

---

# Supplementary Material

## “Survival Permanental Processes for Survival Analysis with Time-Varying Covariates”

---

Anonymous Author(s)

Affiliation

Address

email

### 1 S1 Derivation of MAP Estimator

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

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

4 where the following relation was used,

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

5 Thus the following equation is derived,

$$\underline{\hat{x}}(\mathbf{y}) + 2 \sum_{j=1}^J \Delta_j \hat{x}(\mathbf{y}_j) \delta(\mathbf{y} - \mathbf{y}_j) = \sum_{n=1}^N \frac{2}{\hat{x}(\tilde{\mathbf{y}}_n)} \delta(\mathbf{y} - \tilde{\mathbf{y}}_n), \quad \mathbf{y} \in \mathcal{Y}. \quad (\text{S1})$$

6 By applying operator  $\mathcal{K}$  to (S1), we obtain a linear integral equation that derives the MAP estimator  
 7  $\hat{x}(\mathbf{y})$  as follows,

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

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

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

9 where  $h(\mathbf{y}, \mathbf{y}')$  is a positive semi-definite kernel that solves the following equation,

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

10 Note that the derivations follow Kim [2].

## 11 **S2 Model Configuration**

### 12 **S2.1 Synthetic Data**

#### 13 **Survival Permanental Process (SurvPP)**

14 For all the experiments in Section 4, we set the number of features for Random feature map [4] ( $M$ ),  
15 the learning parameter ( $lr$ ), and the stop condition ( $G$ ) for Adam [3] as follows:

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

16 We applied to SurvPP a multiplicative Gaussian kernel

$$k(\mathbf{y}, \mathbf{y}') = \prod_{d=1}^2 e^{-(\theta(y_d - y'_d))^2}, \quad \mathbf{y} = (y_1, y_2),$$

17 where the hyper-parameter  $\theta$  was optimized for each data by maximizing the marginal likelihood  
18 through grid search. In the experiments with synthetic data, we selected a set of nine values for  $\theta$   
19 for the grid points,

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

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

#### 23 **Cox Proportional Hazards Model (CoxPH)**

24 We implemented CoxPH through package `survival.coxph` in R-4.2.3 (LGPL-3) [5]. The calls in  
25 `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)),
```

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

#### 27 **Generalized Boosted Model (GBM)**

28 We implemented GBM through package `gbm3.gbmt` (GPL) <sup>1</sup> in R-4.2.3. The call in `gbmt` to fit a  
29 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)),
```

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

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

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

---

<sup>1</sup><https://rdr.io/github/gbm-developers/gbm3/>

## 34 **Random Forest-based Model (RFM)**

35 We implemented RFM through package LTRCforests in R-4.2.3 (GPL) [6]. The call to fit a model  
36 was

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

## 37 **S2.2 Real-world Data**

### 38 **Survival Permanental Process (SurvPP)**

39 For the experiments in Section S4, we set the number of features for Random feature map [4] ( $M$ ),  
40 the learning parameter ( $lr$ ), and the stop condition ( $G$ ) for Adam [3] as follows:

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

41 We applied to SurvPP a multiplicative Gaussian kernel

$$k(\mathbf{y}, \mathbf{y}') = \prod_{d=1}^{13} e^{-(\theta(y_d - y'_d))^2}, \quad \mathbf{y} = (y_1, \dots, y_{13}),$$

42 where the hyper-parameter  $\theta$  was optimized for each data by maximizing the marginal likelihood  
43 through grid search. In the experiments with synthetic data, we selected a set of nine values for  $\theta$   
44 for the grid points,

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

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

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

### 48 **Cox Proportional Hazards Model (CoxPH)**

49 We implemented CoxPH through package `survival.coxph` in R-4.2.3 (LGPL-3) [5]. The calls in  
50 `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)),
```

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

### 52 **Generalized Boosted Model (GBM)**

53 We implemented GBM through package `gbm3.gbmt` (GPL) <sup>2</sup> in R-4.2.3. The call in `gbmt` to fit a  
54 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)),
```

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

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

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

<sup>2</sup><https://rdrr.io/github/gbm-developers/gbm3/>

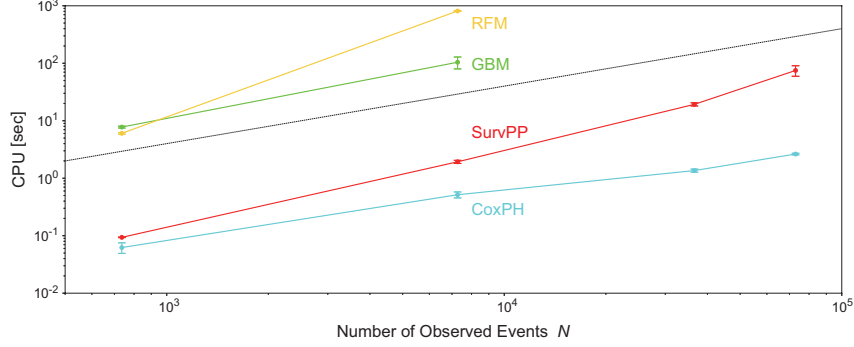


Figure S1: The CPU times demanded for estimating a hazard function regarding the number of observed events. The error bars represent the standard deviations across 10 trials. The dashed line represents a line of  $CPU \propto N$  as reference. For GBM and SurvPP, the average cpu times over 9-points grid 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.

### 59 Random Forest-based Model (RFM)

60 We implemented RFM through package LTRCforests in R-4.2.3 (GPL) [6]. The call to fit a model  
61 was

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

### 62 S3 Experiment on Larger Data Sets

63 To examine the computation scalability of the compared models regarding the number of observed  
64 events  $N$ , we created data sets with the user size  $U \in \{10^3, 10^4, 5 \cdot 10^4, 10^5\}$  according to the  
65 nonlinear scenario (see Section 4.1),

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

66 which resulted in the data sets with  $N \in \{818, 8082, 40660, 81066\}$ . Here, we set  $J_u$  as 10 for  
67 the counting process format of data. For each dataset, we randomly split the  $U$  individuals into 10  
68 subgroups, repeated assigning 9 subgroups to training data, and conducted 10 trials of evaluations  
69 of the CPU times demanded for estimating a hazard function. Figure S1 displays the CPU times  
70 as function of  $N$  of training data. It shows that the computation of CoxPH scaled linearly with  $N$   
71 clearly, while that of SurvPP seems to be a little more than linear. It is because that each iteration  
72 of gradient descent algorithm scaled linearly with  $N$ , but the number of iterations to meet a stop  
73 condition  $G < 10^{-5}$  increased moderately with  $N$ . Among the non-parametric approaches (SurvPP,  
74 GBM, and RFM), our SurvPP achieved the fastest computation, which was hundreds of times faster  
75 than the others regardless of the number of observed events  $N$ .

### 76 S4 Experiment on Real-world Data

77 We examined the validity of SurvPP against the benchmark models on real-world survival data set,  
78 *Mayo Clinic Primary Biliary Cholangitis Data*, provided by R package survival (LGPL-3) [5].  
79 312 patients with primary biliary cirrhosis (PBC) were enrolled in a randomized medical trial at  
80 the Mayo Clinic between 1974 and 1984 [1], where events were the time of death. In the data  
81 set, twelve time-varying covariates were measured at entry and at yearly intervals, which include  
82 age at entry (age), alkaline phosphatase (alk.phos), logarithm of serum albumin (albumin), pres-  
83 ence of ascites (ascites), aspartate aminotransferase (ast), logarithm of serum bilirubin (bili),  
84 serum cholesterol (chol), condition of edema (edema), presence of hepatomegaly or enlarged liver  
85 (hepato), platelet count (platelet), logarithm of prothrombin time (protime) and presence or  
86 absence of spiders (spiders). As with the experiments on synthetic data, we randomly split the

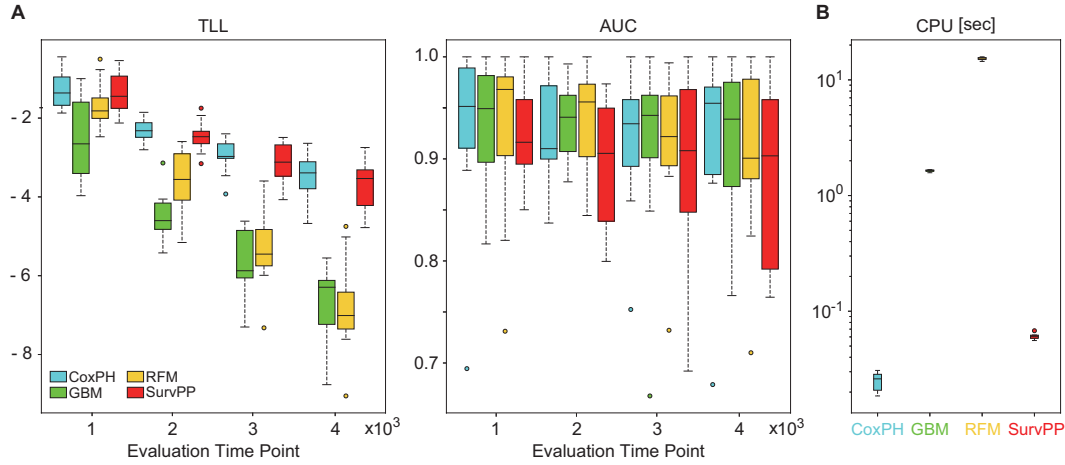


Figure S2: Performances on the real-world dataset. (A) Box plot of TLL and AUC as functions of evaluation time points: the higher, the better. (B) Box plot of CPU times demanded for estimating a hazard function: the lower, the better.

312 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 S2 displays the predictive performances on the real-world dataset, showing that all compared models achieved comparable performances at AUC, while SurvPP and CoxPH outperformed significantly at TLL against the other non-parametric models. We confirmed that among the non-parametric models, our SurvPP was the best on both the predictive accuracy and the computational time.

## References

- [1] 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.
- [2] Hideaki Kim. Fast Bayesian inference for Gaussian Cox processes via path integral formulation. In *Advances in Neural Information Processing Systems 34*, 2021.
- [3] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [4] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems 20*, 2007.
- [5] Terry M. Therneau and Thomas Lumley. Package ‘survival’. *R Top Doc*, 128(10):28–33, 2015.
- [6] Weichi Yao, Halina Frydman, Denis Larocque, and Jeffrey S. Simonoff. Ensemble methods for survival data with time-varying covariates. *arXiv preprint arXiv:2006.00567*, 2020.