
Survival Permanental Processes for Survival Analysis with Time-Varying Covariates

Hideaki Kim

NTT Human Informatics Laboratories
NTT Corporation
hideaki.kin@ntt.com

Abstract

Survival or time-to-event data with time-varying covariates are common in practice, and exploring the non-stationarity in covariates is essential to accurately analyzing the nonlinear dependence of time-to-event outcomes on covariates. Traditional survival analysis methods such as Cox proportional hazards model have been extended to address the time-varying covariates through a counting process formulation, although sophisticated machine learning methods that can accommodate time-varying covariates have been limited. In this paper, we propose a non-parametric Bayesian survival model to analyze the nonlinear dependence of time-to-event outcomes on time-varying covariates. We focus on a computationally feasible Cox process called *permanental process*, which assumes the square root of hazard function to be generated from a Gaussian process, and tailor it for survival data with time-varying covariates. We verify that the proposed model holds with the representer theorem, a beneficial property for functional analysis, which offers us a fast Bayesian estimation algorithm that scales linearly with the number of observed events without relying on Markov Chain Monte Carlo computation. We evaluate our algorithm on synthetic and real-world data, and show that it achieves comparable predictive accuracy while being tens to hundreds of times faster than state-of-the-art methods.

1 Introduction

Survival or time-to-event data analysis has been widely applied for analyzing the dependence of survival time, the time until an event occurs, on covariates. They have a wide ranging list of applications in reliability engineering [27], finance [14], marketing [20], and especially, clinical research [4, 5, 28].

An essential part of the survival data analysis is to estimate a *hazard function*, that is, the instantaneous probability of events occurring at any particular time, as a regression function of covariates. Given a hazard function, we can evaluate the impact of covariates on survival times, and assess the degree of risk of experiencing an event in the future, through the survival function, i.e., the probability of no events occurring during a specified interval. The literature contains a vast number of studies that have proposed survival models to estimate hazard functions, most of which assume that covariates are time-invariant (e.g., gender and age at the time of diagnosis in clinical research): they range from the classical Cox proportional hazards model [5] to modern machine learning models based on generalized boosting machines [38], random forests [17], Gaussian processes [11, 25], and deep neural networks [10, 22, 49]. However, time-varying covariates are common in survival data, and the rapid advances in data collection allow us to access long-term and high temporal resolution covariate data. This is encouraging researchers to explore the non-stationarity in covariates for the accurate estimation/prediction of time-to-event outcomes. An important application of sur-

vival analysis with time-varying covariates is what-if analysis: through simulations of future events under various covariate functions over time, we can find the optimal covariate function or policy that would yield a desirable future state. For example, in reliability engineering applications where the events are the failures of machines, time-varying covariates could be the maintenance schedule and the room temperature/humidity controlled by air conditioners. The maintenance manager can optimize the schedules of maintenance and air conditioning by balancing the risk of failure and the costs of maintenance and air conditioning. If survival models for static covariates were applied to non-stationary data, then the static survival models would fail to estimate the underlying dependence of the hazard function on covariates, resulting in unreliable decision-making. Therefore, survival models that accommodate time-varying covariates are needed. This paper assesses survival data with time-varying covariates to estimate a hazard function that takes covariates into account.

The standard survival model is the Cox proportional hazards model [5, 6, 40] (CoxPH), which forms the logarithm of the hazard function as a linear combination of covariates. The original CoxPH assumes that covariates are time-invariant and have a constant log-linear effect over time on the hazard function. It was extended to address time-varying covariates through the technique of counting process formulation [1]. Although the extended CoxPH is very efficient to compute, it fails to address the nonlinear dependence of survival times on covariates, which is often the case in real-world data. To rectify this limitation, two non-parametric machine learning models have been proposed: generalized boosted models [16] and random forest-based models [47, 48], both of which adopt the extended CoxPH as a learner/tree, and construct non-parametric predictors in an ensemble manner. They succeeded in addressing the nonlinear dependence of survival times on time-varying covariates, but the large number of trees needed for accurate predictions tend to make the algorithm too slow and ineffective for real-time applications.

In this paper, we propose a novel Bayesian survival model to estimate a hazard function as a nonlinear regression function of covariates from survival data with time-varying covariates. To construct a scalable algorithm, we utilize *permanental process*, a doubly-stochastic point process which assumes the square root of hazard function to be generated from a Gaussian process; its computational advantages have been highlighted by recent studies [12, 18, 24, 30, 45]. We tailor it for survival data in the counting process format, and show that the tailored process, which we call the *survival permanental process*, holds the representer theorem [44]: the maximum a posteriori estimator of the latent function can be represented as a finite linear combination of a transformed kernel product evaluated on the event time points. The representer theorem offers us a feasible Bayesian estimation algorithm that scales linearly with the number of events observed without relying on Markov Chain Monte Carlo computation. Furthermore, we derive the predictive distribution and the marginal likelihood in one feasible form, enabling us to implement the uncertainty evaluation and the hyper-parameter optimization in a Bayesian manner. To the best of our knowledge, it is the first time to exploit the representer theorem for survival data analysis. We provide python codes to reproduce the results in this paper¹.

In Section 2, we introduce the survival permanental process (SurvPP) and construct a scalable Bayesian estimation algorithm for it. In Section 3, we outline related work. In Section 4, we compare SurvPP with reference survival models on synthetic and real-world data, and confirm that our algorithm achieves comparable predictive accuracy while being tens to hundreds of times faster than state-of-the-art survival models. Finally, Section 5 states our conclusions.

2 Methods

2.1 Survival Permanental Processes

We assume that a right-censored dataset with time-varying covariates, $\mathcal{D} = \{(\mathbf{y}_u(t), T_u, \eta_u)\}_{u=1}^U$, is observed, where $\mathbf{y}_u(t) : (0, T_u] \rightarrow \mathcal{Y} \subset \mathbb{R}^d$, $T_u \in \mathbb{R}$, and $\eta_u \in \{0, 1\}$ are the map of d -dimensional covariate, the end time of observation, and the indicator that represents whether an individual experienced an event ($\eta_u = 1$) or was right-censored ($\eta_u = 0$) at T_u , for individual $u \in \{1, 2, \dots, U\}$, respectively. We consider tailoring a permanental process for the survival data, where latent function, $x(\mathbf{y}) : \mathcal{Y} \rightarrow \mathbb{R}$, is generated from a Gaussian process (GP) on covariate space \mathcal{Y} , and an event time for each individual u is generated from a point process with hazard function, $\lambda_u(t)$,

¹Code and data are provided at <https://github.com/HidKim/SurvPP>.

such that

$$p(x(\mathbf{y})|\mathcal{D}) = \frac{p(\mathcal{D}|x(\mathbf{y})) \mathcal{GP}(x(\mathbf{y})|k)}{\int \mathcal{D}x(\mathbf{y}) p(\mathcal{D}|x(\mathbf{y})) \mathcal{GP}(x(\mathbf{y})|k)}, \quad (1)$$

$$\log p(\mathcal{D}|x(\mathbf{y})) = \sum_{u=1}^U \left[\eta_u \log \lambda_u(T_u) - \int_0^{T_u} \lambda_u(t) dt \right], \quad \lambda_u(t) = x^2(\mathbf{y}_u(t)), \quad (2)$$

where $\mathcal{GP}(x(\mathbf{y})|k)$ represents a GP with kernel function $k(\mathbf{y}, \mathbf{y}')$, $p(\mathcal{D}|x(\mathbf{y}))$ is the likelihood function of the point process, while $\int \mathcal{D}x(\mathbf{y})$ in the denominator represents the integral over the function or the infinite-dimensional variable $x(\mathbf{y})$. We call the model defined by (1-2) the *survival permanent process* (SurvPP). Hazard function $\lambda_u(t)$ represents an instantaneous probability of events occurring at each point in time, and our goal is to estimate the functional form of the hazard function over covariate domain, $x^2(\mathbf{y})$, from the survival data \mathcal{D} .

Given a hazard function over covariate domain, $x^2(\mathbf{y})$, and a covariate map for individual u , $\mathbf{y}_u(t)$, we can evaluate the survival function and the probability density distribution of event time at arbitrary time point t , denoted by $S(t)$ and $P_e(t)$, respectively, as follows:

$$S(t) = \exp \left[- \int_0^t x^2(\mathbf{y}_u(s)) ds \right], \quad P_e(t) = x^2(\mathbf{y}_u(t)) \exp \left[- \int_0^t x^2(\mathbf{y}_u(s)) ds \right]. \quad (3)$$

By using (3), we can perform survival analyses that include the analysis of the expected duration of event time, and the assessment of the degree of risk that an individual will experience the event before a specified period.

2.2 Counting Process Format of Data

Traditionally, survival data with time-varying covariates have taken the *counting process format* [1]. The format assumes that, for each individual u , covariates are measured at J_u finite representative time points and can be regarded as constant between successive time points,

$$\mathbf{y}_u(t) = \mathbf{y}_u^j, \quad t \in (s_u^j, s_u^{j+1}], \quad j = 0, \dots, J_u - 1, \quad (4)$$

where $(s_u^0, s_u^{J_u}) = (0, T_u)$. It then splits individual u 's record, $(\mathbf{y}_u(t), T_u, \eta_u)$, into J_u pseudo-individual records as in

$$\{(s_u^j, s_u^{j+1}, \xi_u^j, \mathbf{y}_u^j)\}_{j=0}^{J_u-1}, \quad \xi_u^j = \eta_u \cdot \mathbf{I}(j = J_u - 1), \quad (5)$$

where $\mathbf{I}(\cdot)$ represents the indicator. Joining the pseudo-individual records of U individuals together results in the counting process format of data:

$$\mathcal{D} = \{(T_j^0, T_j^1, \xi_j, \mathbf{y}_j)\}_{j=1}^J, \quad J = \sum_{u=1}^U J_u. \quad (6)$$

In this paper, we employ the counting process format of (6), and rewrite the likelihood function of SurvPP (2) as

$$\log p(\mathcal{D}|x(\mathbf{y})) = \sum_{j=1}^J \left[\xi_j \log x^2(\mathbf{y}_j) - (T_j^1 - T_j^0) x^2(\mathbf{y}_j) \right] = \sum_{n=1}^N \log x^2(\tilde{\mathbf{y}}_n) - \sum_{j=1}^J \Delta_j x^2(\mathbf{y}_j) \quad (7)$$

where Δ_j and $\{\tilde{\mathbf{y}}_n\}_{n=1}^N$ represent the durations of observation and the covariates observed at event times, respectively:

$$\Delta_j = T_j^1 - T_j^0, \quad \{\tilde{\mathbf{y}}_n\}_{n=1}^N = \{\mathbf{y}_j | \eta_j = 1, 1 \leq j \leq J\}, \quad N = \sum_{j=1}^J \xi_j = \sum_{u=1}^U \eta_u. \quad (8)$$

Note that the number of events, N , is usually much smaller than the number of pseudo-individual records: $N \ll J$.

In practice, there are two ways of defining the representative time points $\{s_u^j\}_{j=0}^{J_u}$: one is as the points that were observed, which might be sparse over time due to measurement constraints; the other is as the denser points obtained by using interpolation methods (e.g., splines). This paper assumes the latter, and thus various values of J_u were considered in the synthetic data experiment (see Section 4). Although the counting process format of input assumes piecewise stationarity in covariates, this assumption is not so strong because we can adopt the representative time points in any density if necessary.

2.3 Maximum A Posteriori Estimator

We consider the problem of obtaining the maximum *a posteriori* (MAP) estimator of $x(\mathbf{y})$, denoted by $\hat{x}(\mathbf{y})$, to maximize the posterior probability (1). We derive $\hat{x}(\mathbf{y})$ through the approach of using the path integral representation of GP [23],

$$\mathcal{GP}(x(\mathbf{y})|k) \mathcal{D}x(\mathbf{y}) = \sqrt{\frac{1}{|\mathcal{K}|}} \exp\left[-\frac{1}{2} \iint_{\mathcal{Y} \times \mathcal{Y}} k^*(\mathbf{y}, \mathbf{y}') x(\mathbf{y}) x(\mathbf{y}') d\mathbf{y} d\mathbf{y}'\right] \mathcal{D}x(\mathbf{y}), \quad (9)$$

where \mathcal{K} and $\int_{\mathcal{Y}} k^*(\mathbf{y}, \mathbf{y}') \cdot d\mathbf{y}'$ are the integral operator with kernel function $k(\mathbf{y}, \mathbf{y}')$ and its inverse operator, respectively: $\mathcal{K}^{(*)}x(\mathbf{y}) = \int_{\mathcal{Y}} k^{(*)}(\mathbf{y}, \mathbf{y}') x(\mathbf{y}') d\mathbf{y}'$, $\mathcal{K}k^*(\mathbf{y}, \mathbf{y}') = \delta(\mathbf{y} - \mathbf{y}')$. Here $|\mathcal{K}|$ represents the function determinant [15] of \mathcal{K} , defined by the product of its eigenvalues [23]. Using the representation of (9), we write the posterior of SurvPP (1-2) in the following functional form,

$$p(x(\mathbf{y})|\mathcal{D}) \mathcal{D}x = \frac{1}{p(\mathcal{D})} \exp\left[-S(x(\mathbf{y}), \underline{x}(\mathbf{y})) - \frac{1}{2} \log |\mathcal{K}|\right] \mathcal{D}x, \quad (10)$$

where $S(x(\mathbf{y}), \underline{x}(\mathbf{y}))$ is the *action integral*, defined by

$$S(x(\mathbf{y}), \underline{x}(\mathbf{y})) = \int_{\mathcal{Y}} \left[\frac{1}{2} x(\mathbf{y}) \underline{x}(\mathbf{y}) + \sum_{j=1}^J x^2(\mathbf{y}) \Delta_j \delta(\mathbf{y} - \mathbf{y}_j) - 2 \sum_{n=1}^N \log |x(\mathbf{y})| \delta(\mathbf{y} - \tilde{\mathbf{y}}_n) \right] d\mathbf{y}, \quad (11)$$

and $\underline{x}(\mathbf{y}) = \int_{\mathcal{Y}} k^*(\mathbf{y}, \mathbf{y}') x(\mathbf{y}') d\mathbf{y}'$. Then we apply calculus of variations to the action integral, where the functional derivative of $S(x(\mathbf{y}), \underline{x}(\mathbf{y}))$ on the MAP estimator $\hat{x}(\mathbf{y})$ should be equal to zero: $\frac{\delta S}{\delta \hat{x}(\mathbf{y})} \delta \hat{x}(\mathbf{y}) + \frac{\delta S}{\delta \hat{\underline{x}}(\mathbf{y})} \delta \hat{\underline{x}}(\mathbf{y}) = 0$, which results in the exact MAP estimator,

$$\hat{x}(\mathbf{y}) = 2 \sum_{n=1}^N h(\mathbf{y}, \tilde{\mathbf{y}}_n) v_n, \quad v_n = \hat{x}(\tilde{\mathbf{y}}_n)^{-1}, \quad (12)$$

where $h(\mathbf{y}, \mathbf{y}')$ is a transformed kernel function that solves a discretized Fredholm integral equation of the second kind [36],

$$h(\mathbf{y}, \mathbf{y}') + 2 \sum_{j=1}^J k(\mathbf{y}, \mathbf{y}_j) \Delta_j h(\mathbf{y}_j, \mathbf{y}') = k(\mathbf{y}, \mathbf{y}'). \quad (13)$$

See Appendix A for the detailed derivations. Equation (12) shows that the MAP estimator of SurvPP involves the representer theorem under the transformed kernel function $h(\mathbf{y}, \mathbf{y}')$, and thus the Bayesian estimation reduces to a finite-dimensional optimization problem. Here, we call $h(\mathbf{y}, \mathbf{y}')$ the equivalent kernel following studies by Flaxman et al. [12], Walder & Bishop [45], and Kim et al. [24]. Given the equivalent kernel, the unknown coefficients v_n in (12) solve the simultaneous quadratic equations derived from (12),

$$r_n \triangleq 2 v_n \sum_{n'=1}^N h(\tilde{\mathbf{y}}_n, \tilde{\mathbf{y}}_{n'}) v_{n'} - 1 = 0, \quad n = 1, 2, \dots, N. \quad (14)$$

In this paper, we estimate a set of coefficients, $\{v_n\}_{n=1}^N$, by solving a minimization problem of the mean of the squared residuals, $\sum_{n=1}^N |r_n|^2 / N$, with a popular gradient descent algorithm, *Adam* [26].

2.4 Equivalent Kernels

The equivalent kernel $h(\mathbf{y}, \mathbf{y}')$ solves the discretized Fredholm integral equation (13). When the kernel function of GP has degenerate form with rank $M (< \infty)$ such that

$$k(\mathbf{y}, \mathbf{y}') = \sum_{m=1}^M \phi_m(\mathbf{y}) \phi_m(\mathbf{y}') = \boldsymbol{\phi}(\mathbf{y})^\top \boldsymbol{\phi}(\mathbf{y}'), \quad (15)$$

it is easily shown that the discretized Fredholm integral equation (13) can be solved analytically [2, 24] as,

$$h(\mathbf{y}, \mathbf{y}') = \boldsymbol{\phi}(\mathbf{y})^\top (\mathbf{I}_M + 2\mathbf{A})^{-1} \boldsymbol{\phi}(\mathbf{y}'), \quad \mathbf{A} = \sum_{j=1}^J \Delta_j \boldsymbol{\phi}(\mathbf{y}_j) \boldsymbol{\phi}(\mathbf{y}_j)^\top, \quad (16)$$

where \mathbf{I}_M is the $M \times M$ identity matrix, and $\boldsymbol{\phi}(\mathbf{y}) = (\phi_1(\mathbf{y}), \phi_2(\mathbf{y}), \dots, \phi_M(\mathbf{y}))^\top$. Note that \mathbf{A} is defined by a sum of outer products, which can be rewritten through a matrix-matrix multiplication as

$$\mathbf{A} = \mathbf{B}\mathbf{B}^\top, \quad \mathbf{B} = [\sqrt{\Delta_1} \boldsymbol{\phi}(\mathbf{y}_1), \dots, \sqrt{\Delta_J} \boldsymbol{\phi}(\mathbf{y}_J)]. \quad (17)$$

Empirically, implementation by matrix-matrix multiplication (17) is substantially faster than the sum of outer products when $J \gg 1$.

When $k(\mathbf{y}, \mathbf{y}')$ has degenerate form with rank M , the relation (16) shows that the equivalent kernel $h(\mathbf{y}, \mathbf{y}')$ also has degenerate form obtained through Cholesky decomposition:

$$h(\mathbf{y}, \mathbf{y}') = (\mathbf{L}\boldsymbol{\phi}(\mathbf{y}))^\top (\mathbf{L}\boldsymbol{\phi}(\mathbf{y}')), \quad \mathbf{L}^\top \mathbf{L} = (\mathbf{I}_M + 2\mathbf{A})^{-1}. \quad (18)$$

The degenerate equivalent kernel (18) offers fast Bayesian estimation that scales linearly with N (see Section 2.6). In this paper, we used the random feature map [37, 42, 43] of a Gaussian kernel to obtain a degenerate form of kernel ($M=500$).

If the covariate map, $\mathbf{y}_u(t)$, can be assumed to be smooth enough over time, a more accurate evaluation of \mathbf{A} than equation (16) is possible (see Appendix D.1), but this was not exploited in the main experiments.

2.5 Predictive Distribution and Marginal Likelihood

One of the advantages of GP models over non-Bayesian approaches is that they can provide predictive distributions and marginal likelihoods, which enable us to perform uncertainty evaluations, hyper-parameter optimization, and model selection in Bayesian manner. Following the methodology with the path integral representation of GP [23, 24], SurvPP (1-2) adopts a Laplace approximation in the functional space, and finds the approximate form of the predictive distribution and the marginal likelihood. We only show the results due to space limitations. For details, see Appendix B.

The marginal likelihood, $p(\mathcal{D})$, is given as

$$\log p(\mathcal{D}) = \log |\mathbf{Z}| - \frac{1}{2} \log |\mathbf{I}_N + \mathbf{Z}^{-1} \mathbf{H}| - \frac{1}{2} \log |\mathbf{I}_M + 2\mathbf{A}|^{-1} + (\log 2 - 1)N, \quad (19)$$

where \mathbf{I}_N is the identity matrix with size N , and

$$\mathbf{Z}_{nn'} = (2v_n^2)^{-1} \delta_{nn'}, \quad \mathbf{H}_{nn'} = h(\tilde{\mathbf{y}}_n, \tilde{\mathbf{y}}_{n'}), \quad \mathbf{h}(\mathbf{y}) = (h(\mathbf{y}, \tilde{\mathbf{y}}_1), \dots, h(\mathbf{y}, \tilde{\mathbf{y}}_N))^\top. \quad (20)$$

The predictive distribution of intensity function on a covariate value, $p(\lambda = x^2(\mathbf{y}))$, is given as

$$p_{\mu, \nu}(\lambda) = \frac{1}{\Gamma(\nu) \mu^\nu} \lambda^{\nu-1} \exp(-\lambda/\mu), \quad \mu = 2\sigma \frac{2\hat{x}^2 + \sigma}{\hat{x}^2 + \sigma}, \quad \nu = \frac{(\hat{x}^2 + \sigma)^2}{2\sigma(2\hat{x}^2 + \sigma)}, \quad (21)$$

where \hat{x} is the abbreviation of $\hat{x}(\mathbf{y})$, and the predictive variance, σ , is defined by $\sigma = h(\mathbf{y}, \mathbf{y}) - \mathbf{h}(\mathbf{y})^\top (\mathbf{Z} + \mathbf{H})^{-1} \mathbf{h}(\mathbf{y})$. (21) represents the predictive distribution of the hazard function as a regression function of covariates. It is worth noting that the posterior process of SurvPP is given by a permanental process, and we can evaluate the distribution of the survival function and perform a risk-aware survival analysis by generating random samples of the estimated hazard function of covariates, $\lambda(\mathbf{y})$.

2.6 Computational Complexity

The computational complexity of evaluating the equivalent kernel (16) for covariate pair $(\mathbf{y}, \mathbf{y}')$ is $\mathcal{O}(M^3 + dJM + JM^2)$, where $\mathcal{O}(dJM)$ and $\mathcal{O}(JM^2)$ come from the computation of feature maps, $\boldsymbol{\phi}_m(\mathbf{y}_j)$, and the sum of outer products (17), respectively.

When the equivalent kernel is given in degenerate form with rank $M (< N)$ such that $h(\mathbf{y}, \mathbf{y}') = \sum_{m=1}^M \varphi_m(\mathbf{y}) \varphi_m(\mathbf{y}')$, the objective function to be minimized in MAP estimation, $\sum_{n=1}^N |r_n|^2 / N$,

incurs a linear computation with the number of observed events N , that is, $\mathcal{O}(NM)$, for each evaluation in gradient descent algorithms: The vector of residual, $\mathbf{r} = (r_1, \dots, r_N)^\top$, can be expressed as $(\mathbf{R}(\mathbf{R}^\top \mathbf{v})) \odot (2\mathbf{v}) - \mathbf{1}$, where \mathbf{R} is the $N \times M$ matrix defined by $R_{nm} = \varphi_m(\mathbf{y}_n)$, and \odot represents the Hadamard product. \mathbf{R} costs $\mathcal{O}(dNM)$.

Given an equivalent kernel with degenerate form, the evaluation of the predictive variance, $\sigma = h(\mathbf{y}, \mathbf{y}) - \mathbf{h}(\mathbf{y})^\top (\mathbf{Z} + \mathbf{H})^{-1} \mathbf{h}(\mathbf{y})$, needs the computation of $\mathcal{O}(M^3 + NM^2)$: $N \times N$ matrix \mathbf{H} can be decomposed into a product of $N \times M$ matrix \mathbf{R} and its transpose as $\mathbf{H} = \mathbf{R}\mathbf{R}^\top$. The matrix inversion is transformed as $(\mathbf{Z} + \mathbf{R}\mathbf{R}^\top)^{-1} = \mathbf{Z}^{-1} - \mathbf{Z}^{-1}\mathbf{R}(\mathbf{I}_M + \mathbf{R}^\top \mathbf{Z}^{-1}\mathbf{R})^{-1}\mathbf{R}^\top \mathbf{Z}^{-1}$ through the Woodbury matrix identity, which costs $\mathcal{O}(M^3 + NM^2)$. Note that $\mathbf{Z}^{-1}\mathbf{R}$ and $\mathbf{R}^\top \mathbf{Z}^{-1}$ cost $\mathcal{O}(NM)$ because \mathbf{Z} is a diagonal matrix.

In computing the marginal likelihood (19), the complexities of the first, the second, and the third terms are $\mathcal{O}(N)$, $\mathcal{O}(M^3 + NM^2)$, and $\mathcal{O}(M^3 + JM^2)$, respectively, where the matrix determinant lemma is used in computing the second term.

In total, the computational complexity of SurvPP is $\mathcal{O}(NMQ + (N + J)(d + M)M + M^3)$, where Q is the number of gradient descent iterations. This feasible computation, which scales linearly with the number of observed events, N , the data size, J , and the covariate dimensionality, d , is achieved by exploiting the representer theorem. It should be emphasized here that data size J is not part of the gradient descent iteration term. This is a clear advantage over conventional survival models which naively incur $\mathcal{O}((dN + dJ + JN)Q)$ computation: Typically, J is larger than N in a time-varying covariate scenario (see the counting process format in Section 2.2), and the discrepancy between J and N becomes more substantial when the measurements of covariates are made more frequently or/and over a longer period of time.

3 Related Work

Survival Models for Time-invariant Covariates: The most popular survival model is the Cox proportional hazards model [5] (CoxPH), which is a semi-parametric model that forms the hazard function as

$$\lambda(t) = h(t) \exp(\beta_1 y_1 + \beta_2 y_2 + \dots), \quad (22)$$

where the base hazard function, $h(t)$, is obtained by the non-parametric Breslow/Fleming-Harrington estimator [3, 13], and the log-linear regression coefficient (β_1, \dots) is estimated by maximizing the partial likelihood function. CoxPH is very efficient to compute and scales linearly with the number of events [39], but cannot address the nonlinear dependence of survival times on covariates. To overcome this limitation, a vast number of survival models have been proposed that replace the log-linear parametric function (22) with a non-linear one, such as generalized boosted models [38], random survival forests [17], Gaussian process models [11, 25], and deep neural network models [10, 22, 49]. For a comprehensive review, see Wang et al. [46].

Survival Models for Time-varying Covariates: The counting process format of input [1] plays a central role in extending static survival models into those suitable for time-varying covariates. For each individual, the counting process format splits her/his observation period into multiple short intervals, assigns a constant value of covariate to each interval, and marks the intervals as being censored when no event is present, resulting in a set of pseudo-individual right/left-censored observations with a constant covariate (see Section 2.2). CoxPH was extended to accommodate the counting process format, and the extended CoxPH can estimate the hazard function based on survival data with time-varying covariates [1]. To alleviate the simplest assumption in the extended CoxPH, generalized boosted models [16] and random survival forests [47, 48], both of which are ensemble approaches, have been also extended to accommodate the counting process format, where the extended CoxPH is employed as a weak learner or a tree. While they can address the nonlinear dependence of survival times on time-varying covariates, the large number of learners/trees needed for accurate predictions could make the algorithms too slow and ineffective for real-time applications. Also, Cygu et al. reported that random survival forests for time-varying covariates required substantial amounts of computer memory for large datasets [7].

Exogeneity and Endogeneity of Covariates: When considering survival analysis in the presence of time-varying covariates, we need to distinguish between exogenous and endogenous covariates. Kalbfleisch and Prentice [21] define an endogenous (internal) covariate as the output of a stochastic

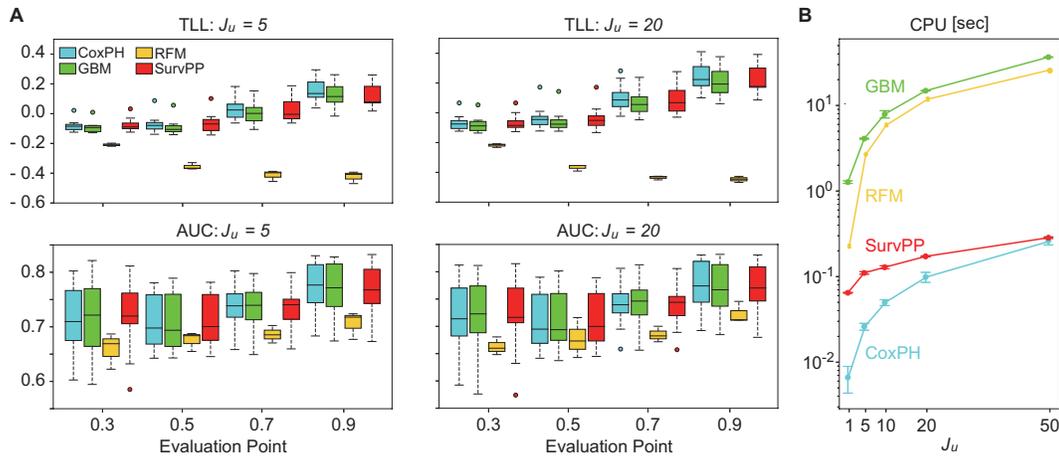


Figure 1: Performance on the dataset of log-linear hazard function $\lambda_{lin}(t)$. (A) Box plot of TLL and AUC as functions of evaluation point: the higher, the better. (B) The CPU times demanded for estimating a hazard function. The error bars represent the standard deviations. For GBM and SurvPP, the average cpu times over 9-points grid search of the hyperparameter are displayed.

process that is generated by the individual under study. In contrast, an exogenous (external) covariate is not influenced by the individual under study. Diggle et al. [9] suggest using similar but slightly different definitions. The primary purpose of survival analysis with exogenous covariates is to estimate a hazard function as an explicit function of covariates, and to make predictions of failure times under various possible *future* covariate functions (i.e., what-if analysis). In contrast, the primary purpose of survival analysis with endogenous covariates is to make predictions of failure times by using the *past* observations of covariates, where joint modeling approaches with recurrent neural networks, which jointly model the stochastic process of covariates and failure times, have been developed intensively [29, 33]. In this paper, we consider exogenous covariates, and survival models for exogenous covariates cannot be compared directly with those for endogenous covariates because the tasks are different. Note that a joint modeling approach for exogenous time-varying covariates has been proposed very recently[32], which was not included in the benchmark models because it was published a month after this paper’s submission.

Permanental Process: The permanental process is a variant of Gaussian Cox process that assumes the square root of hazard function to be generated from a Gaussian process [31]. Its computational advantages have recently been highlighted in machine learning research [12, 18, 24, 25, 30, 45]. In particular, the representer theorem in the permanental process has been exploited through the RKHS theory [12], the Mercer’s theorem [45], and the path integral formulation [23, 24]. Although the key derivation of SurvPP is based on the path integral methodology used for augmented permanental process (APP) [24], SurvPP is a non-trivial extension of APP: (i) SurvPP can accommodate multiple trials of event sequence data, while APP assumes one trial of event sequence; (ii) in SurvPP, the end time of observation is a stochastic variable that depends on the time of event occurrence, while APP assumes that the end time of observation is given. We discovered, for the first time, that the representer theorem holds for such a complicated point process or a survival model.

4 Experiments

We examined the validity of our proposed model by comparing it with conventional survival models on synthetic and real-world data. As benchmark models, we adopted Cox proportional hazards model (CoxPH), generalized boosted model (GBM), and random forest-based model (RFM): We implemented CoxPH and GBM with the established algorithms provided in the packages `survival.coxph` [41] and `gbm3.gbm` [16], respectively; We implemented RFM by using class `ltrcrf` in the open R code provided by Yao et al. [48]. The benchmark models can accommodate time-varying covariates via counting process format of input (`Surv(t0, t1, event)` in `survival`). For our proposal, we implemented SurvPP by using TensorFlow-2.10¹. A MacBook Pro with 12-core CPU (Apple M2 Max) was used, where GPU was set as off (`tf.device('/cpu:0')`) for a

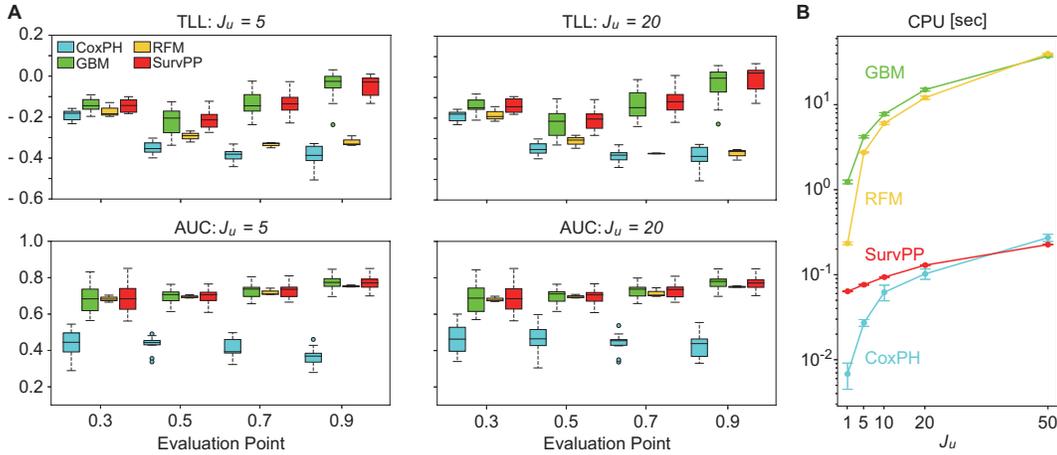


Figure 2: Performance on the dataset of non-linear hazard function $\lambda_{non}(t)$. (A) Box plot of TLL and AUC as functions of evaluation point: the higher, the better. (B) The CPU times demanded for estimating a hazard function. The error bars represent the standard deviations. For GBM and SurvPP, the average cpu times over 9-point grid search of the hyperparameter are displayed.

fair benchmark comparisons. For GBM and SurvPP, the hyper-parameters were optimized through 9-point grid search: the number of trees and the shrinkage for GBM, and the kernel parameters for SurvPP. For details of the model configurations, see Appendix C.

4.1 Synthetic Data

We created two survival datasets with time-varying 2-D covariate. One was generated from a hazard function with a log-linear dependence on covariates, and the other with a nonlinear dependence on covariates as follows:

$$\lambda_{lin}(t) = h(t) \exp[1.5y_1(t) + 3.5y_2(t)], \quad \lambda_{non}(t) = h(t) \exp[2 - 5(y_1^2(t) + y_2^2(t))], \quad (23)$$

where $h(t)$ is the base hazard function defined by a Weibull hazard function, $h(t) = 2 \cdot t^{3/2}$. For each of the datasets, we considered U individuals, each of which had a 2-D covariate function of time:

$$\mathbf{y}_u(t) = (y_1^u(t), y_2^u(t)) = (\alpha_u \cos(2\pi\omega_u t + \pi\gamma_u), \alpha'_u \cos(2\pi\omega'_u t + \pi\gamma'_u)), \quad u = 1, \dots, U, \quad (24)$$

where α_u (α'_u) and γ_u (γ'_u) were sampled uniformly over $[0, 1]$, ω_u was sampled uniformly over $[5, 10]$, and ω'_u was sampled uniformly over $[20, 30]$. We set the censoring time as 1 for all individuals, and an event time generated from (23), denoted by t_* , was censored ($\eta_u = 0$) when t_* was over 1: $T_u = \min(t_*, 1)$.

The predictive performances were evaluated based on the dynamic area under the ROC curve (AUC) [34] and the test log-likelihood (TLL),

$$\text{AUC}(t) = \frac{\sum_{u=1}^U \sum_{u'=1}^U \mathbf{I}(T_u > t) \mathbf{I}(T_{u'} \leq t) w_u \mathbf{I}(S(t|\mathbf{y}_u(\cdot)) \geq S(t|\mathbf{y}_{u'}(\cdot)))}{(\sum_{u=1}^U \mathbf{I}(T_u > t)) (\sum_{u=1}^U \mathbf{I}(T_u \leq t) w_u)}, \quad (25)$$

$$\text{TLL}(t) = \frac{1}{U} \sum_{u=1}^U \left[\mathbf{I}(t \geq T_u) \cdot \eta_u \log \rho(T_u | \mathbf{y}_u(\cdot)) + \log S(\min(t, T_u) | \mathbf{y}_u(\cdot)) \right], \quad (26)$$

where w_u is the inverse probability of censoring weight, and $\rho(t|\mathbf{y}_u(\cdot))$ and $S(t|\mathbf{y}_u(\cdot))$ are the estimated hazard and survival functions given covariate map $\mathbf{y}_u(\cdot)$, respectively. The dynamic AUC essentially estimates the C-index at each time, and is commonly used in the literature of survival analysis with time-varying covariates.

The evaluation point, $0 \leq t_e \leq 1$, was set as $t_e \in \{0.3, 0.5, 0.7, 0.9\}$. For each dataset, we randomly split the U individuals into 10 subgroups, assigned one to test and the others to training data, and

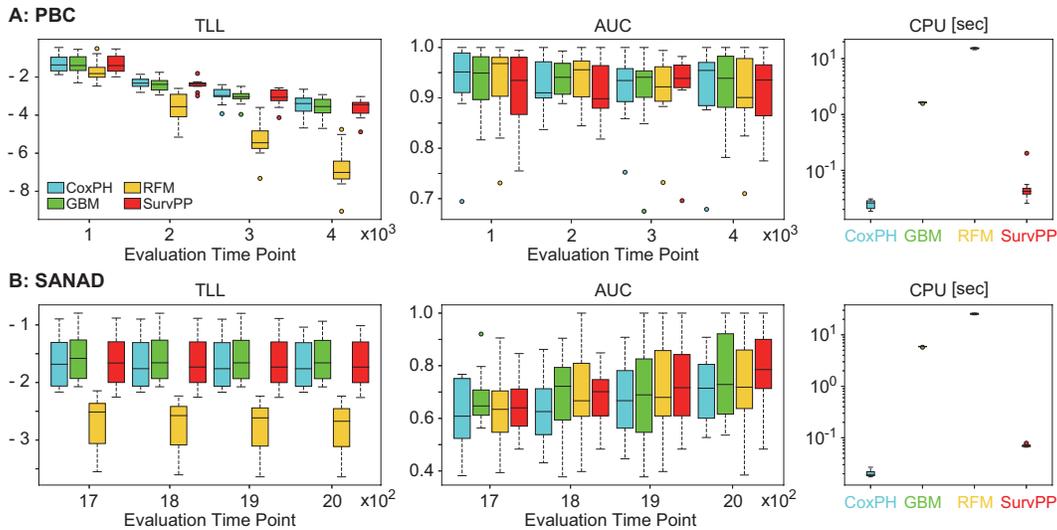


Figure 3: Performances on real-world datasets. TLL (the higher, the better), AUC (the higher, the better), and CPU times (the lower, the better) on PBC dataset (A) and SANAD dataset (B).

conducted 10-fold cross evaluation of the predictive performances. We reformatted the training data into counting process format by equally splitting individual u 's observation periods into J_u sub-region and assuming the covariates to be constant within a sub-region; the survival models were applied to the training data. In this experiment, we considered $U = 10^3$ and $J_u \in \{1, 5, 10, 20, 50\}$: the number of observed events N was 809 for λ_{lin} and 818 for λ_{non} , respectively.

It should be noted here that when time-varying covariates are considered, the time duration from the entry, t , is itself a time-varying covariate. SurvPP considered a 3-D covariate map, $\mathbf{y}_u(t) \equiv (y_1^u(t), y_2^u(t), t)$, and performed hazard function estimation, while the other models assumed the hazard function to be $h(t)f(\mathbf{y}(t))$, and $h(t)$ and $f(\mathbf{y})$ were estimated separately.

Figure 1 displays the predictive performance on the dataset of $\lambda_{lin}(t)$. Figure 1A shows that, except for RFM, there were no significant differences in the performances between the compared models, but CoxPH slightly outperformed the other models in TLL at some evaluation periods. The result is plausible because the underlying generative process was perfectly consistent with the log-linear and proportional assumptions of CoxPH. The comparison in TLL between $J_u = 5$ and $J_u = 20$ suggests that more frequently measured covariates would yield better estimations/predictions in scenarios of time-varying covariates. Therefore, the feasible computational complexity of SurvPP regarding J_u is a clear advantage over the reference models, which is confirmed by Figure 1B.

Figure 2 displays the predictive performance on the dataset of $\lambda_{non}(t)$, showing that the non-parametric models, SurvPP and GBM, significantly outperformed CoxPH, while the performance gaps between SurvPP and GBM were marginal and of no significance. However, Figure 2B shows that SurvPP could estimate the hazard function hundreds of times faster than GBM when covariates were measured frequently ($J_u = 50$), and was even faster than the simplest CoxPH. The preferable scalability of SurvPP's algorithm regarding J_u comes from the fact that J_u (or $J = \sum_u J_u$) is not part of the iterative optimization (Section 2.6), which is a fruit of SurvPP's representer theorem. RFM was comparable with SurvPP and GBM in AUC, but was worse than either of them in TLL. RFM might need more careful tuning of the hyper-parameters, but this was not fully investigated this time.

We conducted additional experiments on independent validation datasets (see Appendix D.2), and on larger synthetic datasets ($U \leq 10^5$) to examine the model's computation scalability regarding the event number N (see Appendix D.3).

4.2 Real-world Data

We examined the validity of SurvPP against the benchmark models on two real-world survival data sets, *Mayo Clinic Primary Biliary Cholangitis data* (PBC) and *Standard And New Antiepileptic*

Drugs study data (SANAD), provided by R packages *survival* (LGPL-3) [41] and *joiner* (GPL-3) [35], respectively. In PBC, 312 patients with primary biliary cirrhosis were enrolled in a randomized medical trial at the Mayo Clinic between 1974 and 1984 [8], where events were the time of death; one static and eleven time-varying covariates were measured at entry and at yearly intervals, which include age at entry, alkaline phosphatase, logarithm of serum albumin, presence of ascites, aspartate aminotransferase, logarithm of serum bilirubin, serum cholesterol, condition of edema, presence of hepatomegaly or enlarged liver, platelet count, logarithm of prothrombin time, and presence or absence of spiders. SANAD was an unblind randomized trial that recruited patients with epilepsy for whom carbamazepine (CBZ) was considered to be standard treatment and they were randomized to CBZ or the newer drug lamotrigine (LTG), where 605 patients were included and event was the time to treatment failure; we adopted calibrated dose as a time-varying covariate, and three static covariates including age of patient at randomization, gender, and randomized treatment (CBZ or LTG); calibrated dose was measured at 166-day intervals on average, and we used a linear interpolation to obtain the values of calibrated dose at which the measurement intervals were regularly trisected. We randomly split the individuals into 10 subgroups, assigned one to test and the others to training data, and conducted 10-fold cross evaluation of the predictive performances.

Figure 3A displays the predictive performance on PBC. It shows that CoxPH achieved very high AUCs (> 0.9), suggesting that the underlying generative process could be consistent with CoxPH's simple assumption of log-linearity and proportionality. Thus the semi-parametric model, CoxPH, is likely to achieve equal or slightly better performances than the nonparametric models, which is consistent with the result. Figure 3B plots the predictive performance on SANAD, where the not so high AUCs of CoxPH suggests that the underlying dependence of intensity on covariates is nonlinear. The figure shows that SurvPP achieved better performance than CoxPH, and achieved comparable performance while being substantially faster than GBM.

5 Conclusions

We have proposed a non-parametric Bayesian survival model to address survival data with time-varying covariates. We tailored a permanental process such that the latent hazard function is defined on covariate space and right-censored observations in a counting process format can be handled, which we call the survival permanental process (SurvPP). Through the path integral formulation of Gaussian process, we showed that SurvPP encompasses a representer theorem, and derived a feasible estimation algorithm that scales linearly with the number of observed events. We evaluated SurvPP on synthetic data, confirming that it achieved comparable predictive accuracy while being tens to hundreds of times faster than state-of-the-art methods.

Limitations & future work: We examined SurvPP for relatively low-dimensional covariates, and its potential suitability for high-dimensional covariates remains to be clarified. Because SurvPP is based on a normal Gaussian process, the (equivalent) kernel function might be too simple to discover meaningful representations in high-dimensional covariate data. A promising direction is to apply deep kernel learning to SurvPP, where high-dimensional covariates are transformed by nonlinear mapping with a deep architecture. A technical issue is that we could not search an appropriate set of (a lot of) parameters of neural networks, because our proposed scheme performs the hyperparameter optimization by grid search, not by gradient descent. Also, Gaussian kernels, which were used in the paper, naively require a scale parameter for each dimension of data, resulting in a high dimensional kernel parameter for high-dimensional covariate scenarios. The grid search becomes prohibitively costly with high dimensional parameters, but we addressed the problem by a well-known approach that normalizes (e.g., centering and scaling) each dimension of the data and puts a common scale parameter across all dimensions of normalized data. This approach empirically works robustly, but a more sophisticated approach could improve the performance of SurvPP. Variational Bayesian approximations with inducing points might address the technical issues of hyperparameter optimization, and thus is an important next step in our study.

As in ordinary permanental processes, the nodal line problem could arise in SurvPP: the posterior distribution of the latent variable $x(\cdot)$ has many local modes since $\pm x(\cdot)$ can lead to similar hazard functions $\lambda = x^2(\cdot)$, and artificial zero crossings of $x(\cdot)$ could happen, especially on locations where the hazard function is low. John and Hensman [18] have proposed to extend the quadratic link function to include an offset parameter β , so that $\lambda(\mathbf{y}) = (x(\mathbf{y}) + \beta)^2$, which is valid for SurvPP.

References

- [1] Per Kragh Andersen and Richard D. Gill. Cox's regression model for counting processes: A large sample study. *The Annals of Statistics*, pages 1100–1120, 1982.
- [2] Kendall Atkinson. A personal perspective on the history of the numerical analysis of Fredholm integral equations of the second kind. In *The Birth of Numerical Analysis*, pages 53–72. World Scientific, 2010.
- [3] Norman E. Breslow. Contribution to discussion of paper by DR Cox. *Journal of the Royal Statistical Society, Series B*, 34:216–217, 1972.
- [4] Taane G. Clark, Michael J. Bradburn, Sharon B. Love, and Douglas G. Altman. Survival Analysis Part I: Basic concepts and first analyses. *British Journal of Cancer*, 89(2):232–238, 2003.
- [5] David R. Cox. Regression models and life-tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, 34(2):187–202, 1972.
- [6] Steve Cygu, Jonathan Dushoff, and Benjamin M. Bolker. pcoxtime: Penalized Cox proportional hazard model for time-dependent covariates. *arXiv preprint arXiv:2102.02297*, 2021.
- [7] Steve Cygu, Hsien Seow, Jonathan Dushoff, and Benjamin M. Bolker. Comparing machine learning approaches to incorporate time-varying covariates in predicting cancer survival time. *Scientific Reports*, 13(1):1370, 2023.
- [8] E. Rolland Dickson, Patricia M. Grambsch, Thomas R. Fleming, Lloyd D. Fisher, and Alice Langworthy. Prognosis in primary biliary cirrhosis: Model for decision making. *Hepatology*, 10(1):1–7, 1989.
- [9] Peter J. Diggle. *Analysis of Longitudinal Data*. Oxford University Press, 2002.
- [10] David Faraggi and Richard Simon. A neural network model for survival data. *Statistics in Medicine*, 14(1):73–82, 1995.
- [11] Tamara Fernández, Nicolás Rivera, and Yee Whye Teh. Gaussian processes for survival analysis. In *Advances in Neural Information Processing Systems 29*, 2016.
- [12] Seth Flaxman, Yee Whye Teh, and Dino Sejdinovic. Poisson intensity estimation with reproducing kernels. In *Artificial Intelligence and Statistics*, pages 270–279. PMLR, 2017.
- [13] Thomas R. Fleming and David P. Harrington. Nonparametric estimation of the survival distribution in censored data. *Communications in Statistics-Theory and Methods*, 13(20):2469–2486, 1984.
- [14] Adrian Gepp and Kuldeep Kumar. The role of survival analysis in financial distress prediction. *International Research Journal of Finance and Economics*, 16(16):13–34, 2008.
- [15] Israel Gohberg, Seymour Goldberg, and Nahum Krupnik. *Traces and Determinants of Linear Operators*, volume 116. Birkhäuser, 2012.
- [16] James Hickey and Greg Ridgeway. *Generalized Boosted Regression Models*, 2023. URL <https://github.com/gbm-developers/gbm3>.
- [17] Hemant Ishwaran, Udaya B. Kogalur, Eugene H. Blackstone, and Michael S. Lauer. Random survival forests. *The Annals of Applied Statistics*, 2(3):841–860, 2008.
- [18] S. T. John and James Hensman. Large-scale Cox process inference using variational Fourier features. In *International Conference on Machine Learning*, volume 80, pages 2362–2370. PMLR, 2018.
- [19] Alistair E.W. Johnson, Tom J. Pollard, Lu Shen, Li-wei H. Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. MIMIC-III, a freely accessible critical care database. *Scientific Data*, 3(1):1–9, 2016.

- [20] Euy-Young Jung, Chulwoo Baek, and Jeong-Dong Lee. Product survival analysis for the app store. *Marketing Letters*, 23:929–941, 2012.
- [21] John D. Kalbfleisch and Ross L. Prentice. *The Statistical Analysis of Failure Time Data*. John Wiley & Sons, 2011.
- [22] Jared L. Katzman, Uri Shaham, Alexander Cloninger, Jonathan Bates, Tingting Jiang, and Yuval Kluger. DeepSurv: personalized treatment recommender system using a Cox proportional hazards deep neural network. *BMC Medical Research Methodology*, 18(1):1–12, 2018.
- [23] Hideaki Kim. Fast Bayesian inference for Gaussian Cox processes via path integral formulation. In *Advances in Neural Information Processing Systems 34*, 2021.
- [24] Hideaki Kim, Taichi Asami, and Hiroyuki Toda. Fast Bayesian estimation of point process intensity as function of covariates. In *Advances in Neural Information Processing Systems 35*, 2022.
- [25] Minyoung Kim and Vladimir Pavlovic. Variational inference for Gaussian process models for survival analysis. In *Uncertainty in Artificial Intelligence*, pages 435–445, 2018.
- [26] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [27] Chin Diew Lai and Min Xie. *Stochastic Ageing and Dependence for Reliability*. Springer Science & Business Media, 2006.
- [28] András Lániczky and Balázs Györfly. Web-based survival analysis tool tailored for medical research (kmlplot): development and implementation. *Journal of Medical Internet Research*, 23(7):e27633, 2021.
- [29] Changhee Lee, Jinsung Yoon, and Mihaela Van Der Schaar. Dynamic-DeepHit: A deep learning approach for dynamic survival analysis with competing risks based on longitudinal data. *IEEE Transactions on Biomedical Engineering*, 67(1):122–133, 2019.
- [30] Chris Lloyd, Tom Gunter, Michael Osborne, and Stephen Roberts. Variational inference for Gaussian process modulated Poisson processes. In *International Conference on Machine Learning*, volume 37, pages 1814–1822. PMLR, 2015.
- [31] Peter McCullagh and Jesper Møller. The permanental process. *Advances in Applied Probability*, 38(4):873–888, 2006.
- [32] Victor Medina-Olivares, Raffaella Calabrese, Jonathan Crook, and Finn Lindgren. Joint models for longitudinal and discrete survival data in credit scoring. *European Journal of Operational Research*, 307(3):1457–1473, 2023.
- [33] Chirag Nagpal, Vincent Jeanselme, and Artur Dubrawski. Deep parametric time-to-event regression with time-varying covariates. In *Survival Prediction-Algorithms, Challenges and Applications*, pages 184–193. PMLR, 2021.
- [34] Norberto Pantoja-Galicia, Olivia I. Okereke, Deborah Blacker, and Rebecca A. Betensky. Concordance measures and time-dependent ROC methods. *Biostatistics & Epidemiology*, 5(2): 232–249, 2021.
- [35] Pete Philipson, Ines Sousa, Peter J. Diggle, Paula Williamson, Ruwanthi Kolamunnage-Dona, Robin Henderson, and Graeme L. Hickey. *joiner: Joint Modelling of Repeated Measurements and Time-to-Event Data*, 2018. URL <https://github.com/graemeleehickey/joiner/>. R package version 1.2.8.
- [36] Andrei D. Polyanin and Alexander V. Manzhirov. *Handbook of Integral Equations*. CRC press, 1998.
- [37] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems 20*, 2007.

- [38] Greg Ridgeway. Generalized Boosted Models: A guide to the gbm package. *Update*, 1(1): 2007, 2007.
- [39] Noah Simon, Jerome Friedman, Trevor Hastie, and Rob Tibshirani. Regularization paths for Cox’s proportional hazards model via coordinate descent. *Journal of Statistical Software*, 39(5):1, 2011.
- [40] Terry Therneau, Cindy Crowson, and Elizabeth Atkinson. Using time dependent covariates and time dependent coefficients in the Cox model. *Survival Vignettes*, 2:3, 2017.
- [41] Terry M. Therneau. *A Package for Survival Analysis in R*, 2023. URL <https://CRAN.R-project.org/package=survival>. R package version 3.5-7.
- [42] Andrea Vedaldi and Andrew Zisserman. Efficient additive kernels via explicit feature maps. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(3):480–492, 2012.
- [43] Sreekanth Vempati, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Generalized RBF feature maps for efficient detection. In *BMVC*, pages 1–11, 2010.
- [44] Grace Wahba. *Spline Models for Observational Data*, volume 59. SIAM, 1990.
- [45] Christian J. Walder and Adrian N. Bishop. Fast Bayesian intensity estimation for the permanent process. In *International Conference on Machine Learning*, volume 70, pages 3579–3588. PMLR, 2017.
- [46] Ping Wang, Yan Li, and Chandan K. Reddy. Machine learning for survival analysis: A survey. *ACM Computing Surveys (CSUR)*, 51(6):1–36, 2019.
- [47] Shannon Wongvibulsin, Katherine C. Wu, and Scott L. Zeger. Clinical risk prediction with random forests for survival, longitudinal, and multivariate (RF-SLAM) data analysis. *BMC Medical Research Methodology*, 20(1):1–14, 2020.
- [48] Weichi Yao, Halina Frydman, Denis Larocque, and Jeffrey S. Simonoff. Ensemble methods for survival function estimation with time-varying covariates. *Statistical Methods in Medical Research*, 31(11):2217–2236, 2022.
- [49] Qixian Zhong, Jonas W. Mueller, and Jane-Ling Wang. Deep extended hazard models for survival analysis. In *Advances in Neural Information Processing Systems 34*, 2021.

A Derivation of MAP Estimator

We detail the derivation of the MAP estimator (12). The functional derivative of $S(x(\mathbf{y}), \underline{x}(\mathbf{y}))$ should be zero on MAP estimator $\hat{x}(\mathbf{y})$:

$$\begin{aligned} \delta S(\hat{x}(\mathbf{y}), \underline{\hat{x}}(\mathbf{y})) &= \int_{\mathcal{Y}} \left[\frac{\delta S}{\delta \hat{x}(\mathbf{y})} \delta x(\mathbf{y}) + \frac{\delta S}{\delta \underline{\hat{x}}(\mathbf{y})} \delta \underline{x}(\mathbf{y}) \right] d\mathbf{y} + O((\delta x)^2) \\ &\simeq \int_{\mathcal{Y}} \left[2 \sum_{j=1}^J \Delta_j \hat{x}(\mathbf{y}_j) \delta(\mathbf{y} - \mathbf{y}_j) - \sum_{n=1}^N \frac{2}{\hat{x}(\mathbf{y})} \delta(\mathbf{y} - \tilde{\mathbf{y}}_n) + \frac{1}{2} \hat{x}(\mathbf{y}) \right] \delta x d\mathbf{y} \\ &\quad + \int_{\mathcal{Y}} \frac{1}{2} \hat{x}(\mathbf{y}) \delta \underline{x} d\mathbf{y} \\ &= \int_{\mathcal{Y}} \left[2 \sum_{j=1}^J \Delta_j \hat{x}(\mathbf{y}) \delta(\mathbf{y} - \mathbf{y}_j) - \sum_{n=1}^N \frac{2}{\hat{x}(\mathbf{y})} \delta(\mathbf{y} - \tilde{\mathbf{y}}_n) + \hat{x}(\mathbf{y}) \right] \delta x d\mathbf{y} = 0, \end{aligned}$$

where the following relation was used,

$$\begin{aligned} \int_{\mathcal{Y}} \hat{x}(\mathbf{y}) \delta \underline{x} d\mathbf{y} &= \int_{\mathcal{Y}} \hat{x}(\mathbf{y}) \int_{\mathcal{Y}} k^*(\mathbf{y}, \mathbf{y}') \delta x(\mathbf{y}') d\mathbf{y}' d\mathbf{y} \\ &= \int_{\mathcal{Y}} d\mathbf{y}' \delta x(\mathbf{y}') \int_{\mathcal{Y}} k^*(\mathbf{y}, \mathbf{y}') \hat{x}(\mathbf{y}) d\mathbf{y} \\ &= \int_{\mathcal{Y}} \hat{x}(\mathbf{y}') \delta x d\mathbf{y}'. \quad (\because k^*(\mathbf{y}, \mathbf{y}') = k^*(\mathbf{y}', \mathbf{y})) \end{aligned}$$

Thus the following equation is derived,

$$\hat{x}(\mathbf{y}) + 2 \sum_{j=1}^J \Delta_j \hat{x}(\mathbf{y}_j) \delta(\mathbf{y} - \mathbf{y}_j) = \sum_{n=1}^N \frac{2}{\hat{x}(\tilde{\mathbf{y}}_n)} \delta(\mathbf{y} - \tilde{\mathbf{y}}_n), \quad \mathbf{y} \in \mathcal{Y}. \quad (\text{A1})$$

By applying operator \mathcal{K} to (A1), we obtain a linear integral equation that derives the MAP estimator $\hat{x}(\mathbf{y})$ as follows,

$$\hat{x}(\mathbf{y}) + 2 \sum_{j=1}^J \Delta_j \hat{x}(\mathbf{y}_j) k(\mathbf{y}, \mathbf{y}_j) = 2 \sum_{n=1}^N k(\mathbf{y}, \tilde{\mathbf{y}}_n) \hat{x}(\tilde{\mathbf{y}}_n)^{-1}, \quad \mathbf{y} \in \mathcal{Y}. \quad (\text{A2})$$

The linearity of the integral equation permits a representation of the form

$$\hat{x}(\mathbf{y}) = 2 \sum_{n=1}^N h(\mathbf{y}, \tilde{\mathbf{y}}_n) \hat{x}(\tilde{\mathbf{y}}_n)^{-1},$$

where $h(\mathbf{y}, \mathbf{y}')$ is a positive semi-definite kernel that solves integral equation (13). Note that the derivations follow Kim [23].

B Derivation of Predictive Distribution and Marginal Likelihood

We now know the mode of the posterior, $\hat{x}(\mathbf{y})$, and consider a Taylor expansion of functional action potential $S(x(\mathbf{y}), \underline{x}(\mathbf{y}))$ centered on the mode such that

$$S(x(\mathbf{y}), \underline{x}(\mathbf{y})) \simeq S(\hat{x}(\mathbf{y}), \underline{\hat{x}}(\mathbf{y})) + \frac{1}{2} \iint_{\mathcal{Y} \times \mathcal{Y}} \sigma^*(\mathbf{y}, \mathbf{y}') (x(\mathbf{y}) - \hat{x}(\mathbf{y})) (x(\mathbf{y}') - \hat{x}(\mathbf{y}')) d\mathbf{y} d\mathbf{y}', \quad (\text{B1})$$

where $\sigma^*(\mathbf{y}, \mathbf{y}') = \left. \frac{\delta^2 S(x, \underline{x})}{\delta x(\mathbf{y}) \delta x(\mathbf{y}')} \right|_{x=\hat{x}}$ is the second derivative of S . The first term in the Taylor expansion vanishes due to the stationary condition. The quadratic approximation of the action integral corresponds to the approximation of the posterior process by a GP, and the predictive covariance or the kernel function for the posterior GP, denoted by $\sigma(\mathbf{y}, \mathbf{y}')$, can be obtained by the functional inversion of $\sigma^*(\mathbf{y}, \mathbf{y}')$, which results in

$$\sigma(\mathbf{y}, \mathbf{y}') = h(\mathbf{y}, \mathbf{y}') - \mathbf{h}(\mathbf{y})^\top (\mathbf{Z} + \mathbf{H})^{-1} \mathbf{h}(\mathbf{y}'), \quad (\text{B2})$$

where the definitions of \mathbf{Z} and \mathbf{H} are taken from (20). The full derivation of (B2) is given in [24]. When the latent function $x(\mathbf{y})$ follows a posterior GP with mean of $\hat{x}(\mathbf{y})$ and kernel $\sigma(\mathbf{y}, \mathbf{y}')$, it is easily verified that the value of the squared function, $\lambda = x^2(\mathbf{y})$, on each point of the covariate domain $\mathbf{y} \in \mathcal{Y}$ follows a Gamma distribution defined by (21).

Furthermore, under Laplace approximation (B1), we can obtain the marginal likelihood, $p(\mathcal{D})$, in (10) by performing the path integral as,

$$\log p(\mathcal{D}) = \log \int \exp[-S(x(\mathbf{y}), \underline{x}(\mathbf{y})) - \frac{1}{2} \log |\mathcal{K}|] \mathcal{D}x \simeq -S(\hat{x}(\mathbf{y}), \hat{\underline{x}}(\mathbf{y})) + \frac{1}{2} \log \frac{|\Sigma|}{|\mathcal{K}|}, \quad (\text{B3})$$

where $|\Sigma|$ is the functional determinant of integral operator $\Sigma = \int_{\mathcal{Y}} \cdot \sigma(\mathbf{y}, \mathbf{y}') d\mathbf{y}'$. We can rewrite the result in a more tractable form by substituting (11, B2) into (B3):

$$\log p(\mathcal{D}) = \log |\mathbf{Z}| - \frac{1}{2} \log |\mathbf{I}_N + \mathbf{Z}^{-1} \mathbf{H}| - \frac{1}{2} \log |\mathbf{I}_M + 2\mathbf{A}|^{-1} + (\log 2 - 1)N.$$

Full derivations are provided in [24].

C Model Configuration

C.1 Synthetic Data

Survival Permanental Process (SurvPP)

We set the number of features for Random feature map [37] (M), learning parameter (lr), and stop condition (G) for Adam [26] as follows:

$$M = 500, \quad lr = 0.05, \quad G < 10^{-5}.$$

We applied to SurvPP a multiplicative Gaussian kernel

$$k(\mathbf{y}, \mathbf{y}') = \prod_{d=1}^3 e^{-(\theta(y_d - y'_d))^2}, \quad \mathbf{y} = (t, y_1, y_2),$$

where hyper-parameter θ was optimized for each data by maximizing the marginal likelihood through grid search. In the experiments on synthetic data, we selected a set of nine values for θ as the grid points,

$$\theta \in \{0.1, 0.2, 0.5, 0.7, 1.0, 2.0, 5.0, 7.0, 10.0\}.$$

We implemented SurvPP by using TensorFlow-2.10. A MacBook Pro with 12-core CPU (Apple M2 Max) was used, with the GPU inactivated (`tf.device('/cpu:0')`) for a fair comparison with the benchmarks.

Cox Proportional Hazards Model (CoxPH)

We implemented CoxPH through package `survival.coxph` (LGPL-3) [41]. The calls in `coxph` to fit a model and compute a base hazard function were

```
> cfit = coxph(Surv(Start, Stop, Event) ~ cov1 + cov2, df)
> sfit = survfit(cfit, list(cov1 = 0, cov2 = 0)),
```

where `df` was the survival data in counting process format.

Generalized Boosted Model (GBM)

We implemented GBM through package `gbm3.gbmt` [16] (GPL). The call in `gbmt` to fit a model was

```
> gfit = gbmt(Surv(Start, Stop, Event) ~ cov1 + cov2, data = df,
              distribution = gbm_dist("CoxPH"),
              cv_folds = 10, train_params = params,
              par_details = gbmParallel(num_threads = 12)),
```

where `params` represents the hyperparameter. We selected a set of nine hyperparameters for the grid search,

$$\text{num_trees} \in \{500, 1000, 2000\} \times \text{shrinkage} \in \{0.001, 0.005, 0.01\},$$

and found the one that minimized the cross validation error (`gfit$valid.error`), where `num_trees` and `shrinkage` represent the number of trees and the shrinkage/learning rate, respectively.

Random Forest-based Model (RFM)

We implemented RFM through package `LTRCforests` (GPL) [48]. The call to fit a model was

```
> rfit = ltrcrrf(Surv(Start, Stop, Event) ~ cov1 + cov2, data = df,
                id = ID, mtry = ceiling(10), ntree = 100).
```

C.2 Real-world Data: PBC

Survival Permanental Process (SurvPP)

We set the number of features for Random feature map [37] (M), learning parameter (lr), and stop condition (G) for Adam [26] as follows:

$$M = 500, \quad lr = 50, \quad G < 10^{-5}.$$

We applied to `SurvPP` a multiplicative Gaussian kernel

$$k(\mathbf{y}, \mathbf{y}') = \prod_{d=1}^{13} e^{-\theta(y_d - y'_d)^2}, \quad \mathbf{y} = (t, y_1, \dots, y_{12}),$$

where the hyper-parameter θ was optimized for each data by maximizing the marginal likelihood through grid search. In the experiments on synthetic data, we selected a set of nine values for θ as the grid points,

$$\theta \in \{0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09\}.$$

Here, we normalized the 13 covariates so that $y_d \rightarrow 0.1(y_d - \text{mean}[y_d])/\text{std}[y_d]$.

As in the experiments on synthetic data, GPU was set off (`tf.device('/cpu:0')`) for a fair comparison with the benchmarks.

Cox Proportional Hazards Model (CoxPH)

The calls in `coxph` to fit a model and compute a base hazard function were

```
> cfit = coxph(Surv(Start, Stop, Event) ~ age + edema + alk.phos + chol + ast
              +platelet + spiders + hepato + ascites + albumin + bili + protime, df)
> sfit = survfit(cfit, list(age = 0, edema = 0, alk.phos = 0, chol = 0, ast = 0,
                           platelet = 0, spiders = 0, hepato = 0, ascites = 0, albumin = 0,
                           bili = 0, protime = 0)),
```

where `df` was the survival data in counting process format.

Generalized Boosted Model (GBM)

The call in `gbmt` to fit a model was

```
> gfit = gbmt(Surv(Start, Stop, Event) ~ age + edema + alk.phos + chol + ast
              +platelet + spiders + hepato + ascites + albumin + bili
              +protime, data = df,
              distribution = gbm_dist("CoxPH"),
              cv_folds = 10, train_params = params,
              par_details = gbmParallel(num_threads = 12)),
```

where `params` represents the hyperparameter. We selected a set of nine hyperparameters for the grid search,

$$\text{num_trees} \in \{500, 1000, 2000\} \times \text{shrinkage} \in \{0.001, 0.005, 0.01\},$$

and found the one that minimized the cross validation error (`gfit$valid.error`), where `num_trees` and `shrinkage` represent the number of trees and the shrinkage/learning rate, respectively.

Random Forest-based Model (RFM)

The call to fit a model was

```
> rfit = ltrcrf(Surv(Start, Stop, Event) ~ age + edema + alk.phos + chol + ast
               +platelet + spiders + hepato + ascites + albumin + bili
               +protime, data = df, id = ID, stepFactor = 1.5).
```

C.3 Real-world Data: SANAD

Survival Permanent Process (SurvPP)

We set the number of features for Random feature map [37] (M), learning parameter (lr), and stop condition (G) for Adam [26] as follows:

$$M = 500, \quad lr = 10, \quad G < 10^{-5},$$

We applied to SurvPP a multiplicative Gaussian kernel

$$k(\mathbf{y}, \mathbf{y}') = 0.1 \prod_{d=1}^5 e^{-\theta(y_d - y'_d)^2}, \quad \mathbf{y} = (t, y_1, \dots, y_4),$$

where hyper-parameter θ was optimized for each data by maximizing the marginal likelihood through grid search. In the experiments on synthetic data, we selected a set of nine values for θ as the grid points,

$$\theta \in \{0.6, 0.7, 0.8, 0.9, 1.0, 1.1, 1.2, 1.3, 1.4\}.$$

Here, we normalized the 5 covariates so that $y_d \rightarrow 0.1(y_d - \text{mean}[y_d])/\text{std}[y_d]$.

As in the experiments on synthetic data, GPU was set off (`tf.device('/cpu:0')`) for a fair comparison with the benchmarks.

Cox Proportional Hazards Model (CoxPH)

The calls in `coxph` to fit a model and compute a base hazard function were

```
> cfit = coxph(Surv(Start, Stop, Event) ~ age + gender + treat + dose, df)
> sfit = survfit(cfit, list(age = 0, gender = 0, treat = 0, dose = 0)),
```

where `df` was the survival data in counting process format.

Generalized Boosted Model (GBM)

The call in `gbmt` to fit a model was

```
> gfit = gbmt(Surv(Start, Stop, Event) ~ age + gender + treat + dose,
              data = df, distribution = gbm_dist("CoxPH"),
              cv_folds = 10, train_params = params,
              par_details = gbmParallel(num_threads = 12)),
```

where `params` represents the hyperparameter. We selected a set of nine hyperparameters for the grid search,

$$\text{num_trees} \in \{500, 1000, 2000\} \times \text{shrinkage} \in \{0.001, 0.005, 0.01\},$$

and found the one that minimized the cross validation error (`gfit$valid.error`), where `num_trees` and `shrinkage` represent the number of trees and the shrinkage/learning rate, respectively.

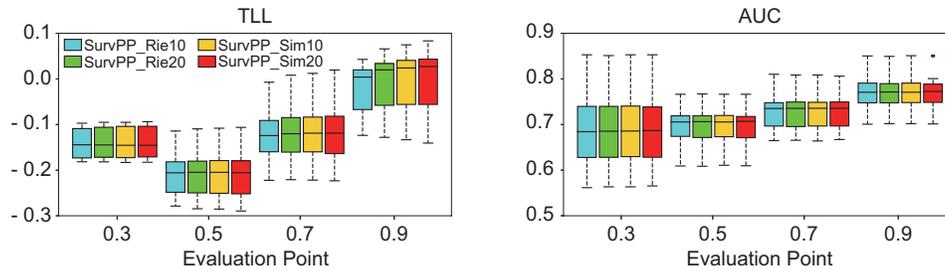


Figure D1: Performances on the dataset of non-linear hazard function $\lambda_{non}(t)$ for different numerical integration rules. SurvPP_Rie# and SurvPP_Sim# represent SurvPPs with Riemann sum and Simpson's rule approximations of integration, respectively. # is the number of representative time points for covariate, J_u .

Random Forest-based Model (RFM)

The call to fit a model was

```
> rfit = ltrcrf(Surv(Start, Stop, Event) ~ age + gender + treat + dose,
               data = df, id = ID, stepFactor = 1.5).
```

D Additional Experiments

D.1 More Accurate Evaluation of Equivalent Kernel

The functional form of covariate map, $\mathbf{y}(t)$, involves the sum of outer products of feature maps, \mathbf{A} , defined by Equation (16), and \mathbf{A} can be considered as a Riemann sum approximation of the time integral as follows:

$$\mathbf{A} = \sum_{u=1}^U \int_0^{T_u} \phi(\mathbf{y}_u(t)) \phi(\mathbf{y}_u(t))^\top dt \sim \sum_{j=1}^J \Delta_j \phi(\mathbf{y}_j) \phi(\mathbf{y}_j)^\top. \quad (\text{D1})$$

Thus, when the covariate map, $\mathbf{y}_u(t)$, and the feature map, $\phi(\mathbf{y})$, are both smooth, the Riemann sum approximation can be replaced by a more accurate approximation,

$$\mathbf{A} = \sum_{u=1}^U \int_0^{T_u} \phi(\mathbf{y}_u(t)) \phi(\mathbf{y}_u(t))^\top dt \sim \sum_{j=1}^J w_j \phi(\mathbf{y}_j) \phi(\mathbf{y}_j)^\top, \quad (\text{D2})$$

where the weights, w_j , depend on the approximation rule. This time, we adopt Simpson's 1/3 rule, and add an experiment on the synthetic nonlinear data, $\lambda_{non}(t)$, to check for an improvement in accuracy. Note that Simpson's 1/3 rule corresponds to quadratic interpolation. Figure D1 shows that adopting Simpson's rule achieved better predictive performance on TLL with a smaller discretization number (J_u).

D.2 Experiment on Independent Validation Data Sets

In Section 4, we evaluated the performances on synthetic data with the cross-validation approach because it is common in the machine learning literature. But synthetic data offers the unique advantage of being able to generate as much data as needed, and we here conducted the performance assessment with the independent validation approach. Figure D2 displays the result on independent validation sets of synthetic data ($U = 900$ for each of 10 training data, and $U = 100$ for each of 10 test data), showing the result similar to the cross-validation approach (Figure 1-2).

D.3 Experiment on Larger Data Sets

To examine the computation scalability of the compared models versus the number of observed events N , we created data sets with user size $U \in \{10^3, 10^4, 5 \cdot 10^4, 10^5\}$ according to the nonlinear

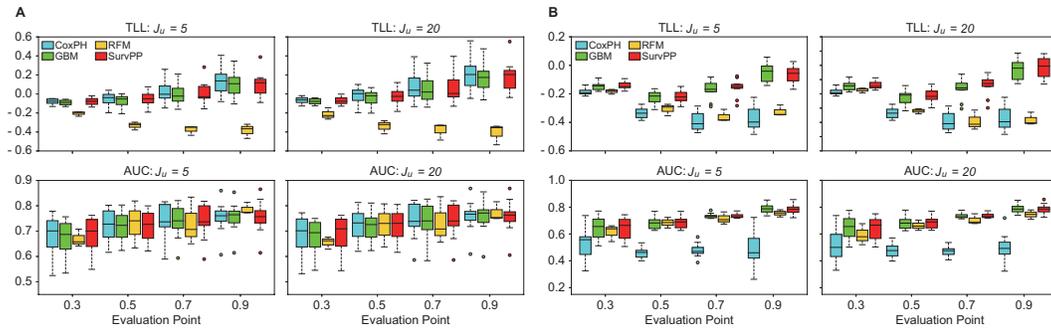


Figure D2: Performance on the independent validation synthetic datasets. (A) Log-linear hazard function $\lambda_{lin}(t)$. (B) Non-linear hazard function $\lambda_{non}(t)$. The results are consistent with Figure 1-2

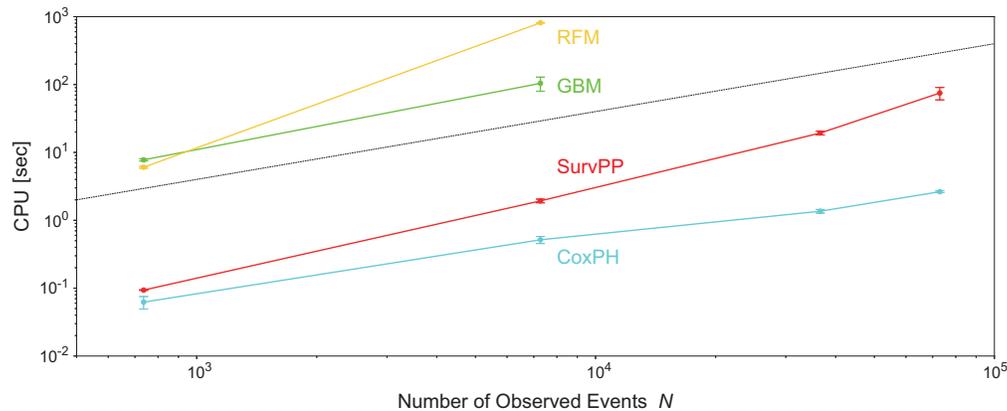


Figure D3: The CPU times demanded for estimating a hazard function versus the number of observed events. The error bars represent the standard deviations across 10 trials. The dashed line represents a line of $\text{CPU} \propto N$ as reference. For GBM and SurvPP, the average cpu times over 9-point grid search of the hyperparameter are displayed. The CPU times of GBM and RFM exceeded 10^3 seconds with $N > 10^4$, and the estimations were given up.

scenario (see Section 4.1),

$$\lambda_{non}(t) = h(t) \exp[2 - 5(y_1^2(t) + y_2^2(t))], \quad h(t) = 2 \cdot t^{3/2},$$

which resulted in the data sets with $N \in \{818, 8082, 40660, 81066\}$. Here, we set J_u to 10 for the counting process format of data. For each dataset, we randomly split the U individuals into 10 subgroups, repeated assigning 9 subgroups to training data, and conducted 10 trials of evaluations of the CPU times demanded for estimating a hazard function. Figure D3 displays the CPU times as function of N of training data. It shows that the CoxPH computation clearly scaled linearly with N , while that of SurvPP seems to be a little more than linear. This is because that each iteration of gradient descent algorithm scaled linearly with N , but the number of iterations to meet stop condition $G < 10^{-5}$ increased moderately with N . Among the non-parametric approaches (SurvPP, GBM, and RFM), SurvPP achieved the fastest computation at hundreds of times faster than the others, regardless of the number of observed events N .

D.4 Sensitivity to Hyper-Parameters

Kernel methods/GP models are generally sensitive to kernel parameter values, and a effective way of optimizing the kernel parameter is essential. Figure D4 plots AUC performance on PBC dataset for SurvPPs with different search ranges of kernel parameters: the range of SurvPP_0 is $\theta \in \{0.1, 0.2, 0.5, 0.7, 1.0, 2.0, 5.0, 7.0, 10.0\}$, which was used in the supplement; the range of SurvPP_1 is $\theta \in \{0.01, 0.02, 0.03, \dots, 0.09\}$. Figure D4 shows that SurvPP_1 achieved substantially

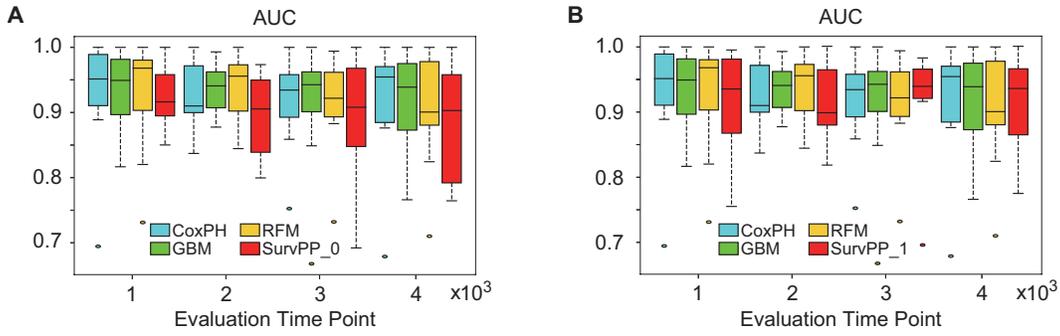


Figure D4: Performances on PBC dataset for SurvPPs with different search ranges of kernel parameters. (A) The search range of SurvPP₀ is $\theta \in \{0.1, 0.2, 0.5, 0.7, 1.0, 2.0, 5.0, 7.0, 10.0\}$. (B) The search range of SurvPP₁ is $\theta \in \{0.01, 0.02, 0.03, \dots, 0.09\}$.

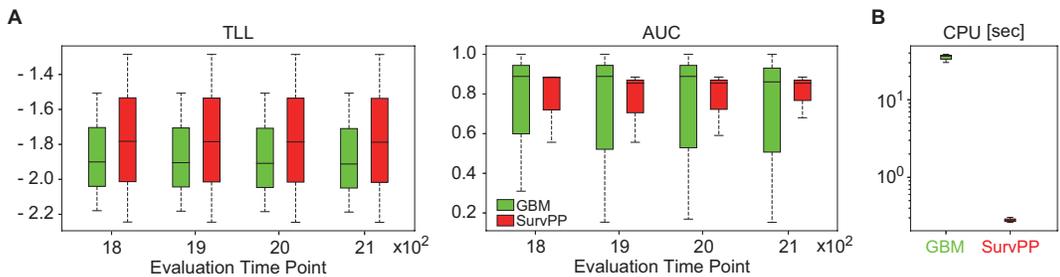


Figure D5: Performances on MIMIC III dataset. (A) Box plot of TLL and AUC as functions of evaluation point: the higher, the better. (B) The CPU times demanded for estimating a hazard function. The error bars represent the standard deviations. The average cpu times over 9-point grid search of the hyperparameter are displayed.

better AUC performance than SurvPP₀, where SurvPP₁ is displayed in Figure 3. A more sophisticated algorithm for hyperparameter optimization could enhance the performance of SurvPP, which is the next step in our study.

D.5 Experiment on ICU Data Set

We examined the validity of SurvPP against GBM on *MIMIC-III Clinical Database* (MIMIC III), the large publicly available dataset of over 50,000 ICU admissions from the Beth Israel Deaconess Medical Center [19], where events were deaths. We extracted admissions and measurements from CareVue (ITEMID) such that admissions shared 14 measured covariates with each other, which resulted in 133 admissions. We adopted sex as a static covariate and the 14 measurements as time-varying covariates. We randomly split the individuals into 3 subgroups, assigned one to test and the others to training data, and conducted 3-fold cross evaluation of the predictive performances. The model configuration follows the experiment on PBC data set (see Appendix C).

Figure D5 displays the predictive performance on MIMIC III. It shows that SurvPP achieved better TLL performance than GBM, and SurvPP achieved comparative AUC performance with GBM. The smaller variance of SurvPP's AUC implies that SurvPP works more robustly than GBM, but further experiments are needed to investigate it.