

Yang and Prentice model with piecewise exponential baseline distribution for modeling lifetime data with crossing survival curves

Fábio N. Demarqui and Vinícius D. Mayrink

Universidade Federal de Minas Gerais

Abstract. Proportional hazards (PH), proportional odds (PO) and accelerated failure time (AFT) models have been widely used to deal with survival data in different fields of knowledge. Despite their popularity, such models are not suitable to handle survival data with crossing survival curves. Yang and Prentice (2005) proposed a semiparametric two-sample approach, denoted here as the YP model, allowing the analysis of crossing survival curves and including the PH and PO configurations as particular cases. In a general regression setting, the present work proposes a fully likelihood-based approach to fit the YP model. The main idea is to model the baseline hazard via the piecewise exponential (PE) distribution. The approach shares the flexibility of the semiparametric models and the tractability of the parametric representations. An extensive simulation study is developed to evaluate the performance of the proposed model. We demonstrate how useful is the new method through the analysis of survival times related to patients enrolled in a cancer clinical trial. Finally, an R package called YPPE was developed to fit the proposed model. The simulation results indicate that our model performs well for moderate sample sizes in the general regression setting. A superior performance is also observed with respect to the original YP model designed for the two-sample scenario.

1 Introduction

Proportional hazards (PH) models have played a central role in the analysis of survival data. Such class of models provides a very flexible framework to model survival data. They further allow an easy interpretation of the parameters from the practical point of view. The main assumption of the PH models is that the hazard ratios are constant over time. When such assumption is not verified by the data, some alternatives such as the proportional odds (PO) and the accelerated failure time (AFT) models can be used in the analysis. However, none of them is suitable to accommodate survival data with crossing survival curves. This type of problem is often related to studies involving treatment and control groups. The survival function for one group may have a fast decay in contrast with a slow decay for the other. The curves tend to intersect at some time point configuring an inversion in terms of which group is on the top/bottom position. Studying this alteration is relevant in many clinical trials, where the identification of the crossing time indicates when the target treatment for a disease can be considered effective.

Survival data with crossing survival curves may arise due to several reasons in practice. For instance, Diao, Zeng and Yang (2013) indicates that this may occur in certain clinical trials related to aggressive treatments such as surgery. Some adverse effects can be observed in an initial stage, but beneficial results may appear in the long run. According to Breslow (1974), another situation connected with crossing survival functions is when a treatment has an early and quick effect and it becomes similar to or worse than the placebo treatment after a certain period.

Key words and phrases. Survival analysis, short-term and long-term hazard ratios, semiparametric modeling, maximum likelihood estimation.

Received February 2019; accepted March 2020.

Several approaches have been proposed in the literature to accommodate this crossing feature in survival data. The most popular ones are based on time-varying regression coefficients; see, for example, the references [Egge and Zahl \(1999\)](#), [Shyur, Elsayed and Luxhøj \(1999\)](#) and [Putter et al. \(2005\)](#). [Zeng and Lin \(2007\)](#) proposes a class of transformation models for counting processes which encompasses linear transformation models and which can handle crossing hazards. Alternatively, [Yang and Prentice \(2005\)](#) presented a semiparametric two-sample model (hereafter denoted as YP model) for this type of problem. The feature “two-sample” refers to the scenario where, for example, there is a treatment group and a control group that can be conveniently represented through a binary variable. The YP proposal is an interesting option, since it includes the PH and PO representations as particular cases. In their model, the baseline hazard function is left unspecified, in fact a counting process is assumed leading to a survival step function. A pair of short-term and long-term hazard ratio parameters is included to accommodate crossing survival curves. In addition, a pseudo maximum likelihood approach is considered for the estimation procedure. Consistency and asymptotic normality of the resulting estimators are demonstrated in the paper.

[Yang and Prentice \(2011\)](#) extended the estimation procedure in [Yang and Prentice \(2005\)](#) to pointwise and simultaneous inference on the hazard ratio function itself. They further proved the consistency and asymptotic normality of the estimates at a fixed time point. [Yang and Zhao \(2012\)](#) proposed two omnibus tests to evaluate the adequacy of the YP model. The first test is based on the martingale residuals and the second one examines the contrast between the non-parametric and model-based estimators of the survival function. [Diao, Zeng and Yang \(2013\)](#) extended the two-sample YP model to a general regression setting with possibly time-dependent covariates; the study developed an efficient likelihood-based estimation procedure. The authors also demonstrated the consistency, asymptotic normality and efficiency of the resulting estimators. [Nieto-Barajas \(2014\)](#) also extended the YP model to accommodate a general regression setting, and proposed a Bayesian nonparametric prior, based on increasing additive processes mixtures, to model the baseline function. [Yang and Prentice \(2015\)](#) presented an alternative formulation of the YP model introduced in [Yang and Prentice \(2005\)](#) by allowing a subset of the explanatory variables to have constant effects over time, *that is*, preserving the proportional hazard structure. The YP model has also been extended by [Tong, Zhu and Sun \(2007\)](#) to accommodate current status survival data. Another extension is found in [Zhang, Wang and Sun \(2018\)](#) to fit case II interval-censored data.

The use of semiparametric methods for univariate survival data started with [Cox \(1972\)](#) on the proportional hazards model. [Breslow \(1972\)](#) and [Breslow \(1974\)](#) are two initial publications proposing the use of the piecewise exponential (PE) distribution to replace the baseline hazard in a survival analysis. The grid configuration for a model with the PE baseline hazard is a central topic in [Kalbfleisch and Prentice \(1973\)](#). Many applications, related to clinical trials and involving the PE distribution, can be found in the literature; some few examples are: leukemia ([Breslow, 1974](#)), gastric cancer ([Gamerman, 1991](#)), kidney infection ([Sahu et al., 1997](#), [Ibrahim, Chen and Sinha, 2001](#)), breast cancer ([Sinha, Chen and Ghosh, 1999](#)), melanoma ([Demarqui et al., 2014](#)) and hospital mortality ([Clark and Ryan, 2002](#)). Although parametric in a strict sense, the model with the PE baseline hazard has a strong nonparametric appeal. The main reason is the fact that assumptions about the shape of the baseline hazard are not required in this approach.

The main contribution of the present paper is to propose a novel fully likelihood-based approach to handle right-censored survival data with crossing survival curves. This is done by assuming the PE distribution to deal with the baseline hazard in the YP model. We emphasize that using the semiparametric PE approach to extend the original YP model has not been seen in the literature. We also developed an R package called YPPE to fit the proposed model. Some important advantages of the methodology proposed here are: (i) it has the tractability

of parametric models; (ii) it provides a continuous survival function being convenient for the detection of the intersection point of two survival curves; (iii) it has the flexibility of a semiparametric model allowing different shapes for the hazard function; (iv) the routine for maximum likelihood estimation and inference is straightforward and easy-to-implement. Another point to be highlighted is the fact that the original reference for the YP model is focused on the two-sample case with the general regression setting being extended in [Diao, Zeng and Yang \(2013\)](#) as previously mentioned. We explore the YP model with the PE baseline hazard, hereafter called “YPPE model”, using categorical and continuous covariates in this paper.

This work is organized as follows. The proposed model is described in Section 2. A comprehensive Monte Carlo simulation study is conducted in Section 3 to evaluate the performance of the YPPE model. Section 4 shows an empirical illustration where the new model is applied to study the survival times of patients enrolled in a gastric cancer clinical trial. Finally, Section 5 presents the main conclusions, final remarks and discusses future research.

2 Model formulation

Let T be a nonnegative random variable representing the time until the occurrence of an event of interest. In order to accommodate survival data with crossing survival curves, [Yang and Prentice \(2005\)](#) proposed the following model in terms of survival function of T :

$$S(t|\mathbf{z}) = \left[1 + \frac{\lambda}{\theta} R_0(t) \right]^{-\theta}, \quad (1)$$

where $\mathbf{z} = (z_1, \dots, z_q)$ is a row vector of explanatory variables, $\lambda = \exp(\mathbf{z}\boldsymbol{\psi})$ and $\theta = \exp(\mathbf{z}\boldsymbol{\phi})$, $\boldsymbol{\psi} = (\psi_1, \dots, \psi_q)^\top$ and $\boldsymbol{\phi} = (\phi_1, \dots, \phi_q)^\top$ are vectors of regression coefficients without intercepts, and $R_0(t) = F_0(t)/S_0(t)$ corresponds to the baseline odds of death (or failure) by time t . The terms $F_0(t)$ and $S_0(t)$ are the baseline cumulative distribution function and the baseline survival function, respectively. Note that $F_0(t) = 1 - S_0(t)$.

The hazard function associated with (1), can be expressed as

$$h(t|\mathbf{z}) = \frac{\lambda\theta}{\lambda F_0(t) + \theta S_0(t)} h_0(t), \quad (2)$$

where $h_0(t) = -\frac{d}{dt} \log(S_0(t)) = h(t|\mathbf{0})$.

The YP model has some interesting and attractive features. First, it is easy to see from (1) and (2) that the PH and PO models arise as particular cases when $\boldsymbol{\psi} = \boldsymbol{\phi}$ and $\boldsymbol{\phi} = \mathbf{0}$, respectively. Another point is that a scenario with crossing survival curves can be obtained when $\psi_j \phi_j < 0$, for any pair of coefficients (ψ_j, ϕ_j) and $j = 1, \dots, q$. Examples of different survival functions obtained by varying the values of $\boldsymbol{\psi}$ and $\boldsymbol{\phi}$ are displayed in Figure 1.

It follows from (2) that

$$\lim_{t \rightarrow 0} \frac{h(t|\mathbf{z})}{h(t|\mathbf{0})} = \lambda \quad \text{and} \quad \lim_{t \rightarrow \infty} \frac{h(t|\mathbf{z})}{h(t|\mathbf{0})} = \theta.$$

The quantities λ and θ can be interpreted as the short-term and long-term hazard ratios, respectively. In addition, the elements of $\boldsymbol{\psi}$ and $\boldsymbol{\phi}$ can be regarded as the short-term and long-term regression coefficients, respectively.

We now describe the main aspects related to the piecewise exponential distribution. Consider a time grid $\rho = \{a_1, \dots, a_{m-1}\}$ inducing the following set of intervals:

$$I_k = \begin{cases} (a_{k-1}, a_k], & k = 1, \dots, m-1, \\ (a_{m-1}, \infty), & k = m, \end{cases}$$

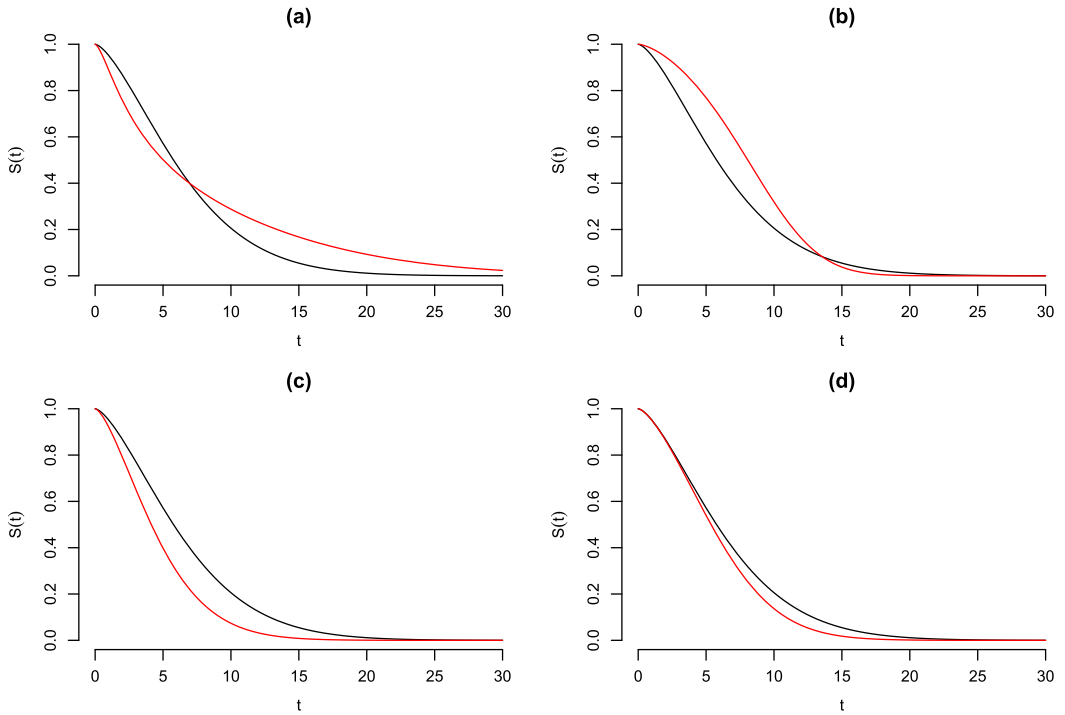


Figure 1 Survival functions according to different choices of ψ and ϕ . Panel (a): $\psi = 1$ and $\phi = -1$ (crossing survivals). Panel (b): $\psi = -1$ and $\phi = 1$ (crossing survivals). Panel (c): $\psi = 0.5$ and $\phi = 0.5$ (PH structure). Panel (d): $\psi = 0$ and $\phi = 0.5$ (PO structure).

with $a_0 = 0$. We shall assume that the baseline hazard function appearing in (2) is constant in each interval induced by ρ , that is

$$h_0(t|\boldsymbol{\xi}, \rho) = \xi_k, \quad \xi_k > 0$$

for $\boldsymbol{\xi} = (\xi_1, \dots, \xi_m)^\top$, $t \in I_k$ and $k = 1, \dots, m$.

The choice of the time grid ρ has a significant impact in terms of goodness-of-fit for the target model. A time grid with a large number of intervals might provide unstable estimates for the failure rates. On the other hand, time grids with few intervals might lead to poor approximations to the true survival function. In practice, the time grid selection must seek a balance in terms of how well the hazard and survival functions can be estimated. Several approaches have been proposed in the literature to address this issue. We shall assume here that the time grid ρ is a known quantity composed by either the set of distinct ordered observed failure times, as suggested by [Breslow \(1974\)](#), or a subset of such observed times, as proposed by [Demarqui et al. \(2012, 2014\)](#). According to [Schneider et al. \(2020\)](#), one advantage of those approaches is that the length of the grid intervals are smaller for time intervals where a large number of time points are observed, and larger for those time intervals with few observations.

Following [Demarqui et al. \(2011\)](#), the baseline survival function $S_0(t|\boldsymbol{\xi}, \rho)$ can be conveniently expressed as:

$$S_0(t|\boldsymbol{\xi}, \rho) = \exp(-H_0(t)),$$

where $H_0(t) = \sum_{k=1}^m \xi_k (t_k - a_{k-1})$ is the baseline cumulative hazard function, and

$$t_k = \begin{cases} a_{k-1}, & \text{if } t < a_{k-1}, \\ t, & \text{if } t \in I_k, \\ a_k, & \text{if } t > a_k, \end{cases}$$

for $k = 1, \dots, m$.

Note that the cumulative hazard function of the PE distribution is a non-decreasing function defined in terms of a sum of positive increments in disjoint intervals. Under the assumption of independent increments (Kalbfleisch, 1978), it can be seen as a realization of a Levy process. This fact can be convenient under the Bayesian framework, since it facilitates the elicitation of prior processes for the cumulative hazard. This is true for the gamma process, which is the most used Levy process to model the cumulative hazard function (Sinha and Dey, 1997). Another attractive characteristic of the PE distribution regards the possibility to add some degree of smoothness on the (baseline) hazard function by the introduction of first-order correlation structures on the failure rates ξ_k 's, such as those proposed by Gamerman (1991) and Arjas and Gasbarra (1994). We emphasize, however, that a Bayesian treatment of the model being proposed here is out of the scope of the paper, and it will be addressed in future works.

Consider a random sample of size n , where all elements are independent, and denote by T_i and C_i the failure and censoring times, respectively. Let \mathbf{z}_i be a row vector of explanatory variables associated with the i -th element in the sample. Assume that the censoring mechanism is non-informative. In addition, the failure times are right-censored so that $Y_i = \min\{T_i, C_i\}$ is the observable failure time. The term $\delta_i = I\{T_i \leq C_i\}$, for $i = 1, \dots, n$, is the failure indicator function. The set of observed data is then denoted by $D = \{(y_i, \delta_i, \mathbf{z}_i); i = 1, \dots, n\}$. Finally, let $\Theta = (\boldsymbol{\psi}^\top, \boldsymbol{\phi}^\top, \boldsymbol{\xi}^\top)$ represent the vector of parameters to be estimated. Since the time grid ρ is regarded as a known quantity in this paper, its notation will be suppressed here for simplicity.

The likelihood function can be expressed as follows:

$$L(\Theta; D) = \prod_{i=1}^n \left[\frac{\lambda_i \theta_i}{\theta_i S_0(y_i | \boldsymbol{\xi}) + \lambda_i F_0(y_i | \boldsymbol{\xi})} h_0(y_i | \boldsymbol{\xi}) \right]^{\delta_i} \left[1 + \frac{\lambda_i}{\theta_i} R_0(y_i | \boldsymbol{\xi}) \right]^{-\theta_i}, \quad (3)$$

where $F_0(y_i | \boldsymbol{\xi}) = 1 - S_0(y_i | \boldsymbol{\xi})$, $\lambda_i = \exp(\mathbf{z}_i \boldsymbol{\psi})$ and $\theta_i = \exp(\mathbf{z}_i \boldsymbol{\phi})$, for $i = 1, \dots, n$. The likelihood function for alternative formulation of the YP model discussed in Yang and Prentice (2015), along with an additional simulation study, is presented in Appendix A.

Maximum likelihood estimates (MLEs) for the parameters are obtained by the direct maximization of the log-likelihood function given in (3) using standard numerical maximization routines. In order to accomplish this task, we developed an R package, called YPPE, which is currently available on CRAN (<https://cran.r-project.org/web/packages/YPPE/index.html>). The YPPE package is based on Stan's Programming Language (Carpenter et al., 2017), and allows the application of three different maximization algorithms (Newton and two related quasi-Newton methods, namely L-BFGS and BFGS). For more details about these optimizers, the reader is referred to Nocedal and Wright (2006). The standard errors of the estimators are computed through the observed Fisher information matrix, which is obtained by inverting the approximated log-likelihood Hessian matrix readily available from the optimizer routines provided by Stan.

In the next section, we empirically investigate some asymptotic properties of the MLEs through a simulation study. This is a comprehensive study based on simulated data replicated in a Monte Carlo (MC) scheme. The main idea is to explore different aspects of the proposed model and compare its results with those from the standard YP model.

3 Simulation study

In this section, we present a Monte Carlo simulation study to evaluate the performance of the model introduced in the previous section. There are two main purposes in this analysis: (i) compare the proposed model with the two-sample semiparametric model in Yang and

Prentice (2005) and (ii) evaluate the performance of the new model in the general regression setting.

In order to generate the simulated data sets, the Weibull baseline survival function $S_0(t|\alpha, \gamma) = \exp(-\gamma t^\alpha)$, with $\alpha = 1.50$ and $\gamma = 0.05$, is assumed to generate the failure times (t_i 's). The censoring times (c_i 's) are obtained from the uniform distribution in the interval $(0, \tau)$, with τ chosen so that the censoring rate is approximately 30% of the observed data. Recall that the final time reported for each sample unit is given by $y_i = \min\{t_i, c_i\}$. We begin the simulation study with the two-sample scenario. The MC schemes are configured with 1000 data sets and they explore three different sample sizes: $n = 50$, $n = 100$ and $n = 200$. In each case, a single binary covariate is included assuming $z_i \sim \text{Bernoulli}(0.5)$, for $i = 1, \dots, n$.

All models were implemented and fitted using the R programming language (R Core Team, 2019). The YPPE model was fitted using our proposed R package YPPE, whereas the YP model in Yang and Prentice (2005) was fitted through the R package YPmodel; see more details in Yang and Prentice (2010, 2011) and Yang and Zhao (2012).

The survival function, defined in the YP model, is a step function with jumps on the observed failure times. In order to ensure a fair comparison between our YPPE model and the original YP model, the endpoints of the intervals forming the grid in the PE distribution are set to be the observed failure times. In other words, each interval contains exactly 1 observation. Naturally, other configurations including more than 1 time point per interval can be applied and this is expected to improve results.

The relative bias reported in Table 1 is calculated according to the following expression:

$$RB(\kappa) = 100(\hat{\kappa} - \kappa_{\text{true}})/|\kappa_{\text{true}}|.$$

In this expression consider that: κ is a generic parameter, $\hat{\kappa}$ is the maximum likelihood estimate and κ_{true} is the true value.

Table 1 Summary for the MC simulation study with 1000 replications and a single binary covariate. Notation: fitted model (Mod), parameter name (Par), true value (True), average point estimate (Est.), average asymptotic standard error (AASE), sample standard deviation of the estimates (SSDE), relative bias (RB), average 95% confidence interval (CI) and coverage probabilities (CP)

Mod	Par	True	Est.	AASE	SSDE	RB(%)	95% CI		CP
							Lower	Upper	
<i>n</i> = 50									
YPPE	ψ	1.0	0.907	0.858	0.870	-9.269	-0.773	2.588	0.940
	ϕ	-1.0	-0.651	3.363	1.788	34.872	-7.242	5.940	0.981
YP	ψ	1.0	1.150	1.419	1.522	14.976	-1.631	3.931	0.863
	ϕ	-1.0	-0.729	1.068	1.147	27.118	-2.822	1.364	0.924
<i>n</i> = 100									
YPPE	ψ	1.0	0.955	0.595	0.606	-4.511	-0.212	2.122	0.947
	ϕ	-1.0	-0.935	0.387	0.384	6.513	-1.694	-0.176	0.969
YP	ψ	1.0	1.138	2.042	1.098	13.787	-2.864	5.140	0.946
	ϕ	-1.0	-0.930	1.701	0.585	6.992	-4.265	2.405	0.993
<i>n</i> = 200									
YPPE	ψ	1.0	0.991	0.417	0.418	-0.944	0.173	1.808	0.949
	ϕ	-1.0	-0.966	0.258	0.265	3.434	-1.471	-0.461	0.956
YP	ψ	1.0	1.039	3.007	0.590	3.939	-4.854	6.933	0.994
	ϕ	-1.0	-0.939	2.412	0.530	6.135	-5.666	3.789	0.995

Table 1 shows the results of the Monte Carlo simulation study. As it can be seen, neither the proposed model nor the YP model performed well when the sample size is small ($n = 50$). In this case, both models show relative biases above $\approx 10\%$ and coverage probabilities far from the nominal level for all parameters. Now looking at the moderate sample sizes ($n = 100$ and $n = 200$) the results in Table 1 change in favour of the proposed YPPE model. It is evident that the YPPE model has a superior performance with respect to the standard semiparametric YP model. Although an improvement in terms of relative bias reduction can be observed for both models under moderate sample sizes, the proposed model provides smaller relative biases and indicates coverage probabilities closer to the nominal level of 95%.

Another important aspect exhibited in Table 1 is the similar results for the AASE and SSDE related to the proposed PE model. This similarity indicates that the method implemented to fit the proposed model is performing properly. Note that this type of result is not true for the standard semiparametric YP model, which seems to overestimate the standard errors of the parameter estimators. In addition, this bad behavior can explain the wider average 95% confidence interval limits and the coverage probabilities above the nominal level observed for this model.

We now turn our attention to the general regression setting. Moderate to large data sets were considered in this analysis to investigate the performance of the YPPE model assuming a regression structure with four covariates. Artificial data sets were simulated taking into account three different sample sizes: $n = 100$, $n = 200$ and $n = 500$. We again use the MC scheme with 1000 replications. The following short-term and long-term linear predictors are explored:

$$\log(\lambda_i) = +2.0z_{1i} - 0.5z_{2i} + 1.5z_{3i} - 1.5z_{4i},$$

$$\log(\theta_i) = -1.0z_{1i} + 1.0z_{2i} - 1.5z_{3i} + 1.5z_{4i}$$

where $z_{1i} \sim \text{Bernoulli}(0.5)$, $z_{2i} \sim N(0, 1)$, $z_{3i} \sim \text{Bernoulli}(0.5)$ and $z_{4i} \sim N(0, 1)$, for $i = 1, \dots, n$ and all of them being independent.

The main reference [Yang and Prentice \(2005\)](#) is entirely focused on the two-sample scenario and does not explore the general regression setting. In fact the paper indicates that the standard YP model can be extended to incorporate covariates, but this was left for future work in the mentioned reference. The corresponding R package `YPmodel` does not allow the analysis using the general configuration. As a consequence of this point, in the next analysis we do not confront the results from the YPPE model and the standard YP case. Recall that, for comparison reasons, the time grid for the YPPE model was initially chosen (analysis of Table 1), with 1 observation per interval. The results presented in Table 2 are obtained by assuming a different grid structure. In this case, the number of intervals is given by $m = \sqrt{n}$. This choice is convenient to reduce the computational burden to fit the model. The endpoints of the intervals are chosen as follows. Given the set ζ of J distinct observed failure times, let k and r be integers such that $J = km + r$. Then, the endpoints of the time grid ρ are chosen among the elements of ζ so that the first $m - r$ intervals have k failure times, and the remaining intervals contain $k + 1$ failures. According to [Demarqui et al. \(2014\)](#), this procedure, which is implemented in the function `timeGrid()` of our package `YPPE`, allows for more failure times to be in the last intervals, where less information is usually available.

As it can be seen from Table 2, relative biases are reasonably low, especially for $n = 200$ and $n = 500$. In addition, the coverage probabilities are, in general, close to the nominal level of 95%. Another important aspect observed here is the fact that both bias and AASE tend to decrease as the sample size increases; this is expected and confirms that the fitting algorithm behaves well. The results displayed in Table 2 also indicate that the standard errors of the parameters are being well estimated, since the AASE and SSDE have similar values for all parameters; this is true regardless of the sample size under investigation. Overall, the

Table 2 Summary for the MC simulation study with 1000 replications and 4 covariates. Notation: parameter name (Par), true value (True), average point estimate (Est.), average asymptotic standard error (AASE), sample standard deviation of the estimates (SSDE), relative bias (RB), average 95% confidence interval (CI) and coverage probabilities (CP)

Par	True	Est.	AASE	SSDE	RB(%)	95% CI		CP
						Lower	Upper	
<i>n</i> = 100 and <i>m</i> = 10								
ψ_1	2.0	1.903	0.614	0.656	-4.855	0.700	3.106	0.931
ψ_2	-0.5	-0.473	0.296	0.315	5.314	-1.053	0.106	0.940
ψ_3	1.5	1.415	0.615	0.655	-5.657	0.209	2.621	0.919
ψ_4	-1.5	-1.477	0.402	0.416	1.562	-2.264	-0.689	0.942
ϕ_1	-1.0	-0.913	0.347	0.397	8.664	-1.593	-0.234	0.927
ϕ_2	1.0	1.047	0.238	0.258	4.738	0.581	1.514	0.951
ϕ_3	-1.5	-1.459	0.350	0.372	2.718	-2.145	-0.773	0.952
ϕ_4	1.5	1.547	0.266	0.295	3.156	1.026	2.068	0.964
<i>n</i> = 200 and <i>m</i> = 15								
ψ_1	2.0	1.946	0.422	0.412	-2.704	1.119	2.773	0.945
ψ_2	-0.5	-0.479	0.200	0.209	4.181	-0.871	-0.087	0.937
ψ_3	1.5	1.430	0.423	0.438	-4.692	0.600	2.259	0.933
ψ_4	-1.5	-1.475	0.276	0.284	1.689	-2.015	-0.935	0.937
ϕ_1	-1.0	-0.953	0.233	0.237	4.657	-1.411	-0.496	0.955
ϕ_2	1.0	1.034	0.160	0.162	3.367	0.720	1.348	0.945
ϕ_3	-1.5	-1.477	0.237	0.251	1.535	-1.941	-1.012	0.939
ϕ_4	1.5	1.528	0.179	0.185	1.848	1.178	1.878	0.937
<i>n</i> = 500 and <i>m</i> = 23								
ψ_1	2.0	1.976	0.264	0.270	-1.223	1.458	2.493	0.936
ψ_2	-0.5	-0.485	0.123	0.124	2.905	-0.727	-0.244	0.948
ψ_3	1.5	1.463	0.264	0.261	-2.441	0.946	1.981	0.949
ψ_4	-1.5	-1.472	0.171	0.174	1.866	-1.807	-1.137	0.936
ϕ_1	-1.0	-0.978	0.144	0.143	2.177	-1.261	-0.695	0.958
ϕ_2	1.0	1.010	0.098	0.097	1.035	0.818	1.202	0.953
ϕ_3	-1.5	-1.482	0.146	0.142	1.218	-1.769	-1.195	0.955
ϕ_4	1.5	1.518	0.110	0.115	1.227	1.303	1.734	0.931

proposed model seems to perform well in the general regression setting for moderate to large data sets.

4 Real data analysis

This section is dedicated to the analysis of a real data set freely available through the R package `Ypmodel` under the label of `gastric`; see also [Gastrointestinal Tumor Study Group \(1982\)](#) as a reference for more details. This gastric cancer data set has become a common application in the literature related to survival analysis and, more specifically, it can be easily found in studies dealing with crossing survival curves; some few references are: [Yang and Prentice \(2005\)](#), [Lee \(2011\)](#), [Yang and Zhao \(2012\)](#), [Diao, Zeng and Yang \(2013\)](#) and [Yang \(2018\)](#). The experiment in this clinical trial involves 90 individuals diagnosed with locally unresectable (advanced) gastric cancer. The participants were randomly assigned to the following groups: (i) the control group composed by 45 patients receiving chemotherapy and (ii) the treatment group including 45 patients receiving a combination of chemotherapy and radiation therapy. These individuals were followed within this study for about 5 years. Three

Table 3 Summary of the models fitted to the gastric cancer data

Model	Par	Est.	SE	95% CI		z	p-value
				Lower	Upper		
YPPE	ψ	1.837	0.648	0.567	3.108	2.834	0.005
	ϕ	-1.017	0.300	-1.606	-0.429	-3.387	0.001
YP	ψ	1.600	0.538	0.547	2.656	2.977	0.003
	ϕ	-0.906	0.248	-1.393	-0.421	-3.650	<0.001

variables are reported in the data set for each patient: the time response representing either a failure (time to death) or a right censoring, a binary failure indicator identifying those patients experiencing the event of interest and, finally, a group binary indicator with 1 meaning the treatment category. Note that this application contains a single binary covariate; therefore, it can be explored and compared via the YPPE and YP models.

Table 3 summarizes the results obtained for both models. As it can be observed, the short-term (ψ) and long-term (ϕ) regression coefficients, within each model, are estimated with opposite signs and they have distinct magnitudes. This can be observed by either looking at the point estimate (column Est.) or the 95% confidence intervals. This behavior is a clear indication of survival curves having an intersection at some intermediate time point between 0 and the maximum. In other words, the top and bottom positioning of the curves are inverted for the intervals below and above the crossing time point; see Figure 2 for a visual idea. This inversion suggests the existence of an alteration in the effectiveness of the treatment at some point during the follow-up period of the study. In general, the results tend to be similar when comparing the corresponding estimates from both models. Note that the standard error related to ψ is larger than the one for ϕ . In addition, all p-values from the z-test are small, indicating significant regression coefficients.

One interesting and attractive feature of the proposed YPPE model is the fact that it provides a continuous survival function. This aspect allows us to apply standard procedures to find the roots of nonlinear equations to determine accurately the time point at which the survival curves intersect each other. One possibility to handle this problem in R is to use the command `uniroot` (Brent, 1973) for unidimensional searches. In line with this idea, right after fitting our YPPE model to the gastric cancer data, we apply the `uniroot` function to conclude that the crossing time occurs, for this application, at the time point given by the day 863 within the full period of the study. From the practical point of view, this means that before the day 863, the patients in the control group (only chemotherapy) have better survival rates than those in the treatment group. On the other hand, the benefits of combining chemotherapy with radiotherapy tend to emerge in a later stage of the study (after day 863).

The 95% confidence interval for the crossing survival time, obtained via nonparametric Bootstrap (based on 1000 resamples), is [1, 1967]. As it can be noted, the confidence interval in this case is quite wide, reflecting the uncertainty associated with the small sample size of the gastric cancer data. Larger clinical trials are expected to provide shorter intervals establishing with higher precision the region of the crossing point.

The panel (a) of Figure 2 shows the Kaplan–Meier curves associated with the two treatments and the survival curves estimated via the semiparametric YPPE and YP models. The panel (b) of Figure 2 displays the estimated survival curves, provided by the YPPE model, along with the estimated time at which the survival curves are expected to cross. Finally, the panel (c) of Figure 2 provides the hazard ratio estimated by the proposed model. As it can be seen in panel (a), both models seem to accommodate and represent well the data, since their estimated survival curves tend to agree with the Kaplan–Meyer survival estimates. Moreover,

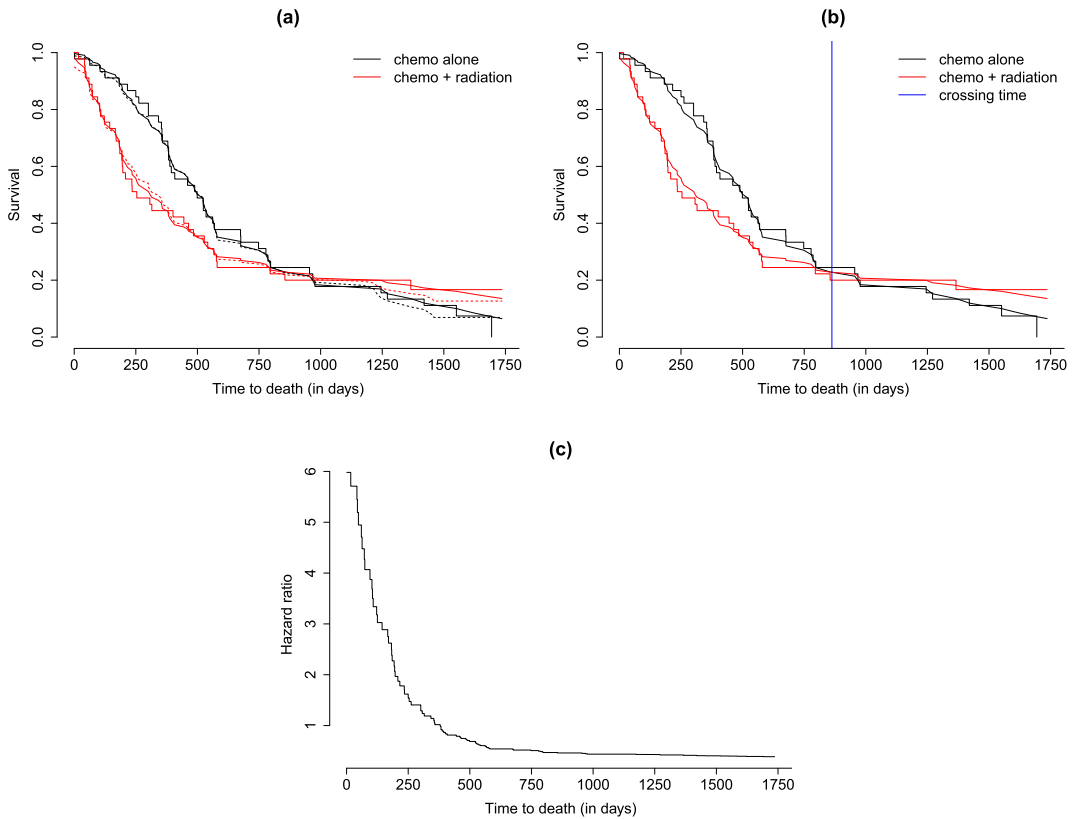


Figure 2 Analysis of the gastric cancer data set. Panel (a): Kaplan–Meier estimates for the survival curves along with the YPPE estimates (solid lines) and the YP estimates (dashed lines). Panel (b): estimated survival curves provided by the YPPE model along with the estimated time point at which the survival curves probably cross each other. Panel (c): estimated hazard ratio provided by the YPPE model.

Table 4 Short and long term hazard ratios summarized by point estimate (Est.) along with the 95% confidence interval (CI)

Model	Hazard ratio	Est.	95% CI	
			Lower	Upper
YPPE	chemo + radiation (short term)	6.280	1.762	22.375
	chemo alone (long term)	2.766	1.535	4.984
YP	chemo + radiation (short term)	4.953	1.728	14.239
	chemo alone (long term)	2.474	1.523	4.027

panel (c) shows a monotonically decreasing behavior of the estimated hazard ratio over time, which is in agreement with the estimated regression coefficients provided by the proposed model.

The short and long term hazard ratios provided by the fitted models are summarized in Table 4. The risk of death for patients under chemo + radiation treatment at the beginning of the therapy is significantly superior (6.28 times greater) than that observed for patients treated with chemotherapy alone. On the long run, however, the combination of chemotherapy and radiation has a beneficial effect, and the risk of death for those patients treated with chemotherapy alone is significantly superior (2.77 times greater) than that observed for the patients receiving the combined treatment. It is also possible to note that the corresponding

hazard ratios estimated by the YP model are slightly smaller than those obtained by the YPPE model.

5 Conclusions

This paper presents a fully likelihood-based approach to deal with crossing survival curves as an extension to the standard YP model proposed in 2005 for a two-sample case. The main difference with respect to other extensions of the YP model is the fact that we take advantage of the piecewise exponential semiparametric modeling to allow a flexible representation of the baseline hazard function. This also configures the main contribution of the paper, since no other study combining these two aspects (YP model structure and PE distribution) can be found in the literature of survival analysis. Using the PE distribution brings some advantages when comparing to other semiparametric options for the YP model. The YPPE model preserves the flexibility of the semiparametric models and the tractability of the parametric ones. In addition, it is relatively easy-to-implement using standard maximization routines. Estimation of parameters, hazard function, survival function and hazard ratios is straightforward. Another important aspect to be emphasized is the fact that the survival function has a continuous representation via the YPPE model; this is not true in the original YP model and other approaches presented in the literature, where a step function is obtained as the survival representation. As a result of this feature, the time in which the survival curves (treatment and control groups) intersect each other can be easily and accurately determined.

A comprehensive MC simulation study was developed to examine the performance of the proposed PE model in comparison with the YP model. The results indicate that the YPPE model provides better results with smaller relative biases being observed for most parameters. Using simulated data sets, the behavior of the YPPE model was also investigated for a general regression setting involving several covariates. The standard YP model can be extended to this context, but the original paper in 2005 does not explore this type of result. Our findings suggest that the YPPE model also has a good performance when dealing with several covariates.

After the simulation study, the present paper presents a real application concerning a well known data set related to a clinical trial for patients detected with advanced gastric cancer. In summary, the results of the YPPE and YP models are similar and they clearly indicate significant regression coefficients with opposite signs, which is expected for the scenario where the survival curves have an intersection.

As supplement to this paper, we developed the R package called YPPE to fit the proposed model in this study; the source code is available at <https://cran.r-project.org/web/packages/YPPE/index.html>. Some R codes used throughout the paper are presented in Appendix B. In terms of future work, we indicate that the approach presented here can also be extended to accommodate survival data with cure fraction and interval-censored observations. We also intend to develop a Bayesian version for our model. Another possible extension is to introduce some degree of smoothness on the PE failure rates in adjacent intervals by using first-order autocorrelated processes.

Appendix A: Additional simulation study

In this section, we present a replication of the simulation study carried out in Yang and Prentice (2015).

Table 5 Bias of the estimators provided by the proposed model based on 1000 MC simulations

		0.9 ↑ 1.2 Censoring (0.9 ↑ 1.2)			1.2 ↑ 0.8 Censoring (1.2 ↑ 0.8)		
		10%	30%	50%	10%	30%	50%
n = 100	ψ	-0.0012	0.0139	0.0098	-0.0015	-0.0250	-0.0021
	ϕ	0.0206	-0.0387	-0.0271	-0.0555	-0.0156	-0.0107
	β	0.0076	0.0169	0.0203	0.0133	0.0185	0.0082
n = 400	ψ	-0.0026	0.0140	0.0007	0.0065	-0.0071	0.0005
	ϕ	-0.0073	-0.0282	0.0009	-0.0268	-0.0328	-0.0360
	β	0.0006	0.0047	0.0045	0.0023	0.0040	0.0009
n = 800	ψ	0.0037	0.0062	-0.0031	-0.0055	-0.0031	-0.0051
	ϕ	-0.0086	-0.0118	-0.0041	-0.0238	-0.0212	-0.0186
	β	-0.0004	0.0003	0.0013	0.0031	0.0039	0.0014

Under the alternative formulation of the YP model proposed by Yang and Prentice (2015), the likelihood function assumes the form:

$$L(\Theta^*; D^*) = \prod_{i=1}^n \left[\frac{\lambda_i \theta_i \kappa_i}{\theta_i S_0(y_i | \xi) + \lambda_i F_0(y_i | \xi)} h_0(y_i | \xi) \right]^{\delta_i} \left[1 + \frac{\lambda_i}{\theta_i} R_0(y_i | \xi) \right]^{-\theta_i \kappa_i},$$

where $D^* = \{(y_i, \delta_i, \mathbf{z}_i, \mathbf{x}_i); i = 1, \dots, n\}$, $\Theta^* = (\psi^\top, \phi^\top, \beta^\top, \xi^\top)$, $h_0(y_i | \xi)$, $F_0(y_i | \xi)$, $S_0(y_i | \xi)$, λ_i and θ_i are defined as in Section 2, $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^\top$ is a row vector containing the subset of covariates with constant effects over time, $\beta = (\beta_1, \dots, \beta_p)^\top$ vector of corresponding regression coefficients, and $\kappa_i = \exp(\mathbf{x}_i \beta)$, for $i = 1, \dots, n$. Inferences on Θ^* are carried out straightforwardly as described in Section 2.

In Table 5 we show the biases of the estimators for the regression coefficients provided by the proposed model, under the alternative formulation of the YP, in a replication of the Monte Carlo simulation study presented in Table 1 of Yang and Prentice (2015). For this particular simulation study we chose to standardize all covariates entering in the linear predictors as a strategy to avoid numerical problems. The corresponding inverse transformation to recover the parameters in their original scale was applied in order to reach the magnitude of the results shown in Yang and Prentice (2015).

By comparing the biases shown in Table 5 with those reported in Table 1 of Yang and Prentice (2015), one may see that the proposed model seems to perform well compared to the model investigated in Yang and Prentice (2015), showing, in general, smaller biases for the regression coefficients estimators.

Appendix B: R code related to the YPPE package

In this section, we provide examples of the R code used throughout the paper to fit the proposed model.

The R code used to fit the gastric cancer data in Section 4 is presented below:

```
> library(YPPE)
> # loading the gastric cancer data:
> data(gastric)
> # fitting the model:
> fit <- yppe(Surv(time, status)~trt, data=gastric, init=0)
> summary(fit)
```

```

Call:
yppe(formula = Surv(time, status) ~ trt, data = gastric, init = 0)
Short-term coefficients:
      Estimate StdErr z.value p.value
trt    1.8373  0.6483   2.834 0.004597 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Long-term coefficients:
      Estimate  StdErr z.value  p.value
trt -1.01753   0.30036 -3.3877 0.0007049 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
---
loglik = 47.65404    AIC = 62.69192

> # Estimating the crossing survival time:
> newdata1=data.frame(trt=0)
> newdata2=data.frame(trt=1)
> crossTime(fit, newdata1, newdata2, nboot=1000)
      Est. 2.5%    97.5%
1 862.363    1 1967.201

> # plotting the estimated survival curves:
> newdata=data.frame(trt=as.factor(0:1))
> St <- survfit(fit, newdata)
> ekm <- survfit(Surv(time, status)~trt, data=gastric)
> time <- sort(gastric$time)
> plot(ekm, col=1:2)
> lines(time, St[[1]])
> lines(time, St[[2]], col="red")

```

In order to fit the proposed model with the alternative formulation of the YP model, the following R code can be used:

```

> library(YPPE)
> simdata <- read.table("simdata.txt", header=TRUE)
> fit <- yppe(Surv(time, status)~x1|x2, data=simdata, init=0)
> summary(fit)
Call:
yppe(formula = Surv(time, status) ~ x1 | x2, data = simdata,
      init = 0)
Short-term coefficients:
      Estimate StdErr z.value p.value
x1  0.10615  0.47385   0.224  0.8228
Long-term coefficients:
      Estimate  StdErr z.value p.value
x1 -0.53777   0.52239 -1.0294  0.3033
Proportional hazards coefficients:
      Estimate StdErr z.value  p.value
x2  0.38583  0.11061  3.4882 0.0004864 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
---
loglik = 74.67624    AIC = 10.64752

```

Acknowledgments

The authors would like to express their gratitude to the Associate Editor and the three anonymous referees for the careful reading of this paper. Their comments and suggestions contributed substantially to its improvement. The second author gratefully acknowledge the support from Fundação de Amparo a Pesquisa do Estado de Minas Gerais (FAPEMIG).

References

- Arjas, E. and Gasbarra, D. (1994). Nonparametric Bayesian inference from right censored survival data using the Gibbs sampler. *Statistica Sinica* **4**, 505–524. [MR1309427](#)
- Brent, R. P. (1973). *Algorithms for Minimization Without Derivatives*. Englewood Cliffs: Prentice-Hall. [MR0339493](#)
- Breslow, N. (1972). Discussion on regression models and life-tables (by D. R. Cox). *Journal of the Royal Statistical Society, Series B* **34**, 216–217. [MR0341758](#)
- Breslow, N. (1974). Covariance analysis of censored survival data. *Biometrics* **30**, 89–99.
- Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M., Guo, J., Li, P. and Riddell, A. (2017). Stan: A probabilistic programming language. *Journal of Statistical Software* **76**.
- Clark, D. E. and Ryan, L. M. (2002). Concurrent prediction of hospital mortality and length of stay from risk factors on admission. *Health Services Research* **37**, 631–645.
- Cox, D. R. (1972). Regression models and life-tables. *Journal of the Royal Statistical Society, Series B* **34**, 187–220. [MR0341758](#)
- Demarqui, F. N., Dey, D. K., Loschi, R. H. and Colosimo, E. A. (2011). Modeling survival data using the piecewise exponential model with random time grid. In *Recent Advances in Biostatistics, Vol. 4* (M. Bhattacharjee, S. K. Dhar and S. Subramanian, eds.) 109–122. Singapore: World Scientific. https://doi.org/10.1142/9789814329804_0006
- Demarqui, F. N., Dey, D. K., Loschi, R. H. and Colosimo, E. A. (2014). Fully semiparametric Bayesian approach for modeling survival data with cure fraction. *Biometrical Journal* **56**, 198–218.
- Demarqui, F. N., Loschi, R. H., Dey, D. K. and Colosimo, E. A. (2012). A class of dynamic piecewise exponential models with random time grid. *Journal of Statistical Planning and Inference* **142**, 728–742. [MR2853579](#) <https://doi.org/10.1016/j.jspi.2011.09.006>
- Diao, G., Zeng, D. and Yang, S. (2013). Efficient semiparametric estimation of short-term and long-term hazard ratios with right-censored data. *Biometrics* **69**, 840–849. [MR3146780](#) <https://doi.org/10.1111/biom.12097>
- EGGE, K. and ZAHN, P. H. (1999). Survival of glaucoma patients. *Acta Ophthalmologica Scandinavica* **77**, 397–401.
- Gamerman, D. (1991). Dynamic Bayesian models for survival data. *Journal of the Royal Statistical Society Series C Applied Statistics* **40**, 63–79.
- Gastrointestinal Tumor Study Group (1982). A comparison of combination chemotherapy and combined modality therapy for locally advanced gastric carcinoma. *Cancer* **49**, 1771–1777.
- Ibrahim, J. G., Chen, M. H. and Sinha, D. (2001). *Bayesian Survival Analysis*. New York: Springer. [MR1876598](#) <https://doi.org/10.1007/978-1-4757-3447-8>
- Kalbfleisch, J. D. (1978). Non-parametric Bayesian analysis of survival time data. *Journal of the Royal Statistical Society, Series B, Methodological* **40**, 214–221. [MR0517442](#)
- Kalbfleisch, J. D. and Prentice, R. L. (1973). Marginal likelihoods based on Cox's regression and life model. *Biometrika* **60**, 267–278. [MR0326939](#) <https://doi.org/10.1093/biomet/60.2.267>
- Lee, S. H. (2011). Maximum of the weighted Kaplan–Meier tests for the two-sample censored data. *Journal of Statistical Computation and Simulation* **81**, 1017–1026. [MR2820063](#) <https://doi.org/10.1080/00949651003627753>
- Nieto-Barajas, L. E. (2014). Bayesian semiparametric analysis of short- and long-term hazard ratios with covariates. *Computational Statistics & Data Analysis* **71**, 477–490. [MR3131984](#) <https://doi.org/10.1016/j.csda.2013.03.012>
- Nocedal, J. and Wright, S. J. (2006). *Numerical Optimization*, 2nd ed. New York: Springer. [MR2244940](#)
- Putter, H., Sasako, M., Hartgrink, H. H., van-de-Velde, C. J. H. and van-Houwelingen, J. C. (2005). Long-term survival with non-proportional hazards: Results from the Dutch gastric cancer trial. *Statistics in Medicine* **24**, 2807–2821. [MR2201852](#) <https://doi.org/10.1002/sim.2143>
- R Core Team (2019). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.

- Sahu, S. K., Dey, D. K., Aslanidou, H. and Sinha, D. (1997). A Weibull regression model with gamma frailties for multivariate survival data. *Lifetime Data Analysis* **3**, 123–137.
- Schneider, S., Demarqui, F. N., Colosimo, E. A. and Mayrink, V. D. (2020). An approach to model clustered survival data with dependent censoring. *Biometrical Journal* **62**, 157–174. MR2635103 <https://doi.org/10.1002/bimj.201800391>
- Shyur, H. J., Elsayed, E. A. and Luxhoj, J. T. (1999). A general model for accelerated life testing with time-dependent covariates. *Naval Research Logistics* **49**, 303–321. MR1677541 [https://doi.org/10.1002/\(SICI\)1520-6750\(199904\)46:3<AID-NAV4>3.3.CO;2-W](https://doi.org/10.1002/(SICI)1520-6750(199904)46:3<AID-NAV4>3.3.CO;2-W)
- Sinha, D., Chen, M. H. and Ghosh, S. K. (1999). Bayesian analysis and model selection for interval-censored survival data. *Biometrics* **55**, 585–590. MR1705161 <https://doi.org/10.1111/j.0006-341X.1999.00585.x>
- Sinha, D. and Dey, D. K. (1997). Semiparametric Bayesian analysis of survival data. *Journal of the American Statistical Association* **92**, 1195–1212.
- Tong, X., Zhu, C. and Sun, J. (2007). Semiparametric regression analysis of two-sample current status data, with applications to tumorigenicity experiments. *Canadian Journal of Statistics* **35**, 575–584. MR2416857 <https://doi.org/10.1002/cjs.5550350408>
- Yang, S. (2018). Improving testing and description of treatment effect in clinical trials with survival outcomes. *Statistics in Medicine* **38**, 530–544. MR3902596 <https://doi.org/10.1002/sim.7676>
- Yang, S. and Prentice, R. L. (2005). Semiparametric analysis of short-term and long-term hazard ratios with two-sample survival data. *Biometrika* **92**, 1–17. MR2158606 <https://doi.org/10.1093/biomet/92.1.1>
- Yang, S. and Prentice, R. L. (2010). Improved logrank-type tests for survival data using adaptive weights. *Biometrics* **66**, 30–38. MR2756688 <https://doi.org/10.1111/j.1541-0420.2009.01243.x>
- Yang, S. and Prentice, R. L. (2011). Estimation of the 2-sample hazard ratio function using a semiparametric model. *Biostatistics* **12**, 354–368.
- Yang, S. and Prentice, R. L. (2015). Assessing potentially time-dependent treatment effect from clinical trials and observational studies for survival data, with applications to the Women’s Health Initiative combined hormone therapy trial. *Statistics in Medicine* **34**, 1801–1817. MR3334693 <https://doi.org/10.1002/sim.6453>
- Yang, S. and Zhao, Y. (2012). Checking the short-term and long-term hazard ratio model for survival data. *Scandinavian Journal of Statistics* **39**, 554–567. MR2971638 <https://doi.org/10.1111/j.1467-9469.2012.00804.x>
- Zeng, D. and Lin, D. Y. (2007). Maximum likelihood estimation in semiparametric regression models with censored data. *Journal of the Royal Statistical Society, Series B* **69**, 507–564. MR2370068 <https://doi.org/10.1111/j.1369-7412.2007.00606.x>
- Zhang, H., Wang, P. and Sun, J. (2018). Regression analysis of interval-censored failure time data with possibly crossing hazards. *Statistics in Medicine* **37**, 768–775. MR3760447 <https://doi.org/10.1002/sim.7538>

Av. Antônio Carlos, 6627
Departamento de Estatística, ICEx
UFMG
Belo Horizonte, MG, 31270-901
Brazil
E-mail: fndemarqui@ufmg.br
vdm@est.ufmg.br