



Assessing model prediction performance for the expected cumulative number of recurrent events

Olivier Bouaziz¹

Received: 27 January 2023 / Accepted: 19 September 2023 / Published online: 17 November 2023
© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2023

Abstract

In a recurrent event setting, we introduce a new score designed to evaluate the prediction ability, for a given model, of the expected cumulative number of recurrent events. This score can be seen as an extension of the Brier Score for single time to event data but works for recurrent events with or without a terminal event. Theoretical results are provided that show that under standard assumptions in a recurrent event context, our score can be asymptotically decomposed as the sum of the theoretical mean squared error between the model and the true expected cumulative number of recurrent events and an inseparability term that does not depend on the model. This decomposition is further illustrated on simulations studies. It is also shown that this score should be used in comparison with a reference model, such as a nonparametric estimator that does not include the covariates. Finally, the score is applied for the prediction of hospitalisations on a dataset of patients suffering from atrial fibrillation and a comparison of the prediction performances of different models, such as the Cox model, the Aalen Model or the Ghosh and Lin model, is investigated.

Keywords Recurrent events · Prediction assessment · Right-censoring · Terminal event · Brier score

1 Introduction

Recurrent event data are often encountered in follow-up studies. They can be seen as a generalisation of the standard time to event data, where individuals may experience the same event repeatedly over time. Typical examples may include HIV studies where patients can experience repeated opportunistic infections, remission data from Leukemia patients who can experience multiple relapses,

✉ Olivier Bouaziz
olivier.bouaziz@parisdescartes.fr

¹ Université Paris Cité, CNRS, MAP5, F-75006 Paris, France

repeated seizures for epileptic patients, or hospitalisation data where the events of interest are the hospitalisations. In those studies, the focus might be on assessing the effect of covariates on the risk of recurrences or on predicting the future recurrences. The first model to deal with recurrent event data was the Andersen-Gill model (Andersen and Gill 1982) which was further extended by Lin et al. (2000) to account for possibly dependent jumps of the recurrent event process. Further models were developed such as in Prentice et al. (1981), Cook and Lawless (1997), Ghosh and Lin (2002), Ghosh and Lin (2003) or Andersen et al. (2019) where the last four papers incorporate the presence of a terminal event in the estimation procedure, or using random effects such as in Hougaard (2000), Liu et al. (2004), or Rondeau et al. (2007). In particular, in Cook and Lawless (1997), Ghosh and Lin (2002), Ghosh and Lin (2003), Andersen et al. (2019) the authors focused on the estimation of the expected cumulative number of recurrent events. This is a marginal quantity that computes the expectation of the number of experienced events of an individual before any time point. This quantity is particularly interesting as it summarises the evolution of the recurrent event process with time. In the presence of a terminal event, it also includes the fact that when the terminal event occurs the patient can no longer experience any further recurrent events.

In some studies the focus is more on the predictiveness ability of a model rather than on the interpretation of the covariates effects. This is the case when clinicians aim at predicting the future repeated events in order to offer the best medical care. Being able to predict the future recurrences of any patients on a short time period also allows to predict the future burden of the disease over the patient's life. Moreover, a predictive model can be an important tool for making medical decisions but also for communicating with the patients about the future course of his/her disease. For instance, in Schroder et al. (2019) the authors studied patients with atrial fibrillation, a well known cardiac disease, in an attempt to predict the future hospitalisations of patients due to their disease. Since the patients suffering from this disease are usually old (the median age in the study was 63 years) and since atrial fibrillation can be a severe disease in some cases, those patients were also at risk of death. Several covariates were collected and a prediction of the expected cumulative number of recurrent events over time was performed using a Cox model with dependence on prior counts.

While such models are certainly of interest for clinicians, it is important to propose relevant diagnosis tools that can evaluate the prediction performance of the proposed model. There already exists several indicators for prediction performances in the standard context of time to event data with only one event per individual. The Brier score was developed in Graf et al. (1999) and in Gerds and Schumacher (2006), which basically is a score for computing the mean squared error of the time to event in the presence of censoring. This score was further developed to deal with random effect models in Van Oirbeek and Lesaffre (2016), or to evaluate the performance of dynamic prediction models in Schoop et al. (2011) where the information available from a longitudinal covariate is updated at each time point. Note also that other types of predictive accuracy measures exist, called discrimination measures,

such as the C-index (see Harrell et al. 1996; Gerds et al. 2013) or the time dependent ROC curve and area under the curve (see for instance Heagerty and Zheng (2005)).

In this paper, the aim is to derive a predictive accuracy measure for recurrent events models where the focus is on predictiveness rather than discrimination. The quantity of interest is solely the expected cumulative mean number of recurrent events. Since no mean squared error measure, such as the Brier score, exists in the context of recurrent events, the goal of this work is to fill in this gap by deriving a new score of this type for recurrent events, which also accommodates for the presence of a terminal event. In this work, we show that this score reduces to the Brier score when only one event per individuals can occur and hence can be seen as a direct generalisation of the standard Brier score. Also, since our prediction criterion focuses on the marginal quantity of the expected cumulative number of recurrent events, it provides a summary score that takes into account the prediction of all recurrent events. In the context of a terminal event, it also incorporates the quality of prediction of the terminal event.

In Sect. 2.1, we introduce the general prediction criterion for recurrent events, denoted $\widehat{\text{MSE}}$. This criterion is very general and can work under right-censoring, for situations with a terminal event and for rate or intensity based models. In Sect. 2.2, we discuss the modelling assumptions in the presence of a terminal event. In Sect. 2.3, we present some existing estimators for the expected cumulative number of recurrent events. In Sect. 3, we derive the main theoretical results of this paper. We first introduce a theoretical criterion and show that it can be decomposed into an inseparability term and an imprecision term, similarly to the results in Gerds and Schumacher (2006). The former does not depend on the model and cannot be removed while the latter is exactly the mean squared error between the recurrent event process and the prediction model of the expected cumulative mean number of recurrent events. We then show that our prediction criterion asymptotically converges towards the theoretical criterion. In Sect. 4, we demonstrate that when individuals can only experience one event, our prediction criterion is equivalent to the standard Brier score. In Sect. 5, a simulation study is conducted. First, the decomposition between inseparability and imprecision terms is illustrated. As the inseparability is, by far, the dominant term, we then recommend to consider as a prediction score, the difference of $\widehat{\text{MSE}}$ s between the considered model and a reference model. Second, we illustrate how this prediction score can be used in order to compare prediction models. In Sect. 6, the atrial fibrillation dataset is studied. We show that the model with dependence on prior counts (which is a multi-state model), stratified with respect to atrial fibrillation type, provides the best prediction performance among all other models considered.

2 Prediction criterion for the expected cumulative number of recurrent events

2.1 The prediction criterion in the general framework

In this section we present a prediction criterion for a recurrent event setting under right-censoring and a terminal event. We define a counting process of interest $N^*(t)$ which counts the number of recurrent events that have occurred before time t . A terminal event T^* is further introduced such that this counting process cannot jump after T^* . We denote $T := T^* \wedge C$ and $N(t) := N^*(t \wedge C)$ where C is a censoring variable and $a \wedge b$ represents the minimum between a and b . We assume that a multivariate left-continuous external time dependent covariate vector $X(t)$ (see Kalbfleisch and Prentice (2002) for the definition of external covariates) is observed and we note $\bar{X}(t) = \{X(u) : 0 \leq u \leq t\}$ the history of the covariate process up until time t . We define \mathcal{M} a class of bounded functions depending on t and $\bar{X}(t)$. For each $t \geq 0$, we denote the support of all possible sample paths for the process $\bar{X}(t)$ by \mathcal{X}_t . We also note τ the endpoint of the study. On the basis of i.i.d. replications $(N_i(t), X_i(t) : 0 \leq t \leq \tau)$, we will assume that an estimator $\hat{\mu} \in \mathcal{M}$ of the expected cumulative number of recurrent events $\mu^*(t | \bar{X}(t)) := \mathbb{E}[N^*(t) | X(u) : 0 \leq u \leq t]$ is available. We will say that this estimator is consistent if there exists $\mu \in \mathcal{M}$ such that for all $t \leq \tau$,

$$\sup_{x \in \mathcal{X}_t} |\hat{\mu}(t | x) - \mu(t | x)| \rightarrow 0, \text{ in probability as } n \rightarrow \infty.$$

The main goal of this paper is to develop a new mean squared error criterion designed to evaluate the performance of this estimator.

We propose to evaluate the prediction ability of a given estimator $\hat{\mu} \in \mathcal{M}$, through the following criterion:

$$\widehat{\text{MSE}}(t, \hat{\mu}) = \frac{1}{n} \sum_{i=1}^n \left(\int_0^t \frac{dN_i(u)}{\hat{G}_c(u | \bar{X}_i(u))} - \hat{\mu}(t | \bar{X}_i(t)) \right)^2, \tag{1}$$

where $\hat{G}_c(u | \bar{X}_i(u)) = 1 - \hat{G}(u- | \bar{X}_i(u))$ is an estimator of $G_c(u | \bar{X}_i(u)) = 1 - G(u- | \bar{X}_i(u))$, the conditional survival function of the censoring variable C given $\bar{X}(\cdot)$. The notation $\hat{G}(u- | \bar{X}_i(u))$ indicates the left limit of the function \hat{G} at u . We will assume uniform consistency of this censoring estimator in the following way.

Assumption 1 Let \mathcal{G} be a model for the conditional censoring distribution. We say that \hat{G} is a uniformly consistent estimator for $G \in \mathcal{G}$ if for all $t \leq \tau$,

$$\sup_{x \in \mathcal{X}_t} |\hat{G}(t | x) - G(t | x)| \rightarrow 0, \text{ in probability as } n \rightarrow \infty.$$

Presentations of different estimators for G are discussed in Sect. 2.2.

We now introduce a theoretical criterion that would be available if the censoring distribution was known. For some function $\mu \in \mathcal{M}$, let:

$$\text{MSE}(t, \mu) = \mathbb{E} \left[\left(\int_0^t \frac{dN(u)}{G_c(u- | \bar{X}(u))} - \mu(t | \bar{X}(t)) \right)^2 \right]. \tag{2}$$

The crucial idea behind our criterion (1) comes from the fact that

$$\mathbb{E} \left[\int_0^t \frac{dN(u)}{G_c(u | \bar{X}(u))} \right] = \mathbb{E}[\mu^*(t | \bar{X}(t))], \tag{3}$$

a relationship that is proved in Sect. 2.2. In Sect. 3, we provide theoretical results that justify the appropriateness of the proposed criterion. In particular, Proposition 1 of Sect. 3 shows that the theoretical criterion can be decomposed in the following way:

$$\text{MSE}(t, \mu) = \mathbb{E} \left[\left(\mu^*(t | \bar{X}(t)) - \mu(t | \bar{X}(t)) \right)^2 \right] + A(t),$$

with $A(t)$ not depending on μ . The first term is an imprecision term and the second term is an inseparability (or residual) term that does not depend on the chosen model. It should be noted that this kind of result is similar to the imprecision/inseparability decomposition of the Brier score (see Gerds and Schumacher 2006). This result shows that it is a relevant purpose to aim at estimating the expectation in Eq. (2) as it is a valid surrogate for the mean squared error between μ^* and μ . Since the inseparability term is not affected by the value of μ , the quantity $\text{MSE}(t, \mu)$ can be used to compare different models in terms of their mean squared error. Of course, this theoretical criterion depends on unknown quantities and needs to be estimated in practice. As a matter of fact, Proposition 2 of Sect. 3 states that if $\hat{\mu}$ is a consistent estimator for some $\mu \in \mathcal{M}$ then as n tends to infinity, our empirical criterion (1) is asymptotically equivalent to the theoretical criterion (2) evaluated at μ . Combining Propositions 1 and 2 thus shows that our empirical prediction criterion (1) can be used to evaluate the mean squared error between μ and μ^* . Such a criterion is especially useful when one wants to compare different models in terms of their prediction performances. One common practice (see Steyerberg et al. 2010) is then to use a reference model, typically a model that does not use any covariate, and to compare the prediction ability of each model relatively to the reference model. To do so, we introduce the score criterion, defined as:

$$\text{Score}(t, \hat{\mu}, \hat{\mu}_0) = \widehat{\text{MSE}}(t, \hat{\mu}_0) - \widehat{\text{MSE}}(t, \hat{\mu}), \tag{4}$$

where $\hat{\mu}_0$ is the reference model. Using this score offers the advantage that the inseparability term will cancel out in the difference and as a result, this score criterion will asymptotically converge towards the difference between the two mean squared errors of the two models. In the rest of the paper, we will call prediction criterion the

quantity $\widehat{\text{MSE}}$ defined in Eq. (1) and prediction score the quantity defined in Eq. (4). Of note, we have chosen to define the score as the difference between the prediction criterion of the null model and of the model of interest such that the larger the score the better the model in terms of prediction accuracy. In the simulation section, we will also show that the inseparability term tends to be very large as compared to the imprecision term which limits the interpretation of the prediction criterion and advocates for the use of the prediction score instead.

When analysing recurrent events it is common to model the recurrent event increments. But because Eq. (3) needs to hold, the validity of the criterion will depend on the modelling assumptions made on those increments. In the next section we will specify those assumptions under which Propositions 1 and 2 hold when a rate model is used, that is when one models the probability that a jump occurs given some covariates. Other models are possible, typically based on the intensity of the recurrent event process, such as models that depend on prior counts. Our criterion (1) will also be valid with such models. However, the required assumptions in this context need to be slightly adapted and are specified in the Supplementary Information. Finally, our prediction criterion is also valid when there is no terminal event and the recurrent events are right-censored. In that case, the precise assumptions needed for our Propositions to hold are also specified in the Supplementary Information.

2.2 Assumptions when a rate model is used

In this section, we consider the following rate model (see e.g. Cook and Lawless 2007; Scheike 2002):

$$\mathbb{E}[dN^*(t) \mid I(T^* \geq t), \bar{X}(t)] = I(T^* \geq t) \lambda^*(t \mid \bar{X}(t)) dt, \quad (5)$$

where $\lambda^*(t \mid \bar{X}(t))$ is the true rate function. In real-data analysis situations, a terminal event often occurs, typically caused by death which precludes the occurrence of further recurrences. Under this model, we observe that

$$\mu^*(t \mid \bar{X}(t)) = \mathbb{E}[N^*(t) \mid \bar{X}(t)] = \int_0^t S(u \mid \bar{X}(u)) \lambda^*(u \mid \bar{X}(u)) du,$$

where $S(t \mid \bar{X}(t)) := \mathbb{P}[T^* \geq t \mid \bar{X}(t)]$ is the conditional survival function of the terminal event. This implies that in order to define an estimator of $\mu^*(t \mid \bar{X}(t)) = \mathbb{E}[N^*(t) \mid \bar{X}(t)]$ one usually needs to also model the hazard rate for the terminal event and to derive an estimator of the conditional survival function. As a result our prediction criterion will both take into account the predictive performance of the survival function and of the rate function of N^* since, if one of those two estimators behaves poorly, the resulting estimator for μ^* is likely to perform badly as well.

We assume independent censoring in the following way:

$$\mathbb{E}[dN^*(t) \mid I(T^* \geq t), \bar{X}(t)] = \mathbb{E}[dN^*(t) \mid I(T \geq t), \bar{X}(t)].$$

We denote $Y(t) = I(T \geq t)$ the observed at-risk process and $N(t) = N^*(T \wedge t)$ the observed counting process. Under the independent censoring assumption, it can be shown that

$$\mathbb{E}[dN(t) \mid Y(t), \bar{X}(t)] = Y(t)\lambda^*(t \mid \bar{X}(t))dt. \quad (6)$$

We assume Assumption 1 and we make the following additional assumption.

Assumption 2 We assume that there exists a constant $\tau > 0$ and a constant $c > 0$ such that

- (1) $\forall t \in [0, \tau], \mathbb{P}[T \geq t \mid \bar{X}(t)] \geq c$ almost surely,
- (2) $N(\tau)$ is almost surely bounded by a constant.

We also assume that T^* is independent of C conditionally on $\bar{X}(\cdot)$.

Those conditions are standard in the context of regression for recurrent events with a terminal event, see Ghosh and Lin (2002) for example. Using Equality (6) one can easily observe that $\mathbb{E}[dN(t) \mid \bar{X}(t)] = S(t \mid \bar{X}(t))G_c(t \mid \bar{X}(t))\lambda^*(t \mid \bar{X}(t))dt$ under the independent censoring hypothesis. We then directly see that Eq. (3) holds.

On the basis of i.i.d. replications $(N_i(t), X_i(t) : 0 \leq t \leq \tau)$, let $\hat{\mu} \in \mathcal{M}$ be an estimator of μ^* where \mathcal{M} is a class of models that are assumed to be bounded. We propose to evaluate the prediction ability of this estimator through criterion $\overline{\text{MSE}}(t, \hat{\mu})$ defined in Eq. (1). This criterion involves an estimator of G , the conditional cumulative distribution function of the censoring variable. If C and $X(\cdot)$ are independent, one can estimate G using the Kaplan-Meier estimator by considering C to be the variable of interest that is incompletely observed due to the terminal event T^* . If C depends on X the conditional distribution of C must also be modelled. Several possible models are presented in Gerds and Schumacher (2006) such as the Cox model, the Aalen additive model, or the kernel type model of Dabrowska (1989). Alternatively, a single-index approach for right-censored data, as in Bouaziz and Lopez (2010), or the random survival forest method developed in Ishwaran et al. (2008) can be used.

Theoretical results on the validity of this criterion are derived in Sect. 3. When dealing with recurrent events, intensity based models are also possible. In particular, in some situations, it might be relevant to consider models with dependence on prior counts, as studied in Cook and Lawless (2007). Such models are multi-state models where the recurrent event intensity is allowed to change after each recurrence. They are detailed in the Supplementary Information with the appropriate assumptions needed in such context.

2.3 Examples of estimators for the expected cumulative number of recurrent events

In this section we present some estimators for the expected cumulative mean number. In the context of right-censored data and no terminal event, the expected cumulative number of recurrent events is simply equal to the cumulative hazard function: $\mu^*(t | \bar{X}(t)) = \int_0^t \lambda^*(u | \bar{X}(u))du$, see the Supplementary Information for more details. It is then common to model the rate function λ^* using the Cox (see Cox 1972), Aalen (see Scheike 2002) or Accelerated Failure Time (see Lin et al. 1998) models. We note $\hat{\Lambda}$ an estimator of Λ^* , where $\Lambda^*(t | \bar{X}(t)) = \int_0^t \lambda^*(u | \bar{X}(u))du$. The estimator of μ^* can then be expressed in the following way (see Cook and Lawless 2007):

$$\hat{\mu}(t | \bar{X}(t)) = \int_0^t d\hat{\Lambda}(u | \bar{X}(u)). \tag{7}$$

In the presence of a terminal event, a common approach is to first model the hazard rate of the terminal event (using again a Cox model for instance) and to derive an estimator of the survival function $\hat{S}(t | \bar{X}(t)) = \exp(-\int_0^t d\hat{\Lambda}^{T^*}(u | \bar{X}(u))du)$ where $\hat{\Lambda}^{T^*}$ is the estimator of the cumulative hazard rate of the terminal event. Then, the final estimator of μ^* is (see Cook and Lawless 2007; Andersen et al. 2019):

$$\hat{\mu}(t | \bar{X}(t)) = \int_0^t \hat{S}(u | \bar{X}(u))d\hat{\Lambda}(u | \bar{X}(u)). \tag{8}$$

Alternative approaches that directly model μ^* also exist. In Ghosh and Lin (2002) the authors consider the following Cox type model: $\mu^*(t | \bar{X}(t)) = \mu_0(t) \exp(X(t)^\top \beta)$. A more general approach consists in using a Single-Index-Model for estimating μ^* : in Bouaziz et al. (2015) the authors assume the existence of a nonparametric function g and a parameter β such that $\mu^*(t | \bar{X}(t)) = g(X(t)^\top \beta)$. Those two approaches provide a direct estimator of the quantity of interest μ^* . However, in the presence of a terminal event, it is no longer possible to disentangle the effects on the recurrent event process or the terminal event. As a result, the regression parameters should be interpreted with caution. See also Cook and Lawless (2007) for a discussion about this issue. Nevertheless, those estimators are appealing in the context of prediction as they do not require to separately model the hazard for the terminal event and the hazard for recurrent events.

A popular modelling approach, when dealing with recurrent events, is to use multi-state models. In particular, the dependence on prior counts model allows to specify a separate hazard risk after each new recurrence. It can also incorporate different hazards for the terminal event that differ according to the number of previous recurrent events that were already experienced. We refer the reader to Cook and Lawless (2007) and to the Supplementary Information for more details on this model. Let $\hat{\Lambda}_l(t | \bar{X}(t))$ be an estimator of the cumulative hazard for the l th recurrent event at time t knowing the covariates and that the individual is at risk (that is, he/she is alive and has already experienced $l - 1$ recurrent events) at time $t-$. Then, an estimator of μ^* can be derived as

$$\hat{\mu}(t | \bar{X}(t)) = \sum_{l=1}^L \int_0^t \hat{Q}_l(u | \bar{X}(u)) d\hat{\Lambda}_l(u | \bar{X}(u)),$$

where $\hat{Q}_l(u | \bar{X}(u))$ is an estimator of the state probability $\mathbb{P}[T^* \geq u, N^*(u-) = l - 1 | \bar{X}(u)]$. These probabilities can be easily estimated using standard R packages for multi-state models (`mstate` or `msm`). On the other hand, the hazard rates for each state can typically be modelled using Cox models and the integral will be evaluated by summing the product of the two quantities evaluated at the jumps, which will occur when an individual experiences his/her l th recurrent event.

3 Theoretical results

In this section we provide theoretical results on the proposed criterion (1) in the context of right-censoring and a terminal event. Two results are obtained. The first one is concerned with the theoretical criterion (2). It shows that this criterion applied to a function $\mu \in \mathcal{M}$ reduces to the mean squared error between μ and μ^* and a term that does not depend on μ . The second result shows the asymptotic consistency between $\widehat{\text{MSE}}(t, \hat{\mu})$ and $\text{MSE}(t, \mu)$ when the estimator $\hat{\mu}$ is an asymptotically consistent estimator of μ .

Proposition 1 *We assume Assumption 2 and independent censoring. We then have for $\mu \in \mathcal{M}$,*

$$\text{MSE}(t, \mu) = \mathbb{E} \left[\left(\mu^*(t | \bar{X}(t)) - \mu(t | \bar{X}(t)) \right)^2 \right] + A(t), \tag{9}$$

where $A(t) \geq 0$ for all $t \geq 0$ and $A(t)$ does not depend on μ .

Proposition 2 *We assume Assumptions 1 and 2, and independent censoring. Then, if the estimator $\hat{\mu} \in \mathcal{M}$ is consistent for $\mu \in \mathcal{M}$, we have*

$$\sup_{t \leq \tau} \left| \widehat{\text{MSE}}(t, \hat{\mu}) - \text{MSE}(t, \mu) \right| \rightarrow 0, \quad \text{in probability as } n \rightarrow \infty.$$

The proofs are provided in the Appendix, in Sections 8.1 and 8.2, with an explicit expression of $A(t)$. The results of the two propositions can be extended to intensity based models and the proof is also provided in the Appendix with the appropriate assumptions. As a special case, they also apply to the situation with no terminal event.

4 Link with the Brier score

The Brier score (see Gerds and Schumacher 2006) is a popular criterion to evaluate the prediction performance of a regression model for the conditional survival function in the context of right-censoring when only a single event can be observed per individual. We show in this section that if we use our criterion when individuals can only experience one event at most, then our theoretical criterion denoted MSE' and defined as follows

$$MSE'(t, \pi) = \mathbb{E} \left[\left(1 - \int_0^t \frac{dN(u)}{G_c(u | \bar{X}(u))} - \pi(t | \bar{X}(t)) \right)^2 \right], \tag{10}$$

where $\pi = 1 - \mu$, reduces to the theoretical Brier score up to a term that does not depend on the model \mathcal{M} . Note that when individuals can only experience one event, the recurrent event process reduces to $N^*(t) = I(T^* \leq t)$ and $\mu^*(t | \bar{X}(t)) = \mathbb{E}[N^*(t) | \bar{X}(t)]$ is the conditional cumulative distribution function of T^* . Since the Brier score has been designed for the prediction of the conditional survival function, we have simply rewritten our criterion in Eq. (2) such that π represents the model for the conditional survival function.

We first recall that the theoretical Brier score is defined as (see Eq. (1) from Gerds and Schumacher (2006)):

$$MSE^{Brier}(t, \pi) = \mathbb{E} [(I(T^* > s) - \pi(t | X))^2],$$

where the expectation is taken with respect to the joint distribution of T^* and X . For simplicity, the covariate X is not time dependent in the formula, as presented in the paper of Gerds and Schumacher (2006), but the results presented in this section are still valid for time dependent covariates. In their work, the authors show similar results as Propositions 1 and 2 of the present paper when the aim is to provide a prediction of the survival function $S(t | X) = \mathbb{P}[T^* > t | X]$. Note that we have suppressed the dependency with respect to S in the definition of the Brier score to stay consistent with the notations used throughout this paper. Also, in the definition of the Brier score, π plays the role of $1 - \mu$ in the present paper, that is, it is the limiting function of a proposed conditional survival estimator $\hat{S}(t | X)$. We have the following result.

Proposition 3 *We assume that only one event per individual can be experienced, that is, $N(t) = I(T \leq t, \Delta = 1)$, with $T = T^* \wedge C$ is the observed time, T^* is the true event time, C is the censoring variable and $\Delta = I(T^* \leq C)$ is the censoring indicator. Then, under independent censoring, we have:*

$$MSE'(t, \pi) = MSE^{Brier}(t, \pi) + B(t),$$

where $B(t) \geq 0$ for all $t \geq 0$ and $B(t)$ does not depend on π .

The proof is provided in the Appendix, in Section 8.3, with an explicit expression of $B(t)$. Since $B(t)$ does not depend on the model π , those two criterions are

completely equivalent. In particular, one may consider the score prediction as defined in Eq. (4) which represents the difference, in terms of MSE, of a regression model to a reference model, such as a model that does not include covariates. In that case, the two criteria will provide exactly the same values since the $B(t)$ term will cancel out in the difference. As we will see in the next section, comparing a model to a reference is typically what we recommend in practice. Since recurrent events are a generalisation of the single-event per individual situation, our criterion can be seen as an extension of the Brier score for recurrent events.

5 Simulations

5.1 A scenario with right-censoring and no terminal event

For $i = 1, \dots, n$, we first simulate a two-dimensional covariate vector $X_i = (X_{i,1}, X_{i,2})^\top$ with $X_{i,1}$ a Bernoulli variable with parameter 0.5 and $X_{i,2}$ a Gaussian variable with expectation 2 and standard deviation 0.5. Conditional on X_i , the recurrent events are generated from a non-homogeneous Poisson process with rate $\lambda^*(\cdot | X_i)$ that follows a Cox model with Weibull baseline and a two-dimensional time independent covariate. More specifically,

$$\lambda^*(t | X_i) = \lambda_0(t) \exp(\theta_0^\top X_i), \quad \lambda_0(t) = \frac{\alpha}{\beta} \left(\frac{t}{\beta}\right)^{\alpha-1},$$

with $\alpha = 2$ the shape parameter, $\beta = 0.39$ the scale parameter and $\theta_0 = (\log(2), \log(0.5))^\top$. Under this simulation setting, the true expected number of recurrent events is equal to (see Section 1 of Supplementary Information for more details):

$$\mu^*(t | X_i) = \int_0^t \lambda^*(u | X_i) du = \left(\frac{t}{\beta}\right)^\alpha \exp(\theta_0^\top X_i). \tag{11}$$

We further simulate a censoring variable C_i that follows a uniform distribution on $[0, 3]$. Using those parameters, we observe 0 or 1 recurrent event for 30% of the individuals, less or equal than 5 events for 54% of the individuals, and less or equal than 12 events for 77% of the individuals. On average, we observe approximately 8 recurrent events per individual.

Based on a single simulated sample, we first illustrate Propositions 1 and 2 when the class of models \mathcal{M} that contains μ assumes no effect of the covariates on the occurrence of the recurrent events. For this purpose, we independently simulate a training and a test samples. The training sample is used to compute $\hat{\mu}^{\text{train}}$ based on Eq. (7) where $\hat{\Lambda}$ does not depend on X and is simply the Nelson-Aalen estimator:

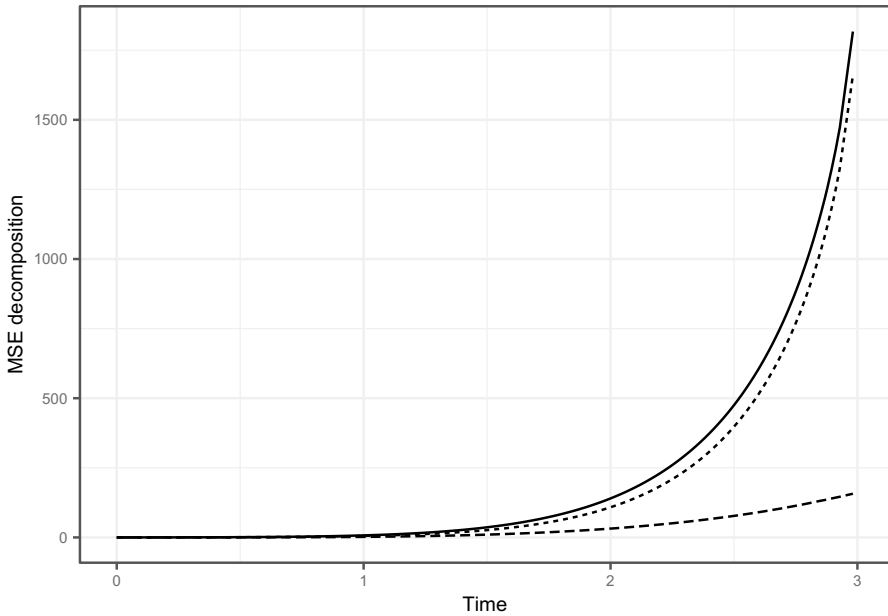


Fig. 1 Decomposition of the MSE (solid line) in Proposition 1 as the sum between the inseparability term $A(t)$ (dotted line) and the imprecision term (dashed line). The data were simulated from a Cox model with two covariates and the expected cumulative number of recurrent events was predicted using the Nelson-Aalen estimator. The train sample ($n_{\text{train}} = 200$) is used for the computation of the Nelson-Aalen estimator, the test sample ($n_{\text{test}} = 1,000$) is used for the computation of the MSE. More simulations under this scenario with no terminal event can be found in the Supplementary Information

$$\hat{\Lambda}(t) = \frac{1}{n} \sum_{i=1}^n \int_0^t \frac{dN_i(u)}{\hat{G}_c(u)},$$

with \hat{G} the Kaplan-Meier estimator of C . We then compute $\widehat{\text{MSE}}(t, \hat{\mu}^{\text{train}})$ from Eq. (1) based on the test sample of size n_{test} , that is, the computation is performed on a sample $(N_1(\cdot), X_1), \dots, (N_{n_{\text{test}}}(\cdot), X_{n_{\text{test}}})$ independent of the training sample. This quantity should provide an accurate estimation of $\text{MSE}(t, \mu)$ from Proposition 2. We then compute the imprecision term $\mathbb{E} \left[\left(\mu^*(t | X) - \mu(t | X) \right)^2 \right]$ in Eq. (9) using the true value of $\mu^*(t | \bar{X}(t))$, replacing μ by $\hat{\mu}^{\text{train}}$ and replacing the expectation by its empirical sum. In other words, we compute

$$\frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \left(\mu^*(t | X_i) - \hat{\mu}^{\text{train}}(t | X_i) \right)^2,$$

which should give a very accurate estimation of the imprecision term. The $A(t)$ term is exactly computed based on its explicit expression (see Section 3 of Supplementary Information). The decomposition of the MSE between its imprecision and inseparability terms is displayed in Fig. 1 using $n_{\text{train}} = 200$ and $n_{\text{test}} = 1,000$.

In Eq. (1) the Kaplan-Meier estimator of G was computed from the combination of the training and test samples. The solid line represents the estimated MSE while the dotted and dashed lines represent the inseparability and imprecision terms, respectively. The inseparability term is seen to be very close to the MSE. In contrast, the imprecision term, which clearly is not null here since the estimated model uses no covariates, is relatively small as compared to the other two terms. This plot suggests that it might be difficult to compare different models as the inseparability term is dominant in the decomposition of the MSE, which implies that two MSEs computed from two different models will tend to look very similar (for instance, for $t = 2.5$, the inseparability term represents approximately 84% of the value of the MSE). As a result, we advocate the use of a reference or null model and to compute the prediction score as defined in Eq. (4): for a given model, this score computes the difference between the MSE of the reference and the MSE of this model. Therefore, this score will represent the prediction gain of the model as compared to the null model. A typical choice of the null model is the one that uses no covariates. Those models will usually be implemented based on a training sample. The idea behind this score is that the inseparability term will cancel out in the difference, and the score is therefore equal to the difference between the imprecision terms of the two models. More simulations that compare the prediction scores between several models in the context of no terminal event are presented in the Supplementary Information. Also, the importance of Assumption 1 is investigated in the Supplementary Information based on data where the censoring distribution is allowed to depend on the covariates. It is seen in particular that the random survival forest model from Ishwaran et al. (2008) for the censoring distribution provides a good approach for computing the prediction criterion when the relationship between censoring and covariates is not known in advance.

5.2 A scenario with right-censoring and a terminal event

We now consider a simulation scenario which also includes a terminal event. The recurrent event process and its covariates are simulated in the same manner as in the previous section, with the same parameter values. The censoring variable is simulated following a uniform variable on $[0, 8]$. The terminal event is simulated according to a Cox model with baseline following a Weibull distribution with shape parameter equal to 5 and scale parameter equal to 1.8. This Cox model also includes the same two covariates as for the recurrent event process with the same effects on the hazard function (i.e. the effects are equal to $\log(2)$, $\log(0.5)$ for the Bernoulli and Gaussian covariates respectively). This setting leads to 8.5 events per individual on average, with 26%, 50%, 77% of individuals that experience less than or equal to 3, 7 and 12 events, respectively. On average 28% of individuals are censored.

We estimate the expected cumulative number of recurrent events based on Eq. (8) where in the formula, the Breslow estimator is used to estimate the conditional survival function of the terminal event, if the estimation model for the terminal event includes covariates. In other words:

$$\hat{S}(t | X) = \exp \left(- \int_0^t \exp(X^\top \hat{\beta}^{T*}) d\hat{\Lambda}_0^{T*} \right),$$

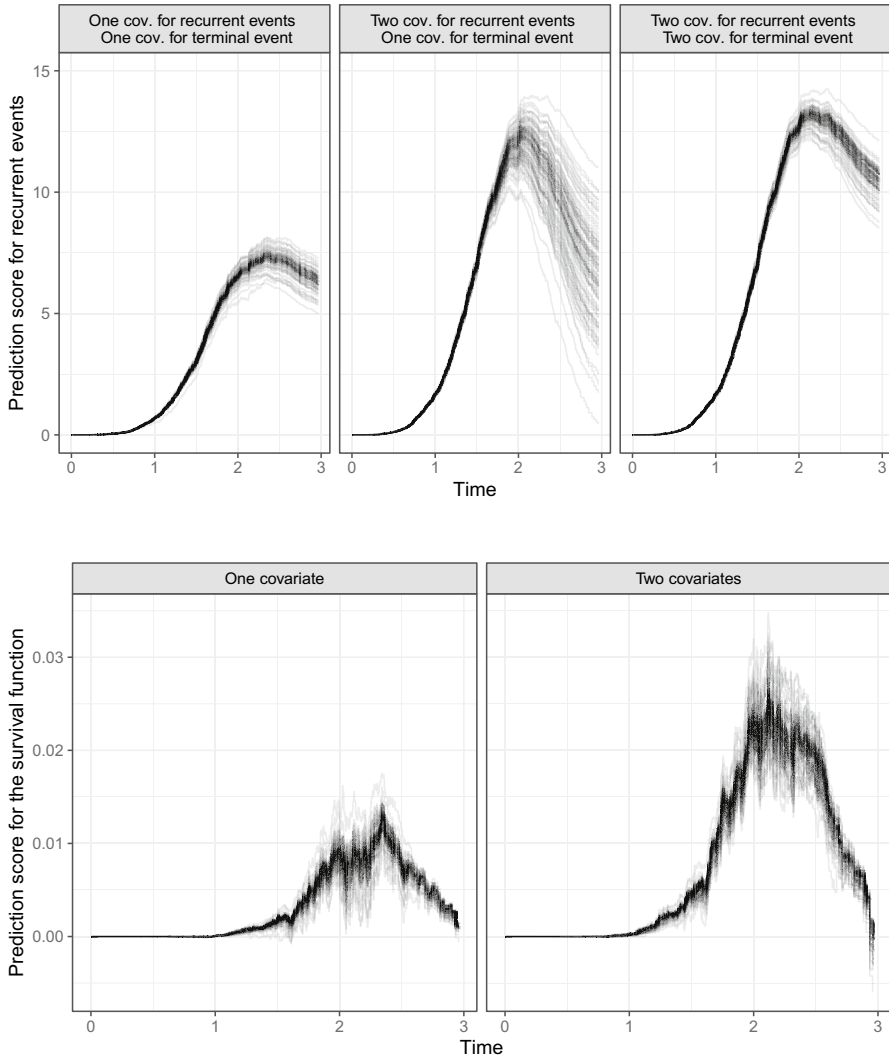


Fig. 2 Prediction scores for the recurrent events (top panels) and the survival function (bottom panels) using different models. The data were generated from a Cox model with two covariates ($n = 200$) for the recurrent event process and with the same two covariates for the terminal event. The expected cumulative number recurrent of events and the survival function of the terminal event were predicted using the Cox model with one or two covariates. The reference model uses no covariates and was estimated from the non-parametric estimator in Eq. (12) and from the Kaplan-Meier estimator in the top and bottom panels, respectively. The prediction scores are computed for 100 training samples of size $n_{\text{train}} = 800$ and a unique test sample of size $n_{\text{test}} = 1000$

with $\hat{\beta}^{T^*}$ is the estimated regression parameter from the Cox model for the terminal event and $\hat{\Lambda}_0^{T^*}$ its corresponding baseline estimator known as the Breslow estimator. If the terminal event model does not contain any covariates, then the Kaplan-Meier estimator is used instead. As previously, we use the score defined in Eq. (4) to evaluate the quality of prediction of a model where the reference model $\hat{\mu}_0$ is defined as

$$\hat{\mu}_0(t) = \int_0^t \hat{S}(u)d\hat{\Lambda}(u), \tag{12}$$

with \hat{S} the Kaplan-Meier estimator of the terminal event and $\hat{\Lambda}$ the Nelson-Aalen estimator of the recurrent event process. The same score is also used for the prediction of the survival function with the Kaplan-Meier estimator as the reference model. We consider four different regression models: a correctly specified model that includes the two covariates for both the recurrent event process and the terminal event in two Cox models, a model where the Gaussian covariate is missing for the Cox model of the terminal event (but the Cox model of the recurrent event is correctly specified) and a model where the Gaussian covariate is missing for both Cox models.

In Fig. 2, we simulated 100 training samples each of size 800 and we evaluated the prediction score on a unique test sample of size 1000. In the bottom panel, we see that including the two covariates in the survival model increases the prediction performance as compared to the model with only one covariate. Also, for both models, the gain in terms of prediction is more important for small time points and is reduced after time 2 approximately, as compared to the Kaplan-Meier estimator. This is due to the fact that the added predictive value of the models as compared to the reference decreases as we reach the tails and equals 0 at time 3 (97% of the terminal events will occur before time 3). This loss in terms of prediction efficiency

Table 1 Means and standard deviations (in bracket) over 500 simulations for the prediction score of the expected number of recurrent events in the presence of a terminal event

	$n_{\text{train}} = 100$			$n_{\text{train}} = 200$		
	$t = 1$	$t = 2$	$t = 2.9$	$t = 1$	$t = 2$	$t = 2.9$
One cov.-one cov	0.68 (0.05)	6.47 (0.43)	6.14 (0.83)	0.69 (0.03)	6.54 (0.24)	6.33 (0.45)
Two cov.-one cov	1.67 (0.07)	12 (0.92)	6.43 (2.71)	1.67 (0.04)	12.11 (0.6)	6.84 (1.77)
Two cov.-two cov	1.67 (0.07)	12.61 (0.59)	10.07 (1.18)	1.68 (0.04)	12.75 (0.3)	10.51 (0.59)
	$n_{\text{train}} = 400$			$n_{\text{train}} = 800$		
	$t = 1$	$t = 2$	$t = 2.9$	$t = 1$	$t = 2$	$t = 2.9$
One cov.-one cov	0.69 (0.01)	6.57 (0.14)	6.42 (0.28)	0.69 (0.01)	6.59 (0.08)	6.47 (0.17)
Two cov.-one cov	1.68 (0.02)	12.11 (0.39)	6.92 (1.22)	1.68 (0.01)	12.15 (0.27)	7.07 (0.85)
Two cov.-two cov	1.68 (0.02)	12.82 (0.17)	10.69 (0.35)	1.68 (0.01)	12.84 (0.09)	10.78 (0.23)

The results are presented for the same three models as in Fig. 2. The reference model uses no covariates and was estimated from the non-parametric estimator in Eq. (12). Large values indicate better predictive performances

Table 2 Means and standard deviations (in bracket) over 500 simulations for the prediction score of the survival function

$\times 10^4$	$n_{\text{train}} = 100$			$n_{\text{train}} = 200$		
	$t = 1$	$t = 2$	$t = 2.9$	$t = 1$	$t = 2$	$t = 2.9$
One cov	0.58 (0.55)	77.22 (34.77)	22.57 (8.29)	0.69 (0.25)	84.71 (15.22)	23.96 (3.81)
Two cov	1.71 (1.49)	214.84 (46.73)	32.3 (18.05)	2.07 (0.73)	227.83 (23.1)	38.3 (8.86)
	$n_{\text{train}} = 400$			$n_{\text{train}} = 800$		
	$t = 1$	$t = 2$	$t = 2.9$	$t = 1$	$t = 2$	$t = 2.9$
One cov	0.74 (0.17)	86.26 (8.65)	24.29 (2.81)	0.77 (0.12)	89.74 (5.44)	27.22 (3.02)
Two cov	2.2 (0.51)	231.51 (13.94)	40.49 (3.74)	2.28 (0.36)	238.41 (9.88)	46.23 (4.27)

The results were multiplied by 10^4 and are presented for the model with one or two covariates. The reference model uses no covariates and is estimated from the Kaplan-Meier estimator. Large values indicate better predictive performance

of the survival function for large time points impacts the prediction of the expected cumulative number of recurrent events. In the top panel, we see that adding the correct covariates in the Cox models of the survival function and of the recurrent event models increases the prediction performances. After time 2, the gain in terms of the performance prediction of the expected cumulative number of recurrent events slightly decreases due to the loss of efficiency in the prediction of the survival function. Tables 1 and 2 provide the mean score of the different models for the recurrent event process and the survival function, based on 500 training samples of size 100, 200, 400 and 800 and one single test sample of size 1000. We see the same trend as in Fig. 2 for all sample sizes. Clearly, increasing the sample size does not provide much gain in terms of average especially for small time points but it does reduce the variability of the predictors.

6 Real data analysis: the atrial fibrillation dataset

In this section, we analyse a dataset on patients with atrial fibrillation (AF). The aim is to compare different regression models for the prediction of the expected cumulative number of atrial fibrillation hospitalisations, using the prediction score developed in this work. Patients were enrolled from January 1st 2008 to December 1st 2012 in the “Atrial Fibrillation Survey-Copenhagen (ATLAS-CPH)” from both the in- and outpatient clinics at the Department of Cardiology at University Hospital Copenhagen, Hvidovre, Denmark. All patients were previously diagnosed with AF and were categorised at baseline, into either suffering from paroxysmal atrial fibrillation (PAF) or persistent atrial fibrillation (PeAF).

In total, 174 patients were enrolled with 50 PAF patients and 124 PeAF patients. Time is measured in days, with a mean follow-up duration of 1279 days. In terms of observed events, the patients experienced a total of 325 AF hospitalisations, with

305 AF hospitalisations in the PeAF group and 20 in the PAF group. A terminal event was defined as either progression to permanent AF or as the occurrence of death. In the dataset, 45 patients experienced a terminal event and the remaining 129 patients were censored. Finally, in top of the AF type, the dataset also includes 11 additional variables: gender, age, alcohol consumption (with two levels 0–5 and > 5), tobacco consumption (with three levels “never smoked”, “ex-smoker”, “current smoker”), presence of hypertension, heart failure, valvular heart disease, ischemic heart disease, diabetes, COPD, antiarrhythmic medication. The data are presented in great details in Schroder et al. (2019). Note also that the data are fully available from the Plos One website.

In Schroder et al. (2019), the authors analysed the data using a multi-state approach with four possible states: no experience of recurrent events yet, 1 recurrent event, 2 or more recurrent events and the absorbing state for the terminal event. The transition intensities were assumed to be proportional with each other using a Cox model, where the number of previous recurrent events was included in the model. Those types of multi-state models with terminal event are described for instance in Cook and Lawless (2007) (see Section 6.6.4 of their book). Those analyses showed a high significant effect of the AF type, the number of previous recurrent events (p -values < 10^{-4}) and of age (p -value = 0.0253) for the risk of future AF hospitalisations. The effect of diabetes had a p -value equal to 0.0955. All other variables were assessed as non significant (p -values > 0.2). A Cox model was also implemented for the terminal event using the multi-state approach (that is including the effect of previous AF hospitalisations through a proportional effect) with all variables. Only the age variable was significant (in the multivariate Cox model, the hazard rate was equal to 1.05 and the p -value was equal to 0.0016). In this previous work, the authors then decided to only include the covariates AF type and age, with a proportional effect of the number of previous AF hospitalisations for the modelisation of the recurrent event process. For the terminal event model, they only included the age variable. Based on those models, it is then possible to produce predictions for the expected cumulative number of future AF hospitalisations, for a given patient based on his/her characteristics. Using the prediction score developed in this paper, we will compare the performance of the model used in Schroder et al. (2019) with several other possible models. Since diabetes is a known risk factor for AF, we will also consider models with this variable, along with AF type and age. As in the simulation section, the prediction score will be computed from a training and a test samples, but this time using 10-fold cross validation, that is one tenth of the observations are used for the test sample and the remaining observations are used for the model estimations and the procedure is repeated and averaged ten times. The aim of this section is to compare different prediction models for the recurrent hospitalisations. Since the data contain a terminal event, a model for the survival function of this terminal event needs first to be proposed. We use a Cox model with age as the only covariate. We refer the reader to the Supplementary Information for a more detailed comparison between different regression models for the prediction of the terminal event.

The prediction performance for the recurrent event process is now investigated. In the definition of our prediction criterion (see Eq. (1)), we will always estimate the

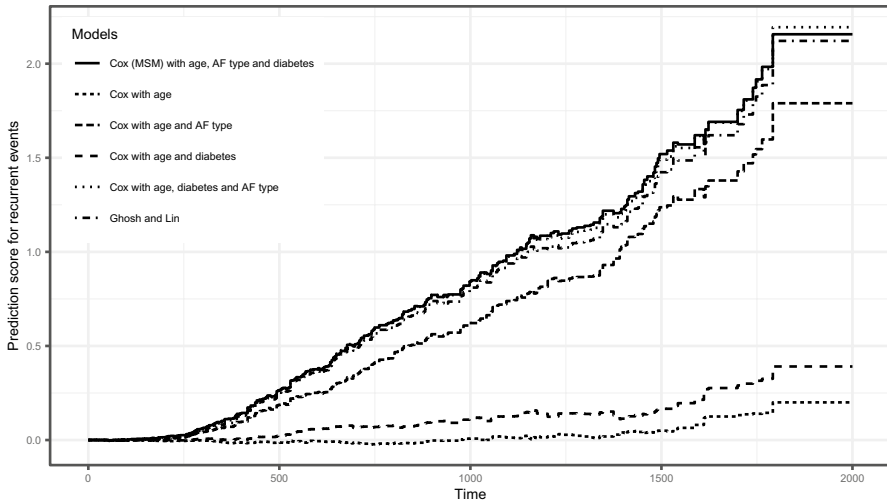


Fig. 3 Prediction scores for the expected cumulative number of recurrent events in the atrial fibrillation dataset. With the non-parametric estimator (see Eq. (12)) as the reference, six different models are compared. All the models use the Cox model with age as the unique covariate for the estimation of the survival function. For ease of visualisation, we describe the six models (for the recurrent events) in increasing order of their scores at time $t = 2000$: the Cox model with covariate age (score = 0.200), the Cox model with covariates age and diabetes (score = 0.391), the Cox model with covariate age and AF type (score = 1.790), the Ghosh and Lin model (score = 2.121), the multi-state (MSM) Cox model with covariates age, AF type and diabetes (score = 2.156) and the Cox model with covariates age, AF type and diabetes (score = 2.194). The MSM Cox model assumes that the transition intensities from 0 event to 1 and to “one event or more” to a new event are proportional

censoring distribution from the survival random forest, based on the `rfsrc` package (see Ishwaran et al. 2008). We consider the following models for the recurrent events:

- Four multivariate Cox models based on the age, diabetes and AF type variables,
- The Cox model stratified with respect to AF type and adjusted for age and diabetes,
- The Aalen model with covariates age, diabetes and AF type,
- The Cox multi-state model with covariates age, diabetes and AF type,
- The Cox multi-state model stratified with respect to AF type and adjusted for age and diabetes,
- The Ghosh and Lin model (Ghosh and Lin 2002) with covariates age, diabetes and AF type.

The reference model is taken as the non-parametric estimator (see Eq. (12)) and the score is again computed using formula (4). The results are displayed in Table 3 and Fig. 3 (in the figure only six different models are represented). The Ghosh and Lin model has been implemented with the `recreg` function from the `metS` package. We observe that the Cox model with the age variable has a poor predictive performance. From this model, adding the diabetes or AF type improves the model, with a

Table 3 Means and 80% intervals (in curly bracket) over 10-folds cross validation for the prediction score of the expected cumulative number of recurrent events in the atrial fibrillation dataset

	$t = 1000$	$t = 1500$	$t = 2000$
Cox with age	0.008 [-0.538, 0.323]	0.050 [-0.790, 0.603]	0.200 [-1.114, 1.407]
Cox with age and diabetes	0.118 [-0.147, 0.560]	0.166 [-0.726, 0.948]	0.391 [-0.807, 2.480]
Cox with age and AF type	0.622 [-0.205, 1.201]	1.237 [0.741, 1.875]	1.790 [0.480, 5.149]
Cox with age, AF type and diabetes	0.847 [0.042, 1.557]	1.490 [0.168, 2.390]	2.194 [0.222, 6.002]
Aalen with age, AF type and diabetes	0.818 [0.381, 1.352]	1.541 [0.556, 2.540]	2.227 [0.182, 5.011]
Cox with age, diabetes and strata(AF type)	0.843 [0.035, 1.559]	1.491 [0.120, 2.412]	2.184 [0.222, 6.000]
Cox (MSM) with age, AF type and diabetes	0.849 [0.123, 1.401]	1.520 [0.221, 2.319]	2.156 [0.306, 5.027]
Cox (MSM/strata) with age, AF type and diabetes	0.853 [0.142, 1.442]	1.526 [0.269, 2.426]	2.160 [0.267, 5.020]
Ghosh and Lin	0.810 [0.015, 1.542]	1.423 [0.177, 2.332]	2.121 [0.150, 6.195]

With the non-parametric estimator (see Eq. (12)) as the reference, nine different models (for the recurrent events) are compared at three different time points: five Cox and Aalen models with covariates age, diabetes and AF type, one Cox model stratified with respect to AF type, two multi-state Cox models, Cox (MSM) and Cox (MSM/strata) and the Ghosh and Lin model. The difference between the two MSM models is that the first one assumes the baseline transition intensities from 0 event to 1 and to “one event or more” to a new event to be proportional while the last one uses two different baselines functions. All the models use the Cox model with age as the unique covariate for the estimation of the survival function

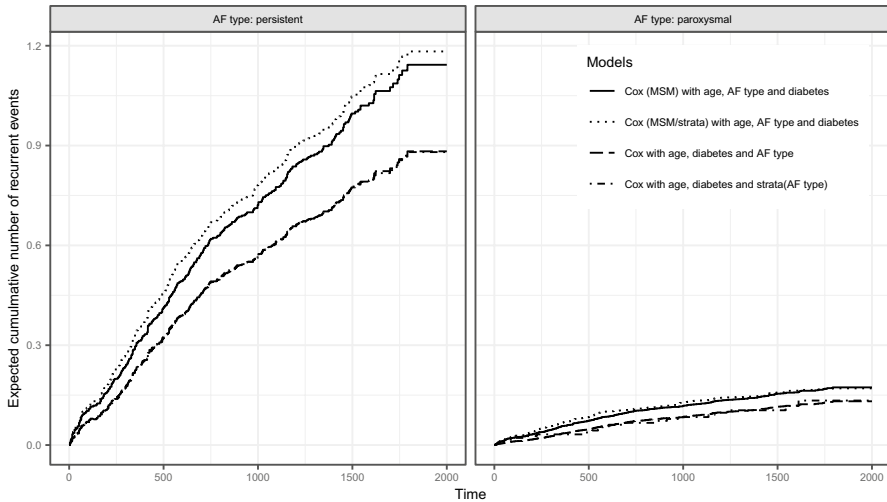


Fig. 4 Expected cumulative number of recurrent events predicted from four different models for two 60 year old patients with diabetes, one with persistent AF and the other with paroxysmal AF. Those models had similar prediction performances (see Table 3 and Fig. 3)

much bigger gain with the AF type variable. Further, combining all three variables in the same model provides a substantial gain with respect to all previous models. On the other hand, the two multi-state models provide only a minor improvement of the predictions with a slight advantage for the stratified model. The Ghosh and Lin model shows a good prediction performance but it is slightly less performant than the Cox model with all three covariates. On the overall, the Aalen model has a very good prediction performance. Finally, the predictions for some of these models on the expected cumulative number of AF hospitalisations are displayed in Fig. 4. In this figure, the predictions are made for two 60 year old patients with diabetes, one with persistent AF and the other with paroxysmal AF. While the different models do not vary much in their predictions for the paroxysmal AF patient, they offer different results for the persistent AF patient. According to the results from Table 3 and Fig. 3, the Aalen model and the Cox multi-state model stratified with respect to AF type have the greatest prediction performances and therefore should be chosen. After 1 500 days after AF diagnosis, the latter model predicts an expected number of AF hospitalisations equal to 1.03 approximately. On the other hand, if one uses the multi-state Cox model, the prediction is equal to 0.98, if one uses the multivariate Cox model with all three variables, the prediction is equal to 0.76.

7 Conclusion

In this work a new prediction criterion was proposed in the context of recurrent event data. The criterion evaluates the prediction performance of the expected cumulative number of recurrent events, while taking into account censoring and a possible terminal event. We showed that it can be decomposed into an inseparability and imprecision terms in the same manner as in Graf et al. (1999). Moreover, the simulations revealed that the inseparability term was largely dominant in the decomposition. As a result, we recommend to use the prediction score defined in Eq. (4), as the difference between the prediction criterion of a given model and of a reference model, typically a model that does make use of the covariates, such that the score provides the absolute gain from the covariates in the proposed model. An alternative score could be derived by computing the relative gain as proposed in Steyerberg et al. (2010). This produces a score that ranges from 0 to 100% and shares similarities with the Pearson's R^2 statistic. However, care should be taken with such a score, due to the fact that we normalise with respect to the prediction criterion of the reference model, which itself can be decomposed into imprecision and inseparability. This criterion could therefore be misleading due to the magnitude of the inseparability term which is unknown in practice.

The proposed prediction criterion is simple to compute and has the advantage to include all the recurrent events. As a result, it can be seen as an overall performance measure that provides information about the global predictive ability of the proposed model in a recurrent event context. Nevertheless, it would be possible to modify the criterion if one is interested into evaluating the performance of a model to only predict further recurrent events after a fixed number of events have already been experienced by a patient. This would amount to conditioning on a given number of experienced recurrent events in a multi-state framework. This type of criterion would be similar to the one developed in Schoop et al. (2011) which conditions on being alive up to a time t^* and evaluate the prediction of the model for a time $s > t^*$. Another improvement would be to allow for frailty models in the manner of Van Oirbeek and Lesaffre (2016). A marginal score that integrates the frailty variable could be derived. Such a score would provide an overall evaluation of the frailty model and would be a natural extension of the score proposed in this paper. Alternatively, a conditional score could be proposed for the conditional (with respect to the frailty) expected cumulative number of recurrent events. More work is needed to develop these two scores. Another extension of interest is the construction of confidence intervals for the prediction score. This would allow to ascertain the sampling variability of the prediction score without performing V-fold cross-validation, as done in the application of this work. In Bradley et al. (2008), the authors have studied the asymptotic distribution of the Brier score when the data are fully observed. This result could be first extended to right-censored data and then to the recurrent events framework. This is left to future research.

8 Appendix: proofs of the convergence of the prediction criterion for the expected cumulative number of recurrent events under the two scenarios

In the proof of Proposition 1, we need to verify the key equality from Eq. (3). This result depends on the modelling assumptions and has already been proved in all three different scenarios, see Section 1 of Supplementary Information, Sect. 2.2 of the main manuscript and Section 2 of Supplementary Information for the right-censoring case with no terminal event, the terminal event case, and the dependence on prior counts case, respectively. In the proof of Proposition 2, we also need to have $\mathbb{E}[\mu^*(\tau | \bar{X}(\tau))] < \infty$ which also depends on the different modelling assumptions made under each scenario.

8.1 Proof of Proposition 1

In all three scenarios, we directly have:

$$\begin{aligned} \text{MSE}(t, \mu) &= \mathbb{E}\left[\left(\mu(t | \bar{X}(t)) - \mu^*(t | \bar{X}(t))\right)^2\right] + \mathbb{E}\left[\left(\int_0^t \frac{dN(u)}{G_c(u | \bar{X}(u))} - \mu^*(t | \bar{X}(t))\right)^2\right] \\ &\quad + 2\mathbb{E}\left[\left(\int_0^t \frac{dN(u)}{G_c(u | \bar{X}(u))} - \mu^*(t | \bar{X}(t))\right)\left(\mu^*(t | \bar{X}(t)) - \mu(t | \bar{X}(t))\right)\right]. \end{aligned}$$

Using the fact that $\mathbb{E}[\int_0^t dN(u)/(G_c(u | \bar{X}(u))) | \bar{X}(t)] = \mu^*(t | \bar{X}(t))$, we conclude that

$$\text{MSE}(t, \mu) = \mathbb{E}\left[\left(\mu(t | \bar{X}(t)) - \mu^*(t | \bar{X}(t))\right)^2\right] + A(t),$$

where

$$A(t) = \mathbb{E}\left[\left(\int_0^t \frac{dN(u)}{G_c(u | \bar{X}(u))}\right)^2\right] - \mathbb{E}\left[\left(\mu^*(t | \bar{X}(t))\right)^2\right]. \tag{13}$$

Now, using the remarkable identity $a^2 - b^2 = (a - b)(a + b)$ and observing that $\int_0^t dN(u)/(G_c(u | \bar{X}(u))) = \sum_{\text{ev} \leq t} \{1/(G_c(u | X(\text{ev})))\}$ either equals 0 if no observed recurrent events occurred before time t or is greater than 1 if at least one recurrent event occurred before time t , we conclude that

$$\begin{aligned} &\left(\int_0^t \frac{dN(u)}{G_c(u | \bar{X}(u))} - \mu^*(t | \bar{X}(t))\right)\left(\int_0^t \frac{dN(u)}{G_c(u | \bar{X}(u))} + \mu^*(t | \bar{X}(t))\right) \\ &\geq \left(\int_0^t \frac{dN(u)}{G_c(u | \bar{X}(u))} - \mu^*(t | \bar{X}(t))\right)\mu^*(t | \bar{X}(t)), \end{aligned}$$

almost surely. Taking the expectation on both sides proves $A(t) \geq 0$.

8.2 Proof of Proposition 2

We start by proving that $\mathbb{E}[\mu^*(\tau \mid \bar{X}(\tau))] < \infty$ in the presence of a terminal event (the scenario without terminal event follows from the same arguments). We have for all $t \in [0, \tau] : \mathbb{P}[C \geq t \mid \bar{X}(t)] \geq \mathbb{P}[T \geq t \mid \bar{X}(t)] \geq c$, from Assumption 2. From the same assumption, $N(\tau)$ is almost surely bounded by a constant. As a consequence,

$$\mu^*(\tau \mid \bar{X}(\tau)) = \int_0^\tau \frac{\mathbb{E}[dN(t) \mid \bar{X}(t)]}{G_c(t \mid \bar{X}(t))}$$

is almost surely bounded, where the equality has been proved in Sect. 2.2. In the dependence on prior counts case, we have for all $t \in [0, \tau] : \mathbb{P}[C \geq t \mid \bar{X}(t)] \geq \mathbb{P}[T \geq t \mid \bar{X}(t)] = \sum_{i=1}^{L+1} \mathbb{P}[T \geq t, N(t-) = i - 1 \mid \bar{X}(t)] \geq (L + 1)c > 0$, where the two last bounds come from Assumption MSM in the Supplementary Information. From the same assumption, $N(\tau)$ is almost surely bounded by a constant. As a consequence,

$$\mu^*(\tau \mid \bar{X}(\tau)) = \int_0^\tau \frac{\mathbb{E}[dN(t) \mid \bar{X}(t)]}{G_c(t \mid \bar{X}(t))}$$

is almost surely bounded, where the equality has been proved in the Supplementary Information. The rest of the proof of Proposition 2 is identical in all three scenarios.

We first note $F_{X(t)}(x) = \mathbb{P}[X(t) \leq x]$, we let $\mathcal{X}_{u,v}$ denote the support of the joint distribution $(X(u), X(v))$ and we note $F_{X(u),X(v)}(x, y) = \mathbb{P}[X(u) \leq x, X(v) \leq y]$. We then introduce the quantity

$$\begin{aligned} \xi(t) &= \int_{0 \leq u, v \leq t} \int_{\mathcal{X}_{u,v}} \frac{\mathbb{E}[dN(u)dN(v) \mid X(u) = x, X(v) = y]}{\hat{G}_c(u \mid x)\hat{G}_c(v \mid y)} dF_{X(u),X(v)}(x, y) \\ &\quad - 2 \int_{\mathcal{X}_t} \hat{\mu}(t \mid x) \mu^*(t \mid x) dF_{X(t)}(x) \\ &\quad + \int_{\mathcal{X}_t} (\hat{\mu}(t \mid x))^2 dF_{X(t)}(x) =: \xi_1(t) + \xi_2(t) + \xi_3(t). \end{aligned}$$

Write:

$$\begin{aligned} \left| \widehat{\text{MSE}}^1(t, \hat{\mu}) - \text{MSE}^1(t, \mu) \right| &\leq \left| \xi(t) - \mathbb{E} \left[\left(\int_0^t \frac{dN(u)}{G_c(u \mid \bar{X}(u))} - \mu(t \mid \bar{X}(t)) \right)^2 \right] \right| \\ &\quad + \left| \frac{1}{n} \sum_{i=1}^n \left(\int_0^t \frac{dN_i(u)}{\hat{G}_c(u \mid \bar{X}_i(u))} - \hat{\mu}(t \mid X_i(t)) \right)^2 - \xi(t) \right| \\ &\leq C(t) + D(t). \end{aligned}$$

By decomposing the square term into three other terms, we bound $C(t)$ in the following way: $C(t) \leq |C_1(t)| + |C_2(t)| + |C_3(t)|$ with

$$\begin{aligned}
 C_1(t) &= \int_{0 \leq u, v \leq t} \int_{\mathcal{X}_{u,v}} \frac{G_c(u | x)G_c(v | y) - \hat{G}_c(u | x)\hat{G}_c(v | y)}{\hat{G}_c(u | x)\hat{G}_c(v | y)G_c(u | x)G_c(v | y)} \\
 &\quad \mathbb{E}[dN(u)dN(v) | X(u) = x, X(v) = y]dF_{X(u),X(v)}(x, y), \\
 C_2(t) &= -2 \int_{\mathcal{X}_t} (\hat{\mu}(t | x) - \mu(t | x))\mu^*(t | x)dF_{X(t)}(x), \\
 C_3(t) &= \int_{\mathcal{X}_t} \left((\hat{\mu}(t | x))^2 - (\mu(t | x))^2 \right) dF_{X(t)}(x).
 \end{aligned}$$

For $C_1(t)$ we have

$$\begin{aligned}
 &G_c(u | x)G_c(v | y) - \hat{G}_c(u | x)\hat{G}_c(v | y) \\
 &= (\hat{G}(u- | x) - G(u- | x)) + (\hat{G}(v- | y) - G(v- | y)) \\
 &\quad + G(u- | x)(G(v- | y) - \hat{G}(v- | y)) + \hat{G}(v- | y)(G(u- | x) - \hat{G}(u- | x)),
 \end{aligned}$$

and we can deal with all four terms in the same fashion. For instance, for the first term,

$$\begin{aligned}
 &\int_{0 \leq u, v \leq t} \int_{\mathcal{X}_{u,v}} \frac{(\hat{G}(u- | x) - G(u- | x))\mathbb{E}[dN(u)dN(v) | X(u) = x, X(v) = y]}{\hat{G}_c(u | x)\hat{G}_c(v | y)G_c(u | x)G_c(v | y)} dF_{X(u),X(v)}(x, y) \\
 &\leq \int_0^t \int_{\mathcal{X}_u} \frac{|\hat{G}(u- | x) - G(u- | x)|\mathbb{E}[dN(u) | X(u) = x]}{\hat{G}_c(u | x)G_c(u | x)} dF_{X(u)}(x),
 \end{aligned}$$

using the fact that $\int_0^t dN(v)/(\hat{G}_c(v | y)G_c(v | y))$ is bounded. Then, since $\int_0^t \mathbb{E}[dN(u)/(1 - G(u- | X(u))) | X(u) = x] = \mu^*(t | \bar{X}(t))$ and $\hat{G}_c(u | x)^{-1}$ is asymptotically uniformly bounded, we conclude that $|C_1(t)|$ tends toward 0 in probability using the uniform consistency of the censoring estimator.

For $C_2(t)$ we use the consistency of $\hat{\mu}$ and the fact that $\mathbb{E}[\mu^*(t | \bar{X}(t))]$ is finite to prove that $|C_2(t)|$ tends towards 0 in probability.

For $C_3(t)$, we directly write $(\hat{\mu}(t | x))^2 - (\mu(t | x))^2 = (\hat{\mu}(t | x) - \mu(t | x))(\hat{\mu}(t | x) + \mu(t | x))$ and we use the fact that $\mu(t | x)$ is bounded and the consistency of $\hat{\mu}$ to prove that $|C_3(t)|$ tends towards 0 in probability.

Similarly to $C(t)$ we obtain the following bound: $D(t) \leq |D_1(t)| + |D_2(t)| + |D_3(t)|$ with

$$\begin{aligned}
 D_1(t) &= \frac{1}{n} \sum_{i=1}^n \int_{0 \leq u, v \leq t} \frac{dN_i(u)dN_i(v)}{\hat{G}_c(u | \bar{X}_i(u))\hat{G}_c(v | \bar{X}_i(v))} - \xi_1(t), \\
 D_2(t) &= -\frac{2}{n} \sum_{i=1}^n \int_0^t \frac{dN_i(u)}{\hat{G}_c(u | \bar{X}_i(u))} \hat{\mu}(t | X_i(t)) - \xi_2(t), \\
 D_3(t) &= \frac{1}{n} \sum_{i=1}^n \left(\hat{\mu}(t | \bar{X}_i(t)) \right)^2 - \xi_3(t).
 \end{aligned}$$

We now use the bound $|D_1(t)| \leq |D_{1,1}(t)| + |D_{1,2}(t)| + |D_{1,3}(t)|$ with

$$\begin{aligned}
 D_{1,1}(t) &= \frac{1}{n} \sum_{i=1}^n \int_{0 \leq u, v \leq t} \frac{dN_i(u)dN_i(v)}{G_c(u | \bar{X}_i(u))G_c(v | \bar{X}_i(v))} \\
 &\quad - \int_{0 \leq u, v \leq t} \int_{\mathcal{X}_{u,v}} \frac{\mathbb{E}[dN(u)dN(v) | X(u) = x, X(v) = y]}{G_c(u | x)G_c(v | y)} dF_{X(u),X(v)}(x, y), \\
 D_{1,2}(t) &= \frac{1}{n} \sum_{i=1}^n \int_{0 \leq u, v \leq t} \chi(u, v, X_i(u), X_i(v)) dN_i(u)dN_i(v) \\
 D_{1,3}(t) &= - \int_{0 \leq u, v \leq t} \int_{\mathcal{X}_{u,v}} \chi(u, v, x, y) \mathbb{E}[dN(u)dN(v) | X(u) = x, X(v) = y] dF_{X(u),X(v)}(x, y)
 \end{aligned}$$

and

$$\begin{aligned}
 \chi(u, v, x, y) &= \{ (\hat{G}(u- | x) - G(u- | x)) + (\hat{G}(v- | y) - G(v- | y)) \\
 &\quad + (\hat{G}(v- | y) - G(v- | y)) + G(u- | x)(G(v- | y) - \hat{G}(v- | y)) \\
 &\quad + \hat{G}(v- | x)(G(u- | y) - \hat{G}(u- | y)) \} \\
 &\quad \times \frac{1}{\hat{G}_c(u | x)\hat{G}_c(v | y)G_c(u- | x)G_c(v | y)}.
 \end{aligned}$$

The term $|D_{1,1}(t)|$ converges towards 0 in probability from the strong law of large numbers. The term $|D_{1,2}(t)|$ is bounded by

$$\sup_{u,v,x,y} |\chi(u, v, x, y)| \frac{1}{n} \sum_{i=1}^n \int_{0 \leq u < v \leq t} dN_i(u)dN_i(v),$$

$\sup_{u,v,x,y} |\chi(u, v, x, y)|$ converges towards 0 from the uniform consistency of \hat{G} while the other term converges towards a bounded quantity from the law of large numbers. The same argument applies to $|D_{1,3}(t)|$ which also converges towards 0 in probability.

For $D_2(t)$ we write $|D_2(t)| \leq |D_{2,1}(t)| + |D_{2,2}(t)| + |D_{2,3}(t)| + |D_{2,4}(t)|$ with

$$\begin{aligned}
 D_{2,1}(t) &= -\frac{2}{n} \sum_{i=1}^n \int_0^t \frac{dN_i(u)}{G_c(u | \bar{X}_i(u))} \mu(t | \bar{X}_i(t)) + 2 \int_{\mathcal{X}_i} \mu(t | x) \mu^*(t | x) dF_{X(t)}(x), \\
 D_{2,2}(t) &= \frac{2}{n} \sum_{i=1}^n \int_0^t \frac{dN_i(u)}{G_c(u | \bar{X}_i(u))} (\mu(t | \bar{X}_i(t)) - \hat{\mu}(t | \bar{X}_i(t))) \\
 D_{2,3}(t) &= 2 \int_{\mathcal{X}_i} (\hat{\mu}(t | x) - \mu(t | x)) \mu^*(t | x) dF_{X(t)}(x), \\
 D_{2,4}(t) &= \frac{2}{n} \sum_{i=1}^n \int_0^t \frac{(G(u- | \bar{X}_i(u)) - \hat{G}(u- | \bar{X}_i(u)))dN_i(u)}{G_c(u | \bar{X}_i(u))\hat{G}_c(u- | \bar{X}_i(u))} \hat{\mu}(t | \bar{X}_i(t)).
 \end{aligned}$$

The $D_{2,1}(t)$ term converges towards 0 in probability from the law of large numbers. For $D_{2,2}(t)$, $D_{2,3}(t)$ and $D_{2,4}(t)$ we use the consistency of $\hat{\mu}$, the convergence in probability of $\sum_i \int_0^t dN_i(u)G_c(u | X_i(u))/n$, the boundedness of $\mathbb{E}[\mu^*(t | \bar{X}(t))]$, the uniform consistency of \hat{G} and the asymptotic boundedness of $\hat{\mu}$ and $\hat{G}_c(u | x)^{-1}$ to prove that all three terms converge towards 0 in probability.

Finally, for $D_3(t)$, we write

$$\begin{aligned}
 D_3(t) &= \frac{1}{n} \sum_{i=1}^n \left(\mu(t \mid \bar{X}_i(t)) \right)^2 - \int_{\mathcal{X}_t} \left(\mu(t \mid x) \right)^2 dF_{X(t)}(x) \\
 &+ \frac{1}{n} \sum_{i=1}^n \left(\left(\hat{\mu}(t \mid \bar{X}_i(t)) \right)^2 - \left(\mu(t \mid \bar{X}_i(t)) \right)^2 \right) \\
 &+ \int_{\mathcal{X}_t} \left(\left(\hat{\mu}(t \mid x) \right)^2 - \left(\mu(t \mid x) \right)^2 \right) dF_{X(t)}(x).
 \end{aligned}$$

Each of the three terms converges towards 0 in probability using the law of large numbers for the first term and the uniform consistency of $\hat{\mu}$ for the other two.

8.3 Proof of Proposition 3

First, note that the Brier score can be written in the following way:

$$\text{MSE}^{\text{Brier}}(t, \pi) = \mathbb{E}[S(t \mid X)] - 2\mathbb{E}[S(t \mid X)\pi(t \mid X)] + \mathbb{E}[(\pi(t \mid X))^2].$$

We now study, our prediction score $\text{MSE}'(t, \pi)$. Using standard martingale properties (see for instance Andersen et al. (1993)), we directly have that $\mathbb{E}[dN(t \mid X)] = H(t \mid X)\lambda^*(t \mid X)dt$, where $H(t \mid X) = \mathbb{P}[T > t \mid X] = S(t \mid X)G_c(t \mid X)$ under independent censoring and λ^* is the hazard rate of T^* . As a consequence,

$$\mathbb{E} \left[\int_0^t \frac{dN(u)}{G_c(u \mid X)} \mid X \right] = \int_0^t S(u \mid X)\lambda^*(u \mid X)du = 1 - S(t \mid X), \tag{14}$$

since $S(u \mid X)\lambda^*(u \mid X)$ is equal to the conditional density function of T^* . Also, it is important to notice that

$$\mathbb{E} \left[\left(\int_0^t \frac{dN(u)}{G_c(u \mid X)} \right)^2 \right] = \mathbb{E} \left[\int_0^t \frac{dN(u)}{(G_c(u \mid X))^2} \right] = \mathbb{E} \left[\int_0^t \frac{S(u \mid X)}{G_c(u \mid X)} \lambda^*(u \mid X)du \right],$$

where the first equality is due to the fact that N can only jump once and thus $(\int_0^t dN(u)/(G_c(u \mid X))^2)$ is simply equal to $\Delta I(T \leq t)/(G_c(T \mid X))^2$. Now,

$$\begin{aligned}
 \text{MSE}'(t, \pi) &= \mathbb{E} \left[\left(1 - \int_0^t \frac{dN(u)}{G_c(u \mid X)} \right)^2 \right] - 2\mathbb{E} \left[\left(1 - \int_0^t \frac{dN(u)}{G_c(u \mid X)} \right) \pi(t \mid X) \right] \\
 &+ \mathbb{E}[(\pi(t \mid X))^2] \\
 &= 1 - 2\mathbb{E}[(1 - S(t \mid X))] + \mathbb{E} \left[\left(\int_0^t \frac{dN(u)}{G_c(u \mid X)} \right)^2 \right] - 2\mathbb{E}[S(t \mid X)\pi(t \mid X)] \\
 &+ \mathbb{E}[(\pi(t \mid X))^2] \\
 &= \text{MSE}^{\text{Brier}}(t, \pi) + B(t),
 \end{aligned}$$

with

$$B(t) = -\mathbb{E}[1 - S(t | X)] + \mathbb{E} \left[\int_0^t \frac{S(u | X)}{G_c(u | X)} \lambda^*(u | X) du \right].$$

Now, using Eq. (14), we can rewrite $B(t)$ in the following way:

$$\begin{aligned} B(t) &= -\mathbb{E} \left[\int_0^t S(u | X) \lambda^*(u | X) du \right] + \mathbb{E} \left[\int_0^t \frac{S(u | X)}{G_c(u | X)} \lambda^*(u | X) du \right] \\ &= \mathbb{E} \left[\int_0^t \frac{G(u-)}{G_c(u | X)} S(u | X) \lambda^*(u | X) du \right]. \end{aligned}$$

This shows that $B(t) \geq 0$ and that this quantity does not depend on π .

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1007/s10985-023-09610-x>.

Acknowledgements We thank the reviewers for their constructive criticisms and comments that have helped improve the paper.

References

- Andersen PK, Angst J, Ravn H (2019) Modeling marginal features in studies of recurrent events in the presence of a terminal event. *Lifetime Data Anal* 25:681–695
- Andersen PK, Borgan Ø, Gill RD, Keiding N (1993) *Statistical models based on counting processes*. Springer series in statistics. Springer-Verlag, New York
- Andersen PK, Gill RD (1982) Cox's regression model for counting processes: a large sample study. *Ann Stat* 10(4):1100–1120
- Bouaziz O, Geffray S, Lopez O (2015) Semiparametric inference for the recurrent events process by means of a single-index model. *Statistics* 49:361–385
- Bouaziz O, Lopez O (2010) Conditional density estimation in a censored single-index regression model. *Bernoulli* 16:514–542
- Bradley AA, Schwartz SS, Hashino T (2008) Sampling uncertainty and confidence intervals for the brier score and brier skill score. *Weather Forecast* 23:992–1006
- Cook RJ, Lawless J (2007) *The statistical analysis of recurrent events*. Springer Science & Business Media, New-York, USA
- Cook RJ, Lawless JF (1997) Marginal analysis of recurrent events and a terminating event. *Stat Med* 16:911–924
- Cox DR (1972) Regression models and life-tables. *J R Stat Soc Ser B (Methodol)* 34:187–202
- Dabrowska DM (1989) Uniform consistency of the kernel conditional Kaplan-Meier estimate. *Ann. Stat* 17(3):1157–1167
- Gerds TA, Kattan MW, Schumacher M, Yu C (2013) Estimating a time-dependent concordance index for survival prediction models with covariate dependent censoring. *Stat Med* 32:2173–2184
- Gerds TA, Schumacher M (2006) Consistent estimation of the expected brier score in general survival models with right-censored event times. *Biom J* 48:1029–1040
- Ghosh D, Lin D (2003) Semiparametric analysis of recurrent events data in the presence of dependent censoring. *Biometrics* 59:877–885
- Ghosh D, Lin DY (2002) Marginal regression models for recurrent and terminal events. *Stat Sin* 12(3):663–688
- Graf E, Schmoor C, Sauerbrei W, Schumacher M (1999) Assessment and comparison of prognostic classification schemes for survival data. *Stat Med* 18:2529–2545
- Harrell FE Jr, Lee KL, Mark DB (1996) Multivariable prognostic models: issues in developing models, evaluating assumptions and adequacy, and measuring and reducing errors. *Stat Med* 15:361–387

- Heagerty PJ, Zheng Y (2005) Survival model predictive accuracy and ROC curves. *Biometrics* 61:92–105
- Hougaard P, Hougaard P (2000) Analysis of multivariate survival data, vol 564. Springer, New York, USA
- Ishwaran H, Kogalur UB, Blackstone EH, Lauer MS (2008) Random survival forests. *Ann Appl Stat* 2:841–860
- Kalbfleisch JD, Prentice RL (2002) The statistical analysis of failure time data, 2nd edn. Wiley series in probability and statistics. Wiley-Interscience (John Wiley & Sons), Hoboken
- Lin D, Wei L, Ying Z (1998) Accelerated failure time models for counting processes. *Biometrika* 85:605–618
- Lin DY, Wei L-J, Yang I, Ying Z (2000) Semiparametric regression for the mean and rate functions of recurrent events. *J R Stat Soc Ser B (Stat Methodol)* 62:711–730
- Liu L, Wolfe RA, Huang X (2004) Shared frailty models for recurrent events and a terminal event. *Biometrics* 60:747–756
- Prentice RL, Williams BJ, Peterson AV (1981) On the regression analysis of multivariate failure time data. *Biometrika* 68:373–379
- Rondeau V, Mathoulin-Pelissier S, Jacqmin-Gadda H, Brouste V, Soubeyran P (2007) Joint frailty models for recurring events and death using maximum penalized likelihood estimation: application on cancer events. *Biostatistics* 8:708–721
- Scheike TH (2002) The additive nonparametric and semiparametric Aalen model as the rate function for a counting process. *Lifetime Data Anal* 8:247–262
- Schoop R, Schumacher M, Graf E (2011) Measures of prediction error for survival data with longitudinal covariates. *Biom J* 53:275–293
- Schroder J, Bouaziz O, Agner BR, Martinussen T, Madsen PL, Li D, Dixen U (2019) Recurrent event survival analysis predicts future risk of hospitalization in patients with paroxysmal and persistent atrial fibrillation. *PLoS One* 14:e0217983
- Steyerberg EW, Vickers AJ, Cook NR, Gerds T, Gonen M, Obuchowski N, Pencina MJ, Kattan MW (2010) Assessing the performance of prediction models: a framework for some traditional and novel measures. *Epidemiology (Cambridge, Mass)* 21:128
- Van Oirbeek R, Lesaffre E (2016) Exploring the clustering effect of the frailty survival model by means of the brier score. *Commun Stat-Simul Comput* 45:3294–3306

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.