0.697 \pm 0.03}$	$0.740^* \pm 0.03$	$0.728 \pm 0.02$	$0.720^* \pm 0$		
		msSplineCox	$0.694 \pm 0.03$	$0.751^* \pm 0.04$	$\textbf{0.734} \pm \textbf{0.01}$	$0.726^{\ddagger} \pm 0$		
		LinearIDNetwork	$0.657^{\ddagger} \pm 0.03$	$0.591^* \pm 0.06$	$0.718 \pm 0.03$	$0.651^* \pm 0$		
		IDNetwork	$0.689 \pm 0.03$	$0.659 \pm 0.05$	$0.717 \pm 0.04$	$0.692 \pm 0$		
	$M_{0.005}^{\mathrm{BH}}$	msCox	$0.702^\dagger \pm 0.02$	$0.750^* \pm 0.03$	$0.722 \pm 0.03$	0.724* ±		
		msSplineCox	$0.690 \pm 0.02$	$0.756^* \pm 0.04$	$\boldsymbol{0.723 \pm 0.02}$	$0.720^{\ddagger} \pm 0$		
		LinearIDNetwork	$0.656^\dagger \pm 0.03$	$0.629^{\ddagger} \pm 0.07$	$0.703^{\circ} \pm 0.03$	$0.653^* \pm 0$		
		IDNetwork	$0.684 \pm 0.03$	$0.677 \pm 0.04$	$0.722 \pm 0.03$	$0.692 \pm 0$		
	$\mathbf{M}_{0.01}^{\mathrm{BH}}$	msCox	$0.693 \pm 0.02$	$0.753^* \pm 0.03$	$0.710 \pm 0.04$	0.713* ± 0		
		msSplineCox	$\boldsymbol{0.694 \pm 0.02}$	$0.763^* \pm 0.03$	$\boldsymbol{0.721 \pm 0.04}$	0.724 <sup>‡</sup> ±		
		LinearIDNetwork	$0.665^\dagger \pm 0.03$	$0.615^{\circ} \pm 0.06$	$0.715 \pm 0.03$	$0.654^{\ddagger} \pm 0$		
		IDNetwork	$0.685 \pm 0.03$	$0.649 \pm 0.06$	$0.717 \pm 0.03$	$0.690 \pm 0.690$		
	$\mathrm{M}_{0.05}^{\mathrm{BH}}$	msCox	$0.696 \pm 0.03$	$0.743^* \pm 0.03$	$0.692^{\ddagger} \pm 0.03$	$0.709^* \pm 0$		
		msSplineCox	$\textbf{0.716} \pm \textbf{0.03}$	$0.722^{\ddagger} \pm 0.03$	$\boldsymbol{0.729 \pm 0.05}$	0.739° ± 0		
		LinearIDNetwork	$0.665^* \pm 0.03$	$0.627 \pm 0.05$	$0.720 \pm 0.02$	$0.667^{\ddagger} \pm 0$		
		IDNetwork	$0.704 \pm 0.03$	$0.640 \pm 0.05$	$0.722 \pm 0.02$	$0.686 \pm 0$		
	${ m M}_{0.1}^{ m BH}$	msCox	$0.681^* \pm 0.04$	$0.722^* \pm 0.03$	$0.671^* \pm 0.04$	$0.691 \pm 0.691$		
		msSplineCox	$0.686 \pm 0.01$	$0.683 \pm 0.02$	$0.688^{\circ} \pm 0.02$	$0.686 \pm 0.0$		
		LinearIDNetwork	$0.670^* \pm 0.03$	$0.634 \pm 0.05$	$0.719 \pm 0.03$	$0.679^{\ddagger} \pm 0$		
		IDNetwork	$0.709 \pm 0.03$	$0.661 \pm 0.05$	$\boldsymbol{0.732 \pm 0.03}$	$0.697 \pm 0.$		

Note: We integrate the AUC and BS measures at all the 30 equidistant time points in  $[90, \tau]$  (ie, at every month from 3 months). The model  $M_0$  uses only the clinical features and each model  $M_{\alpha}^{BH}$ , for  $\alpha \in [0.001, 0.005, 0.01, 0.05, 0.1]$ , uses the clinical features and the genes selected with a corrected P-value less than  $\alpha$ .

 $0 \to 1$  and  $1 \to 2$ . With the model  $M_{0.1}^{BH}$ , IDNetwork outperforms significantly the other algorithms with an averaged iAUC of 0.697 and an averaged iBS of 0.124, an iAUC of 0.714 and an iBS of 0.142 for transition  $0 \to 1$ , an iAUC of 0.730 and an iBS of 0.167 for transition  $1 \to 2$ . IDNetwork show poorer performance for transition  $0 \to 2$ . This is, in our opinion, due to the fact that experiencing transition  $0 \to 2$  means that the patient died from another cause than cancer, but the clinical and biological features provided to the model were not related (excluding the age) to non-cancer causes of death.

# 8 | DISCUSSION

We present IDNetwork a novel deep learning method to model an illness-death process and to predict two-stage evolution of a disease based on baseline covariates. To the best of our knowledge, it is the first deep learning architecture developed in the context of multi-state analysis. It outperforms standard methodologies in this context. In clinical practice, IDNetwork may be useful in personalized medicine by providing predictions of the risks of relapse and death. It could help physicians to adapt the therapeutic guidelines for a specific patient.

**TABLE 10** Integrated BSs (median  $\pm$  SD) on the validation sets (internal validation) for the METABRIC data set, with a uniform subdivision of the time scale in K = 120 (months)

			Transition			
Criteria	Model	Algorithm	0 → 1	$0 \rightarrow 2$	$1 \rightarrow 2$	Average
iBS	$M_0$	msCox	$\textbf{0.145} \pm \textbf{0.01}$	$0.071^* \pm 0.01$	$\textbf{0.156} \pm \textbf{0.01}$	$0.124 \pm 0.0$
		msSplineCox	$0.146 \pm 0.01$	$0.073^* \pm 0.01$	$0.161 \pm 0.01$	$0.126^{\dagger} \pm 0$
		LinearIDNetwork	$0.149 \pm 0.01$	$0.061 \pm 0.01$	$0.175^\dagger \pm 0.02$	$0.127^{\ddagger} \pm 0$
		IDNetwork	$\textbf{0.145} \pm \textbf{0.01}$	$\textbf{0.059} \pm \textbf{0.01}$	$0.164 \pm 0.01$	$0.122 \pm 0.$
	$M_{0.001}^{\mathrm{BH}}$	msCox	$\textbf{0.145} \pm \textbf{0.01}$	$0.067^* \pm 0.01$	$\textbf{0.159}^\dagger \pm \textbf{0.01}$	$\textbf{0.124} \pm \textbf{0}$
		msSplineCox	$0.146 \pm 0.01$	$0.065^\dagger \pm 0.01$	$0.165 \pm 0.01$	$0.125 \pm 0$
		LinearIDNetwork	$0.146 \pm 0.01$	$0.060\pm0.01$	$0.173 \pm 0.01$	$0.126 \pm 0$
		IDNetwork	$0.146 \pm 0.01$	$\boldsymbol{0.059 \pm 0.01}$	$0.166 \pm 0.02$	$\boldsymbol{0.124 \pm 0}$
	$M_{0.005}^{\mathrm{BH}}$	msCox	$\textbf{0.145} \pm \textbf{0.01}$	$0.068^* \pm 0.01$	$\textbf{0.164} \pm \textbf{0.01}$	$0.126 \pm 0$
		msSplineCox	$0.150\pm0.01$	$0.071^{\ddagger} \pm 0.00$	$0.169 \pm 0.01$	$0.131 \pm 0$
		LinearIDNetwork	$0.146 \pm 0.01$	$0.060 \pm 0.01$	$0.178^{\circ} \pm 0.01$	$0.128 \pm 0$
		IDNetwork	$0.146 \pm 0.01$	$\boldsymbol{0.058 \pm 0.01}$	$0.168 \pm 0.02$	$0.125 \pm 0$
	$\mathrm{M}_{0.01}^{\mathrm{BH}}$	msCox	$0.146 \pm 0.01$	$0.068^* \pm 0.01$	$\textbf{0.167} \pm \textbf{0.01}$	$0.126 \pm 0$
		msSplineCox	$0.148 \pm 0.01$	$0.069^* \pm 0.00$	$0.174 \pm 0.01$	$0.130 \pm 0$
		LinearIDNetwork	$0.144 \pm 0.01$	$\boldsymbol{0.060 \pm 0.01}$	$0.175 \pm 0.01$	$0.127 \pm 0$
		IDNetwork	$\textbf{0.145} \pm \textbf{0.01}$	$\boldsymbol{0.060 \pm 0.01}$	$0.171 \pm 0.02$	$\textbf{0.124} \pm \textbf{0}$
	$M_{0.05}^{\mathrm{BH}}$	msCox	$0.149^\dagger \pm 0.01$	$0.071^* \pm 0.01$	$0.185^\dagger \pm 0.02$	0.134* ±
		msSplineCox	$0.144 \pm 0.01$	$0.068^{\circ} \pm 0.00$	$\textbf{0.165} \pm \textbf{0.03}$	$0.126 \pm 0$
		LinearIDNetwork	$0.144 \pm 0.01$	$0.060 \pm 0.01$	$0.175^\dagger \pm 0.01$	$0.127^{\dagger} \pm 0$
		IDNetwork	$\textbf{0.142} \pm \textbf{0.01}$	$\textbf{0.059} \pm \textbf{0.01}$	$0.167 \pm 0.01$	$0.123 \pm 0$
	$\mathbf{M}_{0.1}^{\mathrm{BH}}$	msCox	$0.157^* \pm 0.01$	$0.072^* \pm 0.01$	$0.207^* \pm 0.01$	0.145* ±
		msSplineCox	$0.162^\dagger \pm 0.01$	$0.066 \pm 0.01$	$0.201^{\ddagger} \pm 0.01$	$0.143^{\ddagger} \pm 0$
		LinearIDNetwork	$\textbf{0.145} \pm \textbf{0.01}$	$\boldsymbol{0.060 \pm 0.01}$	$0.178^\dagger \pm 0.02$	$0.127^{\dagger} \pm 0$
		IDNetwork	$\boldsymbol{0.145 \pm 0.01}$	$\boldsymbol{0.060 \pm 0.01}$	$\boldsymbol{0.165 \pm 0.01}$	$\textbf{0.124} \pm \textbf{0}$

*Note*: We integrate the AUC and BS measures at all the 30 equidistant time points in [90,  $\tau$ ] (ie, at every months from 3 months). The model  $M_0$  uses only the clinical features and each model  $M_\alpha^{\rm BH}$ , for  $\alpha \in [0.001, 0.005, 0.01, 0.05, 0.1]$ , uses the clinical features and the genes selected with a corrected *P*-value less than  $\alpha$ .

Prognostication of diseases is a key momentum in the medical decision process of various diseases, for example to identify population at risk of cardiovascular complications or to identify population at risk of cancer relapse. The most used approaches are based on nomograms and scores computed based on a few clinical and biological features. For example, the CHAD2-DS2-VASC score is commonly used by physicians to identify patients who require anticoagulant treatment following the diagnosis of a cardiac atrial arrhythmia. <sup>53</sup> In the same way, the RSClin tool that combines clinical, pathological and genetic information has been developed in oncology to predict the risk of breast cancer relapse and to determine more precisely which patients need chemotherapy in addition to surgery. <sup>54</sup> But a lot of medical information contained in the patients' records is left unexploited. In parallel to the development of new biomarkers relying on highly specialized technologies, another approach is to focus on the optimization of prognostic models based on large but accessible information. In the area of digital medicine which aims to capture and combine a huge amount of data, healthcare workers will need specialized algorithms to support their practices based on standardized guidelines on the one hand, and on personalized assessment on the other hand. IDNetwork has been developed in that sense.

IDNetwork uses a multi-task architecture and transition-specific subnetworks to learn an estimation of the density probabilities of occurrence of state transitions of an illness death process. It uses piecewise approximations to provide

accurate predictions of the cumulative probabilities of state transitions. IDNetwork uses multiple fully connected layers and nonlinear activation functions to model the relationships between covariates and risks of transitions without any assumption. It is trained by minimizing a loss function designed to both capture the relationships between covariates and risks of transitions and provide smooth piecewise approximations of the density probabilities.

Through experiments on simulated data sets, we investigate different configurations of IDNetwork and illustrate the benefit of our loss function. We compare the predictive performances of our method with the state-of-the-art methods using discrimination and calibration criteria. We evaluate IDNetwork in predicting the cumulative probabilities of state transitions on a simulated data set and on real data sets on colon and breast cancer. We show that IDNetwork provides significant improvements compared with the others methods.

Furthermore, medical decision-making requires to combine heterogeneous individual features. On the real data set on breast cancer, we illustrate how IDNetwork can be easily adapted to integrate various types of data (as clinical, biological, molecular, gene expression) and displays in this case coherent and significant improvements in comparison with statistical methods.

Developing a reliable multi-state model with a small training set and few events can be challenging and can result in a poor predictive model. Time-to-event data collected in a clinical setting can suffer from a high censoring rate, especially for rare events. Applying survival models on data sets with fewer events than censored observations can impair the risks of events estimates. In addition, deep learning methods require large amount of training data. This is rarely the case with data sets from clinical trials where the number of patients is relatively low compared to the databases commonly used in deep learning. We adapt the architecture of IDNetwork to handle these limitations by using hyper-parameters to simplify the architecture when available training data may require it for each of the three transitions. Moreover, given the increasing popularity of using real-world data (RWD) collection, such as electronic health records (EHRs) or disease registries, we aim to exploit these data in the future for our problematic to make more efficient and reliable decision-making.

Explaining predictions in deep neural networks is challenging but essential in clinical applications where interpretability and reliability are evaluated to support medical decisions. For the future work, it may be relevant to add a model interpretability functionality to IDNetwork in order to help clinicians in establishing the patient prognosis. It could reveal what are the patient characteristics associated with each transition and increase the understanding of the evolution of the disease. Some methods have already been developed to explain predictions in deep neural networks.<sup>55</sup> But these methods are mainly used in the domain of image classification and are rarely developed in contexts like ours. It will be therefore necessary to adapt these methods for multi-state analysis in order to understand the relationship between covariates and risks of transitions. It would be interesting to illustrate the interaction effects between covariates in IDNetwork as well. In particular, it could reveal the role of gene interactions in the variation of the risks.

IDNetwork is a flexible method developed for an illness-death process and can readily be applied in many cancers and cardiovascular diseases to predict two-stage evolution. It can be generalized to embrace more complex disease evolution patterns by adapting the states and transitions. For example, the evolution of a disease should be modeled with multiple intermediate states to reflect the long-term patient journey composed of sequential treatments. Therefore, IDNetwork could provide a more realistic prediction of disease evolution through different phases. At the end, it could support medical decisions not only based on a single event as done by traditional approaches, but by anticipating relevant outcomes closed to the disease course of a real patient.

In the same way, IDNetwork could be generalized to model time-varying covariates (rather than only baseline covariates) by adapting the model architecture and the loss function. Incorporating time-varying features could provide dynamic updated predictions of the cumulative incidence functions based on all the patients' characteristics measured since the beginning of their clinical follow-up.

# **ACKNOWLEDGEMENTS**

This publication is based on research using information obtained from www.projectdatasphere.org, which is maintained by Project Data Sphere. Neither Project Data Sphere nor the owner(s) of any information from the web site have contributed to, approved or are in any way responsible for the contents of this publication.

#### CONFLICT OF INTEREST

The authors declare no potential conflict of interests.

#### DATA AVAILABILITY STATEMENT

We use two data sets from Phase III clinical trials in non-metastatic colon cancer. (1) The data from the clinical trial NCT00079274<sup>44</sup> is available under request at https://data.projectdatasphere.org/projectdatasphere/html/content/161.

(2) The data from the clinical trial NCT00275210<sup>45</sup> is available under request at https://data.projectdatasphere.org/projectdatasphere/html/content/128.

We use also the data from the Molecular Taxonomy of Breast Cancer International Consortium (METABRIC) cohort. <sup>46</sup> The METABRIC data set is available at https://www.cbioportal.org/study/clinicalData?id=brca metabric.

## **ENDNOTES**

- \*https://www.tensorflow.org/.
- †The clinical trial NCT00079274 data set is available under request at https://data.projectdatasphere.org/projectdatasphere/html/content/161.
- <sup>‡</sup>The clinical trial NCT00275210 data set is available under request at https://data.projectdatasphere.org/projectdatasphere/html/content/128.
- §The METABRIC data set is available at https://www.cbioportal.org/study/clinicalData?id=brca\_metabric.
- The Breast Invasive Carcinoma TCGA PanCancer data set is available at https://www.cbioportal.org/study/clinicalData?id=brca\_tcga\_pan\_can atlas 2018.

#### ORCID

Aziliz Cottin https://orcid.org/0000-0003-4394-7002

#### REFERENCES

- 1. Webster AJ. Multi-stage models for the failure of complex systems, cascading disasters, and the onset of disease. *PLoS One*. 2019;14(5):e0216422.
- 2. Eulenburg C, Mahner S, Woelber L, Wegscheider K. A systematic model specification procedure for an illness-death model without recovery. *PLoS One*. 2015;10(4):e0123489.
- 3. Iacobelli S, Carstensen B. Multiple time scales in multi-state models. Stat Med. 2013;32(30):5315-5327.
- Commenges D, Joly P, Letenneur L, Dartigues J-F. Incidence and mortality of Alzheimer's disease or dementia using an illness-death model. Stat Med. 2004;23(2):199-210.
- 5. Ramezankhani A, Blaha MJ, Mirbolouk M, Azizi F, Hadaegh F. Multi-state analysis of hypertension and mortality: application of semi-Markov model in a longitudinal cohort study. *BMC Cardiovasc Disord*. 2020;20(1):1-13.
- 6. Kaplan EL, Meier P. Nonparametric estimation from incomplete observations. J Am Stat Assoc. 1958;53(282):457-481.
- 7. Aalen O. Nonparametric inference for a family of counting processes. Ann Stat. 1978;6(4):701-726.
- 8. Cox DR. Regression models and life-tables. J Royal Stat Soc Ser B (Methodol). 1972;34(2):187-202.
- 9. De Wreede LC, Fiocco M, Putter H. The mstate package for estimation and prediction in non-and semi-parametric multi-state and competing risks models. *Comput Methods Prog Biomed*. 2010;99(3):261-274.
- 10. Faraggi D, Simon R. A neural network model for survival data. Stat Med. 1995;14(1):73-82.
- 11. Luck M, Sylvain T, Cardinal H, Lodi A, Bengio Y. Deep learning for patient-specific kidney graft survival analysis; 2017. arXiv preprint arXiv:1705.10245.
- 12. 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 Med Res Methodol*. 2018;18(1):24.
- 13. Fotso S. Deep neural networks for survival analysis based on a multi-task framework; 2018. arXiv preprint arXiv:1801.05512.
- 14. Kvamme H, Borgan Ø, Scheel I. Time-to-event prediction with neural networks and Cox regression. J Mach Learn Res. 2019;20(129):1-30.
- 15. Giunchiglia E, Nemchenko A, Schaar M. Rnn-surv: a deep recurrent model for survival analysis. Proceedings of the International Conference on Artificial Neural Networks; 2018:23-32; Springer, Cham.
- 16. Ren K, Qin J, Zheng L, et al. Deep recurrent survival analysis; 2019:4798-4805.
- 17. Gensheimer MF, Narasimhan B. A scalable discrete-time survival model for neural networks. PeerJ. 2019;7:e6257.
- 18. Lee C, Zame WR, Yoon J, Schaar M. Deephit: a deep learning approach to survival analysis with competing risks; 2018.
- 19. Kalbfleisch JD, Prentice RL. The Statistical Analysis of Failure Time Data. Hoboken, NJ: John Wiley & Sons; 2011.
- 20. Friedman M. Piecewise exponential models for survival data with covariates. Ann Stat. 1982;10(1):101-113.
- 21. Kyamme H, Borgan Ø. Continuous and discrete-time survival prediction with neural networks; 2019. arXiv preprint arXiv:1910.06724.
- 22. Andersen PK, Borgan O, Gill RD, Keiding N. Statistical Models Based on Counting Processes. Berlin, Germany: Springer Science & Business Media; 2012.
- 23. Andersen PK, Borgan Ø. Counting process models for life history data: a review. Preprint series statistical research report NBN: no-23420. 1984. http://urn nb no/URN.
- 24. Meira-Machado L, de Uña-Álvarez J, Cadarso-Suárez C, Andersen PK. Multi-state models for the analysis of time-to-event data. *Stat Methods Med Res.* 2009;18(2):195-222.
- 25. Jackson CH. flexsurv: a platform for parametric survival modeling in R. J Stat Softw. 2016;70:1-25.
- 26. Andersen PK, Perme MP. Inference for outcome probabilities in multi-state models. Lifetime Data Anal. 2008;14(4):405.
- 27. Martinussen T, Scheike TH. Dynamic Regression Models for Survival Data. Berlin, Germany: Springer Science & Business Media; 2007.
- 28. Murphy SA, Sen PK. Time-dependent coefficients in a Cox-type regression model. Stoch Process Appl. 1991;39(1):153-180.
- 29. Achab M, Guilloux A, Gaïffas S, Bacry E. SGD with variance reduction beyond empirical risk minimization; 2015. arXiv preprint arXiv:1510.04822.
- 30. Meira-Machado L, Sestelo M. Estimation in the progressive illness-death model: a nonexhaustive review. Biom J. 2019;61(2):245-263.

- 31. Wey A, Salkowski N, Kremers W, Ahn YS, Snyder J. Piecewise exponential models with time-varying effects: estimating mortality after listing for solid organ transplant. *Stat.* 2020;9(1):e264.
- 32. Triebel H. Theory of Function Spaces. Monographs in Mathematics. Vol 78; Switzerland: Birkhäuser Basel; 1983.
- 33. Putter H, Fiocco M, Geskus RB. Tutorial in biostatistics: competing risks and multi-state models. Stat Med. 2007;26(11):2389-2430.
- 34. Ruder S. An overview of multi-task learning in deep neural networks; 2017. arXiv preprint arXiv:1706.05098.
- 35. Lee C, Yoon J, Van Der Schaar M. Dynamic-deephit: a deep learning approach for dynamic survival analysis with competing risks based on longitudinal data. *IEEE Trans Biomed Eng.* 2019;67(1):122-133.
- 36. Möst S. Regularization in Discrete Survival Models. PhD thesis. LMU; 2014.
- 37. Tibshirani R, Saunders M, Rosset S, Zhu J, Knight K. Sparsity and smoothness via the fused lasso. *J Royal Stat Soc Ser B (Stat Methodol)*. 2005;67(1):91-108.
- 38. Aalen OO, Johansen S. An empirical transition matrix for non-homogeneous Markov chains based on censored observations. *Scand J Stat.* 1978;5(3):141-150.
- 39. Jacqmin-Gadda H, Blanche P, Chary E, Touraine C, Dartigues J-F. Receiver operating characteristic curve estimation for time to event with semicompeting risks and interval censoring. *Stat Methods Med Res.* 2016;25(6):2750-2766.
- 40. Spitoni C, Lammens V, Putter H. Prediction errors for state occupation and transition probabilities in multi-state models. *Biom J*. 2018;60(1):34-48.
- 41. Royston P, Altman DG. External validation of a Cox prognostic model: principles and methods. BMC Med Res Methodol. 2013;13(1):33.
- 42. Wilcoxon F. Individual comparisons by ranking methods. Breakthroughs in Statistics. New York, NY: Springer; 1992:196-202.
- 43. Bender R, Augustin T, Blettner M. Generating survival times to simulate Cox proportional hazards models. *Stat Med.* 2005;24(11):1713-1723.
- 44. Alberts SR, Sargent DJ, Nair S, et al. Effect of oxaliplatin, fluorouracil, and leucovorin with or without cetuximab on survival among patients with resected stage III colon cancer: a randomized trial. *JAMA*. 2012;307(13):1383-1393.
- 45. André T, Boni C, Mounedji-Boudiaf L, et al. Oxaliplatin, fluorouracil, and leucovorin as adjuvant treatment for colon cancer. *N Engl J Med*. 2004;350(23):2343-2351.
- 46. Curtis C, Shah SP, Chin S-F, et al. The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups. *Nature*. 2012;486(7403):346-352.
- 47. Van Wieringen WN, Kun D, Hampel R, Boulesteix A-L. Survival prediction using gene expression data: a review and comparison. *Comput Stat Data Anal.* 2009;53(5):1590-1603.
- 48. Yousefi S, Amrollahi F, Amgad M, et al. Predicting clinical outcomes from large scale cancer genomic profiles with deep survival models. *Sci Rep.* 2017;7(1):1-11.
- 49. Ching T, Zhu X, Garmire LX. Cox-nnet: an artificial neural network method for prognosis prediction of high-throughput omics data. *PLoS Comput Biol.* 2018;14(4):e1006076.
- 50. Kuleshov MV, Jones MR, Rouillard AD, et al. Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Res.* 2016;44(W1):W90-W97.
- 51. Liberzon A, Birger C, Thorvaldsdóttir H, Ghandi M, Mesirov JP, Tamayo P. The molecular signatures database hallmark gene set collection. *Cell Syst.* 2015;1(6):417-425.
- 52. Wu Y, Sarkissyan M, Vadgama JV. Epithelial-mesenchymal transition and breast cancer. J Clin Med. 2016;5(2):13.
- 53. Lip GYH, Nieuwlaat R, Pisters R, Lane DA, Crijns HJGM. Refining clinical risk stratification for predicting stroke and thromboembolism in atrial fibrillation using a novel risk factor-based approach: the euro heart survey on atrial fibrillation. *Chest.* 2010;137(2):263-272.
- 54. Sparano JA, Crager MR, Tang G, Gray RJ, Stemmer SM, Shak S. Development and validation of a tool integrating the 21-gene recurrence score and clinical-pathological features to individualize prognosis and prediction of chemotherapy benefit in early breast cancer. *J Clin Oncol.* 2021;39(6):557-564.
- 55. Ancona M, Ceolini E, Öztireli C, Gross M. Towards better understanding of gradient-based attribution methods for deep neural networks; 2017. arXiv preprint arXiv:1711.06104.
- 56. Jackson C. Multi-state Modelling with R: The MSM Package. Cambridge, UK: Cambridge University Press; 2007:1-53.

#### SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of this article.

**How to cite this article:** Cottin A, Pecuchet N, Zulian M, Guilloux A, Katsahian S. IDNetwork: A deep illness-death network based on multi-state event history process for disease prognostication. *Statistics in Medicine*. 2022;41(9):1573-1598. doi: 10.1002/sim.9310