



City Research Online

City, University of London Institutional Repository

Citation: Bishopberger, S., Hiabu, M., Mammen, E. & Perch Nielsen, J. (2025). Smooth Backfitting for Additive Hazard Rates. *Scandinavian Journal of Statistics*, sjos.70004. doi: 10.1111/sjos.70004

This is the accepted version of the paper.

This version of the publication may differ from the final published version.

Permanent repository link: <https://openaccess.city.ac.uk/id/eprint/35380/>

Link to published version: <https://doi.org/10.1111/sjos.70004>

Copyright: City Research Online aims to make research outputs of City, University of London available to a wider audience. Copyright and Moral Rights remain with the author(s) and/or copyright holders. URLs from City Research Online may be freely distributed and linked to.

Reuse: Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

Smooth Backfitting for Additive Hazard Rates

Stephan M. Bischofberger^{*1}, Munir Hiabu³, Enno Mammen⁴, and Jens Perch Nielsen²

¹Staburo GmbH, Munich, Germany

²Bayes Business School, City, University of London, United Kingdom

³Department of Mathematical Sciences, University of Copenhagen, Denmark

⁴Institute for Applied Mathematics, Heidelberg University, Germany

March 7, 2025

Smooth backfitting was first introduced in an additive regression setting via a direct projection alternative to the classic backfitting method by Buja, Hastie and Tibshirani. This paper translates the original smooth backfitting concept to a survival model considering an additively structured hazard. The model allows for censoring and truncation patterns occurring in many applications such as medical studies or actuarial reserving. Our estimators are shown to be a projection of the data into the space of multivariate hazard functions with smooth additive components. Hence, our hazard estimator is the closest nonparametric additive fit even if the actual hazard rate is not additive. This is different to other additive structure estimators where it is not clear what is being estimated if the model is not true. We provide full asymptotic theory for our estimators. We propose an implementation of estimators that show good performance in practice.

Keywords: additive hazard model; local linear kernel estimation; smooth backfitting; survival analysis.

1. Introduction

This paper introduces a fundamental model and estimator for structured multivariate marker dependent hazards: the smooth backfitting of additive hazards. In structured non-parametric regression, Mammen et al. [1999b] modelled and estimated the additive structure by projecting data onto the appropriate additive subspace. The resulting projection estimator is known as the smooth backfitting estimator. The name comes from the fact that when calculating the projection estimator iteratively, then one must not only smooth the component that is being updated, but all components. This is different to classical backfitting [Buja et al., 1989a] where only the component that is being updated is smoothed. It has been shown that smooth backfitting performs much better than previous comparable smoothing kernel based backfitting approaches, in particular in high dimensional problems and with correlated covariates, see Nielsen and Sperlich [2005]. A theoretical comparison between classical and smooth backfitting for

^{*}Corresponding author: Stephan M. Bischofberger, e-mail: bischofberger@staburo.de, address: Aschauer Str. 26a, 81549 Munich, Germany

additive regression models was recently done in [Huang and Yu, 2019], explaining why smoothing of all components has a better adaption. Since the initial smooth backfitting paper many variations and extensions have been developed using smooth backfitting to tackle more sophisticated problems in mathematical statistics, Mammen and Nielsen [2003], Yu et al. [2008], Mammen and Yu [2009], Mammen et al. [2014], Han et al. [2018], Mammen and Sperlich [2021], Bissantz et al. [2016], Han et al. [2020], Jeon et al. [2020], Hiabu et al. [2021a] and Gregory et al. [2020].

The aim of the current paper is to transfer the original approach of additive non-parametric structures to marker dependent hazard estimation and to allow for a potentially high number of covariates with possibly correlated markers. It turns out that when the original estimation problem is phrased as a minimisation problem in the correct way via a counting process formulation, then our smooth backfitting additive hazard approach can be implemented and analysed in a very similar way to smooth backfitting in regression. We see this as a necessary step to understand more complicated structures in marker dependent hazards. The additive subspace is closed making analysis more accessible and the additive structure allows for a more immediate interpretation than more complicated models of structured hazards. One important alternative structure is the multiplicative or proportional hazard model. Survival analysis practitioners often work with such multiplicative marker dependent hazard models, including the Cox model. Smooth backfitting for the multiplicative model was recently analysed in Hiabu et al. [2021a], where the analysis was challenged by the shape of the multiplicative subspace that is not closed like the additive subspace is and where some tricks had to be developed, e.g. a solution weighted optimization, to arrive at a tractable estimation method and analysis. The additive approach developed in this paper does not face these two latter challenges and it had perhaps been more natural to have developed this current paper first and then Hiabu et al. [2021a] afterwards. Both this current paper and Hiabu et al. [2021a] arrive at the same conclusion for smooth backfitting of marker dependent hazard estimators as the authors in Nielsen and Sperlich [2005] did for smooth backfitting of non-parametric regression: Smoothing all components in every iteration step and not only smoothing the component that is being updated is important. Otherwise the estimator breaks down in many cases - in particular in high dimensions - where smooth backfitting still works. Smooth backfitting seems more reliable than classical backfitting of kernel estimators and we expect that the additive marker dependent hazard model and estimator of this paper can be an important starting point for further developments of structured marker dependent hazard approaches in survival analysis, just like the many developments we have seen in non-parametric regression. In the next section we give some insight on the additive model itself and its role in marker dependent hazard models as a practical survival analysis tool.

2. Additive structured hazards and related literature

One well known model in hazard regression is the proportional hazards model of Cox [1972] and it has been seen as the natural equivalent to additive regression functions in linear and nonparametric regression. As pointed out in [Martinussen and Scheike, 2006, p. 103], additive hazard models have been “somewhat overlooked in practice” although they share the same advantages of additive regression models concerning both theoretical properties and implementation. To the best knowledge of the authors, this is still the case, with some exceptions [Tchetgen Tchetgen et al., 2015, Aalen et al., 2019, Dukes et al., 2019]. However, in certain applications an additive relationship in the hazard function is indeed more plausible than a proportional one [Beslow and Day, 1987, Lin and Ying, 1994, Kravdal, 1997, McDaniel et al., 2019]. Moreover, [Aalen et al., 2008, pp.155f] provides a variety of reasons for additive risk factors.

In the original additive hazards model [Aalen, 1980], the intensity of a counting process $\{N(t) : t \in [0, 1]\}$, conditional on the d -dimensional covariate $Z(t) = (Z_1(t), \dots, Z_d(t))^T$, satisfies

$$\lambda(t) = Z^T(t)\beta(t)Y(t) \quad (1)$$

at time t with a regression coefficient $\beta(t) = (\beta_1(t), \dots, \beta_d(t))^T$ and exposure Y which is equal to

unity when an individual is at risk. An overview about this model is given in Martinussen and Scheike [2006] in which the authors praise it as a simple nondistributional model that is easy to implement. Nonparametric estimators of the cumulative regression coefficient $B(t) = \int_0^t \beta(s)ds$ in model (1) have been examined in McKeague [1988] and Huffer and McKeague [1991] among others.

Model (1) imposes a linear relationship between the intensity and the value of the covariates through $Z^T(t)\beta(t)$. We loosen up the assumption of linearity. Before introducing the model we investigate in this article, we describe the most general model and its disadvantages and explain why we assume certain additive constraints. The completely nonparametric conditional intensity model

$$\lambda(t) = \alpha(t|Z)Y(t) \quad (2)$$

for a conditional hazard function α generalizes model (1) making it the most flexible model. As it is common, we assume $\alpha(t|Z) = \alpha(t, Z(t))$ in this paper, i.e. that the conditional hazard at time t given the covariates only depends on the values of the covariates at time t and not on the values of the past.

Model (2) has first been introduced for time-constant covariates in Beran [1981]. Time dependent covariates were considered in McKeague and Utikal [1990] and Nielsen and Linton [1995]. Other examples from the vast literature on nonparametric hazard estimators for this model include Van Keilegom and Veraverbeke [2001] or Spierdijk [2008]. Without further structural restrictions, estimators of (2) suffer from the curse of dimensionality: The rate of convergence decreases exponentially. This is a well known issue for unstructured nonparametric estimators, making them in many cases in-practical already in dimensions higher than, say, three. That one can not do better in the unstructured nonparametric case is known at least since Stone [1980] who provided formulas for the best possible rate of convergence for nonparametric estimators. Accordingly, the aforementioned nonparametric hazard estimators were only illustrated for the case with one-dimensional covariate Z .

To overcome this issue, one has to focus on a model that is more restrictive than the unstructured nonparametric hazard model (2). We restrict our assumptions on an additive model which is nested in (2). However, instead of the original additive Aalen model (1), we assume that the hazard rate consists of additive nonparametric components,

$$\alpha(t, z) = \alpha^* + \alpha_0(t) + \alpha_1(z_1) + \dots + \alpha_d(z_d), \quad (3)$$

with smooth, but not further restricted, components α_k , $k = 1, \dots, d$, depending on covariate values z_1, \dots, z_d . The constant α^* is a norming constant making the decomposition unique, as will later be further specified.

The additive model (3) is both more general but also more restrictive than the additive Aalen model (1). It is more restrictive because it does not allow the effect of covariates on the hazard to change with time. It is more general because the effect of the covariates on the hazard do not need to be linear. A very interesting model that generalises both models is to replace each component $\alpha_k(z_k)$, $k \geq 1$, in (3) by a two-dimensional components $\alpha_k(t, z_k)$ capable of capturing a covariate effect that changes with time. While we do not consider this more general setting in this paper, we see the work done in this paper as a crucial step towards developing methods of such a more general kind. Another possible generalisation is to consider multiple time scales, see e.g. Hiabu et al. [2021b].

To estimate the components in (3), we propose a local polynomial least squares minimisation under the constraint (3). The solution can be identified with the projection of the observation into the space of local polynomial additive hazard functions and can be calculated through a simple iterative procedure. We call the resulting estimator additive smooth backfitting hazard estimator.

When estimating the hazard function $\alpha(t, z)$, by the nature of equation (3), it can happen that the estimate is negative at certain points. This is expected to happen especially more if the underlying hazard function is far from being additive. However, it is reassuring that the smooth backfitting components $\hat{\alpha}_k$ will still have a clear interpretation as approximation of the closest additive fit. In practice, if probabilities need to be calculated, one ad-hoc solution is to use the non-additive adjusted hazard

$$\alpha^{adj} = \max(\alpha, \varepsilon), \varepsilon \geq 0.$$

Indeed, this is also what we do in the application Section 6.1.1 for $\varepsilon = 0$ with satisfying results.

3. The additive hazard model

Let $\mathcal{T} > 0$. We observe n *i.i.d.* copies of the stochastic processes $\{(N(t), Y(t), Z(t)) : t \in [0, \mathcal{T}]\}$ where N is a right-continuous counting process which is zero at time zero and which has jumps of size one. We assume that Y is a left-continuous stochastic process with values in $\{0, 1\}$ and which equals unity if the observed individual is at risk. Moreover, let Z be a d -dimensional left-continuous stochastic process with $Z(t) \in [0, R]^d$, $t \in [0, \mathcal{T}]$, for some $R > 0$. The multivariate process $((N_1, Y_1, Z_1), \dots, (N_n, Y_n, Z_n))$ is assumed to be adapted to the filtration $\{\mathcal{F}_t : t \in [0, \mathcal{T}]\}$ which satisfies the *usual conditions* [Andersen et al., 1993, p. 60].

In the following, we assume that for each $i = 1, \dots, n$, the process N_i satisfies Aalen's multiplicative intensity model, i.e. that its intensity λ_i satisfies

$$\lambda_i(t) = \lim_{h \downarrow 0} h^{-1} \mathbb{E}[N_i((t+h)-) - N_i(t-)| \mathcal{F}_{t-}] = \alpha(t, Z_i(t))Y_i(t), \quad (4)$$

where $Y_i(t)$ is indicating if individual i is at risk at time t . The function $\alpha(t, Z(t))$ is the conditional hazard rate given the covariates Z at time t . Furthermore, we assume that α satisfies the additive structure of model (3), which we write as

$$\alpha(t, Z_i(t)) = \alpha^* + \sum_{j=0}^d \alpha_j(X_{ij}(t))$$

with the notation $X_i(t) = (t, Z_{i1}(t), \dots, Z_{id}(t)) \in \mathcal{X}$ for $\mathcal{X} = [0, \mathcal{T}] \times [0, R]^d$. In the sequel, we will also write $x = (t, z_1, \dots, z_d) \in \mathcal{X}$ and henceforth $\alpha(x) = \alpha(t, z)$ for short.

Each component of the additive hazard α is only identifiable up to an additive shift. Later, we will give conditions under which each component is uniquely identified.

Model (4) allows for different kind of filtered data making it very flexible. These filterings include left-truncation and right-censoring which occurs in many applications of survival analysis [Martinussen and Scheike, 2006]. We now illustrate how to embed left-truncated covariates and right-censored survival time into model (4). Let T denote the survival time. Left-truncation means that we observe copies of (T, Z) only on a compact subset $\mathcal{I} \subseteq \mathcal{X}$ with the property that $(t_1, Z(t_1)) \in \mathcal{I}$ and $t_2 \geq t_1$ imply $(t_2, Z(t_2)) \in \mathcal{I}$ almost surely. We allow \mathcal{I} to be random but assume it is independent from T given Z . The survival time T can also be subject to right censoring with censoring time C as long as C is conditionally independent from T given the covariate process Z . This condition holds in particular if the censoring time equals one of the components of Z . Hence, under this filtering scheme, we observe n *i.i.d.* copies of $(\tilde{T}, Z^*, \mathcal{I}, \delta)$, where $\delta = \mathbb{1}(T^* < C)$, $\tilde{T} = \min(T^*, C)$, and (T^*, Z^*) is the truncated version of (T, Z) , i.e. (T^*, Z^*) arises from (T, Z) by conditioning on the event $\{(T, Z(T)) \in \mathcal{I}\}$.

We can now define a counting process N_i for each individual $i = 1, \dots, n$, via

$$N_i(t) = \mathbb{1}\{\tilde{T}_i \leq t, \delta_i = 1\},$$

with respect to the filtration $\mathcal{F}_{i,t} = \sigma\left(\left\{\tilde{T}_i \leq s, Z_i^*(s), \mathcal{I}_i, \delta_i : s \leq t\right\} \cup \mathcal{N}\right)$, for a class of null-sets \mathcal{N} , which completes the filtration. In this setting it can be easily shown that, under above assumption of $\alpha(t|Z) = \alpha(t, Z(t))$, Aalen's multiplicative intensity model (4) is satisfied with hazard rate

$$\alpha(t, z) = \lim_{h \downarrow 0} h^{-1} \mathbb{P}(T_i \in [t, t+h) | T_i \geq t, Z_i(t) = z),$$

and exposure

$$Y_i(t) = \mathbb{1}\{(t, Z_i^*(t)) \in \mathcal{I}_i, t \leq \tilde{T}_i\},$$

for individual i . The sets \mathcal{I}_i are allowed to be independent random copies of \mathcal{I} .

4. The smooth backfitting estimator of additive hazards

4.1. Smooth backfitting hazard estimator as projection

In this and the next section we illustrate the equivalence of projections and estimators that minimize squared errors following the line of Mammen et al. [1999b] where smooth backfitting was first introduced for nonparametric regression. The idea of describing smoothing estimators as projections in a regression setting is explained in great detail in Mammen et al. [2001]. In the following we introduce this projection principle for a counting process framework.

We will introduce our estimators as a projection from a functional space \mathcal{H} onto a certain subspace. The choice of the subspace, implies the class of functions that can be estimated and also the class of estimators to be considered. We now specify these functional spaces as well as (semi-)norms.

We define the unrestricted functional space as

$$\mathcal{H} = \{(f^{i,j})_{i=1,\dots,n,j=0,\dots,d+1}; f^{i,j} : \mathbb{R}^{d+2} \rightarrow \mathbb{R}\},$$

and subsets $\mathcal{H}_{full}^{LC} \subseteq \mathcal{H}_{full}^{LL} \subseteq \mathcal{H}$ via

$$\begin{aligned}\mathcal{H}_{full}^{LL} &= \{f \in \mathcal{H} : f^{i,j}(s, x) \text{ does not depend on } i, s\}, \\ \mathcal{H}_{full}^{LC} &= \{f \in \mathcal{H} : f^{i,j}(s, x) \text{ does not depend on } i, s; \\ &\quad f^{i,j}(s, x) \equiv 0 \text{ for } j = 1, \dots, d+1\}.\end{aligned}$$

Furthermore, for additive hazard functions we define additive subsets

$$\begin{aligned}\mathcal{H}_{add}^{LL} &= \{f \in \mathcal{H}_{full}^{LL} : f^{i,0}(s, x) = \sum_{j=0}^d g_j(x_j); f^{i,j}(s, x) = h_j(x_j), j = 1, \dots, d+1, \\ &\quad \text{for some functions } g_j, h_j : \mathbb{R} \rightarrow \mathbb{R}\}, \\ \mathcal{H}_{add}^{LC} &= \{f \in \mathcal{H}_{full}^{LC} : f^{i,0}(s, x) = \sum_{j=0}^d g_j(x_j) \text{ for some functions } g_j : \mathbb{R} \rightarrow \mathbb{R}\},\end{aligned}$$

that contain the class of local linear and local constant hazard estimators, respectively. Moreover, we define a semi-norm $\|\cdot\|$ on \mathcal{H} through

$$\begin{aligned}\|f\|^2 &= \int \int \frac{1}{n} \sum_{i=1}^n \left[f^{i,0}(s, x) + \sum_{j=0}^d f^{i,j+1}(s, x) \left(\frac{x_j - X_{i,j}(s)}{h} \right) \right]^2 \\ &\quad \times Y_i(s) K_h(x - X_i(s)) ds d\nu(x),\end{aligned}$$

for $f \in \mathcal{H}$ and where ν is a measure with strictly positive density. This semi-norm will be used to define the projection in the sequel.

Next we will illustrate how \mathcal{H} contains both hazard functions and the observations $(N_i), i = 1, \dots, n$. For every $\varepsilon > 0$, the data can be identified with an element $\Delta_\varepsilon N \in \mathcal{H}$ via

$$\Delta_\varepsilon N^{i,0}(s, x) = \frac{1}{\varepsilon} \int_s^{s+\varepsilon} dN_i(s), \quad \Delta_\varepsilon N^{i,j}(s, x) \equiv 0, \quad j = 1, \dots, d+1.$$

We define the unstructured local constant and local linear hazard estimator as

$$\lim_{\varepsilon \rightarrow 0} \arg \min_{\theta \in \mathcal{H}_{full}^{LC}} \|\Delta_\varepsilon N - \theta\|, \quad \lim_{\varepsilon \rightarrow 0} \arg \min_{\theta \in \mathcal{H}_{full}^{LL}} \|\Delta_\varepsilon N - \theta\|, \quad (5)$$

respectively. One can easily verify that these estimators coincide with the well known local constant and local linear hazard marker dependent hazard estimators introduced in Nielsen and Linton [1995] and Nielsen [1998].

For $\varepsilon \rightarrow 0$, each element $\Delta_\varepsilon N^{i,0}$ converges to a Dirac delta function. Hence, we write

$$\min_{\theta \in \mathcal{G}} \|\Delta N - \theta\| := \lim_{\varepsilon \rightarrow 0} \min_{\theta \in \mathcal{G}} \|\Delta_\varepsilon N - \theta\|,$$

for $\mathcal{G} \subseteq \mathcal{H}$.

We define the local constant and local linear nonparametric additive hazard estimator respectively as

$$\arg \min_{\theta \in \mathcal{H}_{add}^{LC}} \|\Delta N - \theta\|, \quad \arg \min_{\theta \in \mathcal{H}_{add}^{LL}} \|\Delta N - \theta\|. \quad (6)$$

For the minimisation over all additive hazard functions, we can either use a direct projection into \mathcal{H}_{add}^P , $P \in \{LC, LL\}$ which is given by $\min_{\theta \in \mathcal{H}_{add}^P} \|\Delta N - \theta\|$ or we use a Pythagorean argument to project in two steps: For $\hat{\alpha} \in \mathcal{H}_{add}^P$, it holds $\|\Delta N - \hat{\alpha}\|^2 = \|\Delta N - \tilde{\alpha}\|^2 + \|\tilde{\alpha} - \hat{\alpha}\|^2$ with $\tilde{\alpha} \in \mathcal{H}_{full}^P$. The last identity holds because the elements $\Delta N - \tilde{\alpha}$ and $\tilde{\alpha} - \hat{\alpha}$ are orthogonal [Mammen et al., 2001]. In additive marker dependent hazard estimation, the unrestricted marker dependent hazard estimators can be understood as intermediate in an iterative projection procedure that first projects to the unrestricted space and then to the additive space.

4.2. Smooth backfitting hazard estimator via least squares

In the previous section, we introduced the local constant estimator as a projection from \mathcal{H} . In this section, we show how this connects to the more known least squares criteria, and thereby also state the estimator in a way that is more directly mathematically tractable. We first consider the unstructured local polynomial hazard estimators. For a general understanding, we write down the general formulation for polynomials of order p , but in this paper we will only consider the local constant and the local linear case, $p = 0, 1$.

We will estimate the additive components of the hazard function via kernel smoothers. Let $k : \mathbb{R} \rightarrow \mathbb{R}$ be a symmetric and continuous kernel function such that $\int k(u) du = 1$. We define $K(u_0, \dots, u_d) = \prod_{j=0}^d k(u_j)$. For a smoothing parameter $h > 0$, $K_h(u) = \prod_{j=0}^d k_h(u_j) = \prod_{j=0}^d h^{-1} k(h^{-1} u_j)$. In the sequel, we will use a modification of the kernel function to ensure that the kernel always integrates to unity. We replace $k_h(u - v)$ by

$$k_h(u, v) = I_{(u, v \in [0, 1])} \left(\int k_h(s - v) ds \right)^{-1} k_h(u - v) \quad (7)$$

for every $h > 0$ to correct for normalization at the boundaries from now on. Furthermore, we define the multivariate kernel

$$K_h(u, v) = \prod_{j=0}^d k_h(u_j, v_j),$$

for $u = (u_0, \dots, u_d)$ and $v = (v_0, \dots, v_d)$.

The unstructured p th order local polynomial estimator of the hazard function in x is defined as the

first component of

$$\begin{aligned}
& \lim_{\varepsilon \rightarrow 0} \arg \min_{\substack{\theta_0: \mathbb{R}^{d+1} \rightarrow \mathbb{R} \\ \theta_j: \mathbb{R}^{d+1} \rightarrow \mathbb{R}^{d+1} \\ j=1, \dots, p}} \sum_{i=1}^n \int \int \left\{ \frac{1}{\varepsilon} \int_s^{s+\varepsilon} dN_i(u) - \theta_0(x) \right. \\
& \quad - \theta_1^T(x) \left(\frac{x_0 - X_{i0}(s)}{h}, \dots, \frac{x_d - X_{id}(s)}{h} \right)^T - \dots \\
& \quad \left. - \theta_p^T(x) \left(\left(\frac{x_d - X_{id}(s)}{h} \right)^p, \dots, \left(\frac{x_d - X_{id}(s)}{h} \right)^p \right)^T \right\}^2 \\
& \quad \times K_h(x, X_i(s)) Y_i(s) ds d\nu(x),
\end{aligned} \tag{8}$$

The cases $p = 0, 1$ are exactly the local constant and local linear projection estimator defined in (5).

For the rest of this paper, we limit ourselves to the same kernel k and bandwidth h for each dimension to keep the notation simple. Henceforth, if there is no confusion about the boundaries of the integrals, \int denotes integration over the whole support $[0, \mathcal{T}] \times [0, R]^d$. The measure ν has to have a strictly positive density but the estimator does not depend on the specific choice of ν if we don't have restrictions on the functions θ_j . We will specify a weighting function w such that $d\nu(x) = w(x)dx$. Note that this estimator allows for local polynomial approximation at degree p but it is not additive yet.

The nonparametric additive hazard estimator we investigate in this paper is defined by the minimisation in equation (8) under the following constraints on the structural form of θ . For $p = 0$, the constraint $\theta_0(x) = \bar{\alpha}^* + \sum_{j=0}^d \bar{\alpha}_j(x_j)$ for some functions $\bar{\alpha}_0, \dots, \bar{\alpha}_d$ and a constant $\bar{\alpha}^*$, leads to the local constant estimator as introduced in (6):

$$\begin{aligned}
& \lim_{\varepsilon \rightarrow 0} \arg \min_{\substack{\bar{\alpha}^* \in \mathbb{R}, \\ \bar{\alpha}_j: \mathbb{R} \rightarrow \mathbb{R}, \\ j=0, \dots, d}} \sum_{i=1}^n \int \int \left\{ \frac{1}{\varepsilon} \int_s^{s+\varepsilon} dN_i(u) - [\bar{\alpha}^* + \bar{\alpha}_0(t) + \bar{\alpha}_1(z_1) + \dots + \bar{\alpha}_d(z_d)] \right\}^2 \\
& \quad \times K_h(x, X_i(s)) Y_i(s) ds d\nu(x).
\end{aligned} \tag{9}$$

For the unique identification of the constant component α^* and the components α_j , $j = 0, \dots, d$, we will set further constraints in equation (13).

The local linear additive hazard estimator as defined in (6) arises by setting $\theta_0(x) = \bar{\alpha}^* + \sum_{j=0}^d \bar{\alpha}_j(x_j)$ and $\theta_1(x) = (\partial/\partial x_0 \theta_0(x), \dots, \partial/\partial x_d \theta_0(x))$.

$$\begin{aligned}
& \lim_{\varepsilon \rightarrow 0} \arg \min_{\substack{\bar{\alpha}^* \in \mathbb{R}, \\ \bar{\alpha}_j: \mathbb{R} \rightarrow \mathbb{R}, \\ \bar{\alpha}'_j: \mathbb{R} \rightarrow \mathbb{R}, \\ j=0, \dots, d}} \sum_{i=1}^n \int \int \left\{ \frac{1}{\varepsilon} \int_s^{s+\varepsilon} dN_i(u) - \left[\bar{\alpha}^* + \bar{\alpha}_0(t) + \bar{\alpha}_1(z_1) + \dots + \bar{\alpha}_d(z_d) \right. \right. \\
& \quad \left. \left. + \bar{\alpha}'_0(x_0) \left(\frac{x_0 - X_{i0}(s)}{h} \right) + \dots + \bar{\alpha}'_d(x_d) \left(\frac{x_d - X_{id}(s)}{h} \right) \right] \right\}^2 \\
& \quad \times K_h(x, X_i(s)) Y_i(s) ds d\nu(x).
\end{aligned} \tag{10}$$

Existence and uniqueness of the minimizers of (9) and (10) will be established later.

4.3. The local constant smooth backfitting additive kernel hazard estimator

The minimisation in equation (8) for $p = 0$ leads to the unstructured local constant estimator $\hat{\alpha}^{LC}$ defined via $\hat{\alpha}^{LC}(x) = \hat{O}(x)/\hat{E}(x)$ with

$$\hat{O}(x) = \frac{1}{n} \sum_{i=1}^n \int K_h(x, X_i(s)) dN_i(s),$$

$$\hat{E}(x) = \frac{1}{n} \sum_{i=1}^n \int K_h(x, X_i(s)) Y_i(s) ds.$$

for $x \in \mathcal{X}$. The estimators \hat{O} and \hat{E} estimate the occurrence and exposure of the observations. The exposure E is defined via $E(x) = f_t(z) \mathbb{E}[Y(t)]$ where $f_t(z)$ is the conditional density of $(Z_1(t), \dots, Z_d(t))$ given $Y(t) = 1$. The occurrence is defined as $O(x) = \alpha(x)E(x)$ for $x = (t, z) \in \mathcal{X}$. The structure of a hazard estimator as an estimator of occurrence divided by an estimator of exposure is in line with piece-wise constant hazard estimators in Martinussen and Scheike [2002].

To define the local constant smooth backfitting additive hazard estimators we proceed as follows. Following the derivation in Section 4.2, the estimator is defined through equation (9). The solution $\bar{\alpha} = (\bar{\alpha}^*, \bar{\alpha}_0, \dots, \bar{\alpha}_d)$ satisfies the first order conditions

$$\bar{\alpha}^* = \frac{\int_{\mathcal{X}} [\hat{\alpha}^{LC}(x) - \sum_{j=0}^d \bar{\alpha}_j(x_j)] w(x) dx}{\int_{\mathcal{X}} w(x) dx} \quad (11)$$

and

$$\bar{\alpha}_k(x_k) = \int_{\mathcal{X}_{x_k}} \hat{\alpha}^{LC}(x) \frac{w(x)}{w_k(x_k)} dx_{-k} - \sum_{j \neq k} \int_{\mathcal{X}_{x_k}} \bar{\alpha}_j(x_j) \frac{w(x)}{w_k(x_k)} dx_{-k} - \bar{\alpha}^*, \quad (12)$$

for $k = 0, \dots, d$, where we write $w_k(x_k) = \int_{\mathcal{X}_{x_k}} w(x) dx_{-k}$ for the marginals of w using the notation $\mathcal{X}_{x_k} = \{y \in \mathcal{X} : y_x = x_k\}$ and dx_{-k} denoting integration over all components except for k . For the unique identification of the solution we also set the conditions

$$\int_{\mathcal{X}_k} \bar{\alpha}_k(x_k) w_k(x_k) dx_k = 0, \quad k = 0, \dots, d. \quad (13)$$

These identification conditions enable us further to get

$$\bar{\alpha}^* = \frac{\int_{\mathcal{X}} \hat{\alpha}^{LC}(x) w(x) dx}{\int_{\mathcal{X}} w(x) dx} = \frac{\int_{\mathcal{X}} \hat{O}(x) dx}{\int_{\mathcal{X}} \hat{E}(x) dx}$$

from equation (11), where the second equality arises from the definition of $\hat{\alpha}$ and if we set the weighting to $w(x) = \hat{E}(x)$. One can further reduce the estimator to

$$\bar{\alpha}^* = \frac{\sum_{i=1}^n \int dN_i(s)}{\sum_{i=1}^n \int Y_i(s) ds}. \quad (14)$$

This simplification is due to the normalization $\int K_h(x, X_i(s)) dx = 1$ of the kernel function K_h in (7). The estimator $\bar{\alpha}^*$ is the additive hazard equivalent of the intercept in nonparametric regression. Note that in backfitting of the regression function m in Mammen et al. [1999b], the estimator for the additive constant m_0 of the conditional mean m is given as $\tilde{m}_0 = \bar{Y}_n$. Our result for $\bar{\alpha}^*$ is the total number of occurrences divided by the average exposure time. In the case of non-filtered data, $\int dN_i(s)$ equals unity for every i and thus $\bar{\alpha}^* = \left(\frac{1}{n} \sum_{i=1}^n \int Y_i(s) ds\right)^{-1}$. This term is the natural survival analysis equivalent to what the empirical mean is in regression.

The constant component α^* and all components α_j of the unknown underlying hazard α are uniquely identified through

$$\int \alpha_j(x_j) E_j(x_j) dx_j = 0 \quad (15)$$

with $E_j(x_j) = \int E(x) dx_{-j}$ for all j . This motivates the choice $w(x) = \hat{E}(x)$ in equation (13) and the notation $\hat{E}_k(x_k)$ instead of $w_k(x_k)$ for this choice of weighting from now on.

For the same data-adaptive weighting we simplify the terms in equation (12) with some new notation. Analogously to the one-dimensional marginals, we write $\hat{E}_{k,j}(x_k, x_j) = \int_{\mathcal{X}_{x_k, x_j}} \hat{E}(x) dx_{-(k,j)}$

for $x_{-(k,j)} = (x_0, \dots, x_{j-1}, x_{j+1}, \dots, x_{k-1}, x_{k+1}, \dots, x_d)$ and $\mathcal{X}_{x_k, x_j} = \{(x'_0, \dots, x'_d) \in \mathcal{X} : x'_k = x_k, x'_j = x_j\}$, i.e. we integrate over all components except for x_j and x_k which are fixed values. Analogously, we define the marginal occurrence estimator $\hat{O}_k(x_k) = \int_{\mathcal{X}_{x_k}} \hat{O}(x) dx_{-k}$.

In the local constant case investigated here, it can be easily shown that it holds

$$\hat{O}_k(x_k) = \frac{1}{n} \sum_{i=1}^n \int k_h(x_k, X_{ik}(s)) dN_i(s), \quad (16)$$

$$\hat{E}_k(x_k) = \frac{1}{n} \sum_{i=1}^n \int k_h(x_k, X_{ik}(s)) Y_i(s) ds, \quad (17)$$

$$\hat{E}_{j,k}(x_j, x_k) = \frac{1}{n} \sum_{i=1}^n \int k_h(x_j, X_{ij}(s)) k_h(x_k, X_{ik}(s)) Y_i(s) ds, \quad (18)$$

for $j \neq k$ if each pair of covariates has a rectangular support. Thus, these estimators are indeed just one- and two-dimensional marginal estimators and can be computed efficiently for high dimensions $d > 2$.

Now equation (12) implies the backfitting equation

$$\bar{\alpha}_k(x_k) = \hat{\alpha}_k(x_k) - \sum_{j \neq k} \int_{\mathcal{X}_j} \bar{\alpha}_j(x_j) \frac{\hat{E}_{k,j}(x_k, x_j)}{\hat{E}_k(x_k)} dx_j - \bar{\alpha}^*, \quad (19)$$

for the notation $\hat{\alpha}_k(x_k) = \hat{O}_k(x_k) / \hat{E}_k(x_k)$.

Using the last expression, we can get estimators for $\alpha_0, \dots, \alpha_d$ through iterative backfitting via

$$\begin{aligned} \bar{m}_k^{[r+1]}(x_k) &= \hat{\alpha}_k(x_k) - \sum_{j < k} \int \bar{\alpha}_j^{[r+1]}(x_j) \frac{\hat{E}_{k,j}(x_k, x_j)}{\hat{E}_k(x_k)} dx_j - \sum_{j > k} \int \bar{\alpha}_j^{[r]}(x_j) \frac{\hat{E}_{k,j}(x_k, x_j)}{\hat{E}_k(x_k)} dx_j, \\ \bar{\alpha}_k^{[r+1]}(x_k) &= \bar{m}_k^{[r+1]}(x_k) - \left(\int \hat{E}_k(x_k) dx_k \right)^{-1} \int \bar{m}_k^{[r+1]}(x_k) \hat{E}_k(x_k) dx_k, \end{aligned} \quad (20)$$

for $k = 1, \dots, d$ in step $r+1$. Recall that $\hat{\alpha}_k, k = 0, \dots, d$, are the (non-additive) estimators which were defined via $\hat{\alpha}_k(x_k) = \hat{O}_k(x_k) / \hat{E}_k(x_k)$. We suggest to start with the initialization $\bar{\alpha}_k^{[0]}(x_k) = \hat{\alpha}_k(x_k)$, that is related to the one-dimensional local linear hazard estimator, see Nielsen and Tanggaard [2001]. However, these pilot estimators can be set to different estimators. The asymptotic theory we present here is illustrated for the choice $\hat{\alpha}_k$. In Section A.3 of the appendix, we illustrate how one can obtain the same estimator $\bar{\alpha}_k$ by first minimizing (8) without an additive constraint, yielding the pilot estimator $\hat{\alpha}_k$ and then running an additive minimisation of $\hat{\alpha}_k$.

The complete smooth backfitting algorithm for the local constant additive hazard estimator $\bar{\alpha}$ is as follows.

1. Compute \hat{O}_k, \hat{E}_k , and $\hat{E}_{j,k}$ from equations (16)–(18) and set $\hat{\alpha}_k(x_k) = \hat{O}_k(x_k) / \hat{E}_k(x_k)$ for $k, j = 0, \dots, d$.
2. Set $r = 0$ and $\bar{\alpha}_k^{[r]} = \hat{\alpha}_k$ for $k = 0, \dots, d$.
3. For $k = 0, \dots, d$, compute $\bar{\alpha}_k^{[r+1]}(x_k)$ via equation (20) for all points x_k .
4. If the convergence criterion

$$\frac{\sum_{k=0}^d \int \left(\bar{\alpha}_k^{[r+1]}(x_k) - \bar{\alpha}_k^{[r]}(x_k) \right)^2 dx_k}{\sum_{k=0}^d \int \left(\bar{\alpha}_k^{[r+1]}(x_k) \right)^2 dx_k + 0.0001} < 0.0001$$

is fulfilled, stop; otherwise set r to $r + 1$ and go to step 3.

5. After convergence in step r , set $\bar{\alpha}_k = \bar{\alpha}_k^{[r+1]}$ for $k = 0, \dots, d$, and $\bar{\alpha}^* = \sum_{i=1}^n \int dN_i(s) / \sum_{i=1}^n \int Y_i(s) ds$.

Note that the quantities $\hat{E}_{j,k}(x_j, x_k)$, $\hat{E}_k(x_k)$, $\hat{\alpha}(x_k)$, and $\bar{\alpha}^*$ can be calculated once in the beginning and they are not updated during the iteration process. This is a computational advantage. However, we want to emphasize that the downside of the analogue local linear approach to this section is that the local linear pilot estimator does not necessarily exist for low numbers of observations in high dimensions. The local constant estimator on the other hand suffers from bad performance at boundaries.

4.4. Asymptotic properties of the local constant smooth backfitting additive kernel hazard estimator

We now derive the asymptotic behavior of the local constant estimator under weak assumptions. Indeed, we don't assume existence of \hat{O} , \hat{E} but only existence of some one- and two-dimensional marginal estimators $\hat{O}_k, \hat{O}_{k,j}, \hat{E}_k, \hat{E}_{k,j}, j, k = 0, \dots, d$, which is satisfied under the conditions illustrated below.

The following conditions are sufficient to derive asymptotic normality of the resulting smooth backfitting estimators $\bar{\alpha}_j, j = 0, \dots, d$.

A1 The exposure satisfies $\inf_{x \in \mathcal{X}} E(x) > 0$ and its marginals E_j are differentiable for every j . Moreover, the conditional density f_t of Z given $Y(t) = 1$ is continuous for every $t \in [0, T]$ and it holds $\sup_{x \in \mathcal{X}} f_t(x) < C_f$ for some constant C_f .

A2 There exists a function $\gamma \in C^2([0, T])$ such that it holds $n^{-1} \sum_{i=1}^n Y_i(t) \rightarrow \gamma(t)$ in probability as $n \rightarrow \infty$ for every $t \in [0, T]$.

A3 The function k is a second order kernel, that is it satisfies $\int k(u) du = 1, \int uk(u) du = 0$. Furthermore, k is a symmetric and Lipschitz continuous function with support $[-1, 1]$.

A4 It holds $n^{1/5}h \rightarrow c_h$ for a constant $0 < c_h < \infty$ as $n \rightarrow \infty$.

A5 The hazard α is two times continuously differentiable in every component of $x \in \mathcal{X}$.

Note that in our notation $\gamma(t)$ from A2 and $E_0(t)$ are almost surely identical. However, the definition of E_0 does not assure $E_0 \in C^2([0, T])$ without A2.

Theorem 1 (Local constant smooth backfitting estimator). *Let $\hat{\alpha}_j = \hat{O}_j / \hat{E}_j$ be the pilot estimator for $j = 0, \dots, d$. Under Assumptions A1–A5, with probability tending to 1, there exists a unique solution $\{\bar{\alpha}^*, \bar{\alpha}_j : j = 0, \dots, d\}$ to (9), and the backfitting algorithm converges to it:*

$$\int \left[\bar{\alpha}_j^{[r]}(x_j) - \bar{\alpha}_j(x_j) \right]^2 E_j(x_j) dx_j \rightarrow 0.$$

For $x_0 \in (0, T)$ and $x_l \in (0, R)$, $l = 1, \dots, d$, the solution satisfies

$$n^{2/5} \left\{ \begin{pmatrix} \bar{\alpha}_0(x_0) - \alpha_0(x_0) \\ \vdots \\ \bar{\alpha}_d(x_d) - \alpha_d(x_d) \end{pmatrix} \right\} \rightarrow \mathcal{N} \left(\begin{pmatrix} c_h^2 b_0(x_0) \\ \vdots \\ c_h^2 b_d(x_d) \end{pmatrix}, \begin{pmatrix} v_0(x_0) & 0 & \cdots & 0 \\ 0 & \ddots & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \cdots & 0 & v_d(x_d) \end{pmatrix} \right),$$

and in particular $\bar{\alpha}(x) = \bar{\alpha}^* + \sum_{j=0}^d \bar{\alpha}_j$ with $\bar{\alpha}^*$ from equation (14) satisfies

$$n^{2/5} \{ \bar{\alpha}(x) - \alpha(x) \} \rightarrow \mathcal{N} \left(c_h^2 \sum_{j=0}^d b_j(x_j), \sum_{j=0}^d v_j(x_j) \right),$$

for $n \rightarrow \infty$, where

$$v_j(x_j) = c_h^{-1} \int k(u)^2 du \sigma_j^2(x_j) E_j(x_j)^{-1},$$

$$\sigma_j^2(x_j) = \alpha^* E_j(x_j)^{-1} + \sum_{l \neq j} \int \alpha_l(u) E_{jl}(x_j, u) E_j(x_j)^{-1} du + \alpha_j(x_j).$$

and where b_j is given through

$$(b_0, b_1, \dots, b_d) = \arg \min_{\mathcal{B}} \int [\beta(x) - \beta_0 - \beta_1(x_1) - \dots - \beta_d(x_d)]^2 E(x) dx,$$

for

$$\beta(x) = \sum_{j=0}^d \int u^2 k(u) du \left[\alpha'_j(x_j) \frac{\partial \log E(x)}{\partial x_j} + \frac{1}{2} \alpha''_j(x_j) \right],$$

and $\mathcal{B} = \{\tilde{\beta} = (\beta_0, \beta_1, \dots, \beta_d) : \int \beta_j(x_j) E_j(x_j) dx_j = 0; j = 0, \dots, d\}$.

The proof of Theorem 1 is given in Appendix A.1.

Remark 1. Define the martingale $M_i = N_i - \Lambda_i$ where Λ_i is the compensator of N_i . The term $\int k(u)^2 du \sigma_j^2(x_j) E_j(x_j)$ occurs as the asymptotic variance of the martingale $\int k_h(x_j, X_{ij}(s)) dM_i(s)$. The convergence rate is the same as for a one-dimensional local constant hazard estimator, see e.g. Nielsen and Tanggaard [2001]. In the nonparametric regression setting $Y = m(X) + \varepsilon$ of Mammen et al. [1999b], and in contrast to our hazard estimator, the asymptotic variance under certain regularity conditions is specified through $\sigma_j^2(x_j) = \text{Var}(Y - m(X) | X_j = x_j)$ without any closed form expression.

Remark 2. By Lemma 1 in the appendix, $\tilde{\alpha}^*$ is an unbiased estimator of α^* if the identification conditions $\int \alpha_j(x_j) E_j(x_j) dx_j = 0$ hold for $j = 0, \dots, d$.

4.5. The local linear smooth backfitting additive kernel hazard estimator

The local linear smooth backfitting estimator $\tilde{\alpha}_j(x_j)$ for $j = 0, \dots, d$, can be described by the minimisation in equation (10). As described in Section 4.2, this is equivalent to the minimisation in (8) for $p = 1$ with respect to $(\hat{\alpha}, \hat{\alpha}^{(1)})$ under the constraints $\theta_0(x) = \hat{\alpha}^* + \sum_{j=0}^d \hat{\alpha}_j(x_j)$, $\theta_{1,j}(x_j) = \hat{\alpha}_j^{(1)}(x_j)$ for a certain weighting function w .

Denoting the estimator of derivatives α'_j by $\tilde{\alpha}^j$ in the following, the first order conditions for the minimisation in $\tilde{\alpha}_j(x_j) + \tilde{\alpha}^*$ and $\tilde{\alpha}^j(x_j)$ can be written as

$$[\tilde{\alpha}_j(x_j) + \tilde{\alpha}^*] \hat{V}^j(x_j) + \tilde{\alpha}^j(x_j) \hat{V}_{j,j}^j(x_j) = \frac{1}{n} \sum_{i=1}^n \int k_h(x_j, X_{ij}(s)) dN_i(s)$$

$$- \sum_{l \neq j} \int \tilde{\alpha}_l(x_l) \hat{V}^{l,j}(x_l, x_j) dx_l$$

$$- \sum_{l \neq j} \int \tilde{\alpha}^l(x_l) \hat{V}_l^{l,j}(x_l, x_j) dx_l \quad (21)$$

$$[\tilde{\alpha}_j(x_j) + \tilde{\alpha}^*] \hat{V}_{j,j}^j(x_j) + \tilde{\alpha}^j(x_j) \hat{V}_{j,j}^j(x_j) = \frac{1}{n} \sum_{i=1}^n \int \left(\frac{x_j - X_{i,j}(s)}{h} \right) k_h(x_j, X_{ij}(s)) dN_i(s),$$

$$- \sum_{l \neq j} \int \tilde{\alpha}_l(x_l) \hat{V}_{j,j}^{l,j}(x_l, x_j) dx_l \quad (22)$$

$$- \sum_{l \neq j} \int \tilde{\alpha}^l(x_l) \hat{V}_{l,j}^{l,j}(x_l, x_j) dx_l,$$

with the new notation

$$\hat{V}^j(x_j) = \frac{1}{n} \sum_{i=1}^n \int k_h(x_j, X_{ij}(s)) Y_i(s) ds, \quad (23)$$

$$\hat{V}^{l,j}(x_l, x_j) = \frac{1}{n} \sum_{i=1}^n \int k_h(x_l, X_{il}(s)) k_h(x_j, X_{ij}(s)) Y_i(s) ds,$$

$$\hat{V}_j^j(x_j) = \frac{1}{n} \sum_{i=1}^n \int \left(\frac{x_j - X_{i,j}(s)}{h} \right) k_h(x_j, X_{ij}(s)) Y_i(s) ds,$$

$$\hat{V}_l^{l,j}(x_l, x_j) = \frac{1}{n} \sum_{i=1}^n \int \left(\frac{x_l - X_{i,l}(s)}{h} \right) k_h(x_l, X_{il}(s)) k_h(x_j, X_{ij}(s)) Y_i(s) ds,$$

$$\hat{V}_j^{l,j}(x_l, x_j) = \frac{1}{n} \sum_{i=1}^n \int \left(\frac{x_j - X_{i,j}(s)}{h} \right) k_h(x_l, X_{il}(s)) k_h(x_j, X_{ij}(s)) Y_i(s) ds,$$

$$\hat{V}_{j,j}^j(x_j) = \frac{1}{n} \sum_{i=1}^n \int \left(\frac{x_j - X_{i,j}(s)}{h} \right)^2 k_h(x_j, X_{ij}(s)) Y_i(s) ds,$$

$$\hat{V}_{l,j}^{l,j}(x_l, x_j) = \frac{1}{n} \sum_{i=1}^n \int \left(\frac{x_l - X_{i,l}(s)}{h} \right) \left(\frac{x_j - X_{i,j}(s)}{h} \right) k_h(x_l, X_{il}(s)) k_h(x_j, X_{ij}(s)) Y_i(s) ds. \quad (24)$$

Here, x_{-k} denotes $(x_0, \dots, x_{k-1}, x_{k+1}, \dots, x_d)$ and \mathcal{X}_{x_k} denotes the set $\{(x'_0, \dots, x'_d) \in \mathcal{X} : x'_k = x_k\}$.

Note that $\hat{V}^j(x_j)$ and $\hat{V}^{l,j}(x_l, x_j)$ are identical to the one- and two-dimensional local constant fits $\hat{E}_j(x_j)$ and $\hat{E}_{j,k}(x_j, x_k)$ from the local constant estimator. For simplicity of notation, we relabel them in the sequel. The terms $\hat{V}_j^j(x_j)$, $\hat{V}_l^{l,j}(x_l, x_j)$, $\hat{V}_j^{l,j}(x_l, x_j)$, $\hat{V}_{j,j}^j(x_j)$ and $\hat{V}_{l,j}^{l,j}(x_l, x_j)$ contain linear and quadratic components, which distinguish this approach from the one in the last section.

Furthermore, for $j = 0, \dots, d$ we introduce the same identification condition as equation (13) in the local constant case and require

$$\int \tilde{\alpha}_j(x_j) \hat{V}^j(x_j) dx_j = 0 \quad (25)$$

to get a unique solution of (21) and (22).

We can derive a local constant estimator from the same conditions (21) and (22) for $\hat{\alpha}_k(x_k)$ but with $\hat{\alpha}_j'(x_j)$ set to zero for every j . If we choose $w \equiv 1$, this local constant estimator coincides with the one from Section 4.3.

Conditions (21)–(25) uniquely define our estimator and for the derivation of asymptotic theory (21)–(22) can be written in one equation as

$$\hat{M}_j(x_j) \begin{pmatrix} \tilde{\alpha}_j(x_j) - \hat{\alpha}_j(x_j) \\ \tilde{\alpha}_j'(x_j) - \hat{\alpha}_j'(x_j) \end{pmatrix} = -\tilde{\alpha}^* \begin{pmatrix} \hat{V}^j(x_j) \\ \hat{V}_j^j(x_j) \end{pmatrix} - \sum_{l \neq j} \int \hat{S}_{l,j}(x_l, x_j) \begin{pmatrix} \tilde{\alpha}_l(x_l) \\ \tilde{\alpha}_l'(x_l) \end{pmatrix} dx_l, \quad (26)$$

where we have used the matrices

$$\hat{M}_j(x_j) = \begin{pmatrix} \hat{V}^j(x_j) & \hat{V}_j^j(x_j) \\ \hat{V}_j^j(x_j) & \hat{V}_{j,j}^j(x_j) \end{pmatrix}, \quad (27)$$

$$\hat{S}_{l,j}(x_l, x_j) = \begin{pmatrix} \hat{V}^{l,j}(x_l, x_j) & \hat{V}_l^{l,j}(x_l, x_j) \\ \hat{V}_j^{l,j}(x_l, x_j) & \hat{V}_{l,j}^{l,j}(x_l, x_j) \end{pmatrix}, \quad (28)$$

and the one-dimensional local linear fit of the observations

$$\begin{pmatrix} \hat{\alpha}_j(x_j) \\ \hat{\alpha}^j(x_j) \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n \int \hat{M}_j(x_j)^{-1} \begin{pmatrix} 1 \\ h^{-1}(x_j - X_{ij}(s)) \end{pmatrix} k_h(x_j, X_{ij}(s)) dN_i(s).$$

Note, that we would get the same asymptotic result for any estimator which arises from equation (26) by replacing $\hat{V}_{0,0}^j$, $\hat{V}_{0,0}^j$ and $(\hat{\alpha}_j, \hat{\alpha}^j)$ with asymptotically equivalent estimators that satisfy the same regularity conditions in Appendix A.2.

For the implementation as an iterative algorithm, step $r + 1$ of the backfitting algorithm is given by:

$$\begin{pmatrix} \hat{m}_j(x_j) \\ \tilde{\alpha}^{[r+1],j}(x_j) \end{pmatrix} = \begin{pmatrix} \hat{\alpha}_j(x_j) \\ \hat{\alpha}^j(x_j) \end{pmatrix} - \hat{M}_j(x_j)^{-1} \sum_{l \neq j} \int \hat{S}_{l,j}(x_l, x_j) \begin{pmatrix} \tilde{\alpha}_l^{[r]}(x_l) \\ \tilde{\alpha}^{[r],l}(x_l) \end{pmatrix} dx_l, \quad (29)$$

$$\tilde{\alpha}_j^{[r+1]}(x_j) = \hat{m}_j(x_j) - \left(\int \hat{V}^j(u_j) du_j \right)^{-1} \int \hat{m}_j(u_j) \hat{V}^j(u_j) du_j, \quad (30)$$

for $r = 0, 1, 2, \dots$

Note that $\tilde{\alpha}^*$ from equation (26) vanishes in the component $\alpha^{[r+1],j}(x_j)$ and it is made redundant in the other component by the norming condition (30). Theorem 2 assures the convergence of this estimator.

We recommend avoiding the inverse of the matrices \hat{M}_j in the implementation for computational stability. Solving equations (21)–(22) for $\tilde{\alpha}_j(x_j)$ and $\tilde{\alpha}^j(x_j)$, respectively, and first replacing $\tilde{\alpha}^j(x_j)$ in (21) by its latest fit $\tilde{\alpha}^{[r],j}(x_j)$ and then $\tilde{\alpha}_j(x_j)$ in (22) by $\tilde{\alpha}_j^{[r+1]}(x_j)$ in step $r + 1$, we get the asymptotically equivalent, more stable backfitting equations

$$\begin{aligned} \tilde{\alpha}_j^{[r+1]}(x_j) &= \hat{V}_j^j(x_j)^{-1} \left(\hat{U}_j^j(x_j) - \tilde{\alpha}^{[r],j}(x_j) \hat{V}_j^j(x_j) - \tilde{\alpha}^* \hat{V}_j^j(x_j), \right. \\ &\quad \left. - \sum_{l \neq j} \int \tilde{\alpha}_l^{[r]}(x_l) \hat{V}_l^{l,j}(x_l, x_j) dx_l - \sum_{l \neq j} \int \tilde{\alpha}^{[r],l}(x_l) \hat{V}_l^{l,j}(x_l, x_j) dx_l \right), \end{aligned} \quad (31)$$

$$\begin{aligned} \tilde{\alpha}^{[r+1],j}(x_j) &= \hat{V}_{j,j}^j(x_j)^{-1} \left(\hat{U}_j^j(x_j) - \tilde{\alpha}_j^{[r]}(x_j) \hat{V}_j^j(x_j) - \tilde{\alpha}^* \hat{V}_j^j(x_j) \right. \\ &\quad \left. - \sum_{l \neq j} \int \tilde{\alpha}_l^{[r+1]}(x_l) \hat{V}_l^{l,j}(x_l, x_j) dx_l - \sum_{l \neq j} \int \tilde{\alpha}^{[r],l}(x_l) \hat{V}_l^{l,j}(x_l, x_j) dx_l \right), \end{aligned} \quad (32)$$

for step $r + 1$ with the notation

$$\hat{U}_j^j(x_j) = \frac{1}{n} \sum_{i=1}^n \int k_h(x_j, X_{ij}(s)) dN_i(s), \quad (33)$$

$$\hat{U}_j^j(x_j) = \frac{1}{n} \sum_{i=1}^n \int \left(\frac{x_j - X_{ij}(s)}{h} \right) k_h(x_j, X_{ij}(s)) dN_i(s). \quad (34)$$

Note that $\hat{U}_j^j(x_j)$ is identical to $\hat{O}_j(x_j)$, the local constant occurrence estimator described in Section 4.3. We set the initialization in step $r = 0$ to $(\tilde{\alpha}_j^{[0]}(x_j), \tilde{\alpha}^{[0],j}(x_j)) = (0, 0)$.

The complete smooth backfitting algorithm for the local linear additive hazard estimator $\tilde{\alpha}$ is as follows.

1. Compute \hat{V}^j , $\hat{V}^{l,j}$, \hat{V}_j^j , $\hat{V}_l^{l,j}$, $\hat{V}_{j,j}^j$, and $\hat{V}_{l,j}^{l,j}$ from equations (23)–(24) and set $\hat{\alpha}(x_k) = \hat{O}_k(x_k) / \hat{E}_k(x_k)$ for $k, j = 0, \dots, d$.
2. Set $r = 0$ and $\tilde{\alpha}_k^{[r]} = \hat{\alpha}_k$ for $k, j = 0, \dots, d$.

3. For $k = 0, \dots, d$, calculate for all points x_k Set $r = 1$, compute $\tilde{\alpha}_k^{[r+1]}(x_k)$ via equations (31) and (32). Then replace $\tilde{\alpha}_j^{[r+1]}$ by

$$\tilde{\alpha}_j^{[r+1]} - \left(\int \hat{V}^j(u_j) du_j \right)^{-1} \int \tilde{\alpha}_j^{[r*]}(u_j) \hat{V}^j(u_j) du_j.$$

4. If the convergence criterion

$$\frac{\sum_{k=0}^d \int \left(\tilde{\alpha}_k^{[r+1]}(x_k) - \tilde{\alpha}_k^{[r]}(x_k) \right)^2 dx_k}{\sum_{k=0}^d \int \left(\tilde{\alpha}_k^{[r+1]}(x_k) \right)^2 dx_k + 0.0001} < 0.0001$$

is fulfilled, stop; otherwise set r to $r + 1$ and go to step 3.

5. After convergence in step r , set $\tilde{\alpha}_k = \tilde{\alpha}_k^{[r+1]}$ for $k = 0, \dots, d$, and $\tilde{\alpha}^* = \sum_{i=1}^n \int dN_i(s) / \sum_{i=1}^n \int Y_i(s) ds$.

4.6. Asymptotic properties of the local linear smooth backfitting additive kernel hazard estimator

For the asymptotic behavior of $\tilde{\alpha}_j$, we assume the same Assumptions A1–A5 as for the local constant estimator.

Theorem 2 (Local linear smooth backfitting estimator). *Under Assumptions A1–A5, with probability tending to 1, there exists a unique solution $\{\tilde{\alpha}_j, \tilde{\alpha}^j : j = 0, \dots, d\}$ to (10) and the backfitting algorithm (29) converges to it:*

$$\begin{aligned} \int \left[\tilde{\alpha}_j^{[r]}(x_j) - \tilde{\alpha}_j(x_j) \right]^2 E_j(x_j) dx_j &\rightarrow 0, \\ \int \left[\tilde{\alpha}^{j,[r]}(x_j) - \tilde{\alpha}^j(x_j) \right]^2 E_j(x_j) dx_j &\rightarrow 0. \end{aligned}$$

For $x_0 \in (0, \mathcal{T})$ and $x_l \in (0, R)$, $l = 1, \dots, d$, the solution satisfies

$$n^{2/5} \left\{ \begin{pmatrix} \tilde{\alpha}_0(x_0) - \alpha_0(x_0) + \nu_{n,0} \\ \vdots \\ \tilde{\alpha}_d(x_d) - \alpha_d(x_d) + \nu_{n,d} \end{pmatrix} \right\} \rightarrow \mathcal{N} \left(\begin{pmatrix} c_h^2 b_0(x_0) \\ \vdots \\ c_h^2 b_d(x_d) \end{pmatrix}, \begin{pmatrix} v_0(x_0) & 0 & \cdots & 0 \\ 0 & \ddots & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \cdots & 0 & v_d(x_d) \end{pmatrix} \right),$$

for $n \rightarrow \infty$, where

$$\begin{aligned} \nu_{n,j} &= \int \int \alpha_j(x_j) k_h(x_j, u) E_j(u) du dx_j, \\ b_j(x_j) &= \frac{1}{2} \int u^2 k(u) du \left[\alpha_j''(x_j) - \int \alpha_j''(x_j) E_j(x_j) dx_j \right], \\ v_j(x_j) &= c_h^{-1} \int k(u)^2 du \sigma_j^2(x_j) E_j(x_j)^{-1}, \\ \sigma_j^2(x_j) &= \alpha^* E_j(x_j)^{-1} + \sum_{l \neq j} \int \alpha_l(u) E_{jl}(x_j, u) E_j(x_j)^{-1} du + \alpha_j(x_j). \end{aligned}$$

This result yields in particular

$$n^{2/5} \{ \tilde{\alpha}(x) - \alpha(x) \} \rightarrow \mathcal{N} \left(c_h^2 \sum_{j=0}^d b_j(x_j), \sum_{j=0}^d v_j(x_j) \right),$$

for $\tilde{\alpha}(x) = \tilde{\alpha}^* + \sum_{j=0}^d \tilde{\alpha}_j(x_j)$ with $\tilde{\alpha}^* = \sum_{i=1}^n \int dN_i(s) / \sum_{i=1}^n \int Y_i(s) ds$.

The proof of Theorem 2 is given in Appendix A.2.

Remark 3. Note that the convergence rate is the same as for a one-dimensional local linear hazard estimator, see e.g. Nielsen and Tanggaard [2001]. Furthermore, $\tilde{\alpha}_j(x_j)$ estimates $\alpha_j(x_j) - \int \alpha_j(x_j) \hat{V}^j(x_j) dx_j$ instead of $\alpha_j(x_j)$. The terms $\nu_{n,j}$ correct for this shift in the estimation of each component. The sum $\sum_{j=0}^d \nu_{n,j}$ vanishes as the additive adjustments cancel each other off.

The component $\tilde{\alpha}^*$ of the estimator $\tilde{\alpha}$, which estimates α^* , is identical to $\bar{\alpha}^*$ from the local constant case. Its asymptotic behavior is explained in Remark 2.

5. Simulation Study

5.1. Simulation Setting

We assume that the survival times T_i follows a Gompertz–Makeham distribution, with hazard function is given by

$$\alpha(t, Z_i) = \alpha_0(t) + \sum_{k=1}^d \alpha_k(Z_{ik}) = e^{0.01t} + \frac{4}{\sqrt{d}} \sum_{k=1}^d (-1)^{k+1} \sin(\pi Z_{ik}),$$

($i = 1, \dots, n$). We add right censoring with censoring variables C_i that follows the same distribution as T_i , except the scale parameter being divided by 1.75. The factor $4d^{-1/2}$ is chosen so that the distribution of T_i doesn't much vary in the number of covariates d . Note that for convenience the components are differently identified then in (25), (15). We now describe how the covariates (Z_{i1}, \dots, Z_{id}) are generated. We first simulate $(\tilde{Z}_{i1}, \dots, \tilde{Z}_{id})$ from a d -dimensional multi-normal distribution with mean equal 0 and $\text{Corr}(Z_{ij}, Z_{il}) = \rho$ if $j \neq l$, else 1. Afterwards we set

$$Z_{ik} = 2.5\pi^{-1} \arctan(\tilde{Z}_{ik}).$$

We repeat the procedure and take the first $i = 1, \dots, n$ observations such that $4d^{-1/2} \sum_{k=1}^d (-1)^{k+1} \sin(\pi Z_{ik})$ is positive. Technically, the values of the covariates are conditioned such that the resulting hazard is positive, and hence well defined.

As kernel function k , we used the Epanechnikov kernel. Performance is measured via the integrated squared error evaluated ,

$$MISE_k = n^{-1} \sum_i (\eta_k(Z_{ik}) - \hat{\eta}_k(Z_{ik}))^2.$$

5.2. Simulation Results

We compare the performance of the local linear smooth backfitting estimator to the local constant smooth backfitting estimator. We also compare those proposed estimators to a version a classical backfitting equivalent where only the updated component is smoothed, see Buja et al. [1989b].

Figure 1 shows the estimation results for the first component from 100 simulations in a setting with sample size $n = 5000$, dimension $d = 3$, and correlation $\rho = 0.5$, calculated with a MISE optimal bandwidth. We find that the classical backfitting estimators are much more noisy than their smooth backfitting counterpart. The local constant smooth backfitting estimator is more wiggly than the local linear version. This first impression can be further verified in Table 1: classical backfitting estimators perform much worse than the smooth alternatives. The local linear classical backfitting estimator only gives sensible results in the most easy settings, that is when $n = 5000$ and or $d = 3$, while breaking down in all other cases. Another observation is that that the local linear smooth backfitting estimator is nearly always to be preferred to local constant smooth backfitting estimator. Only in the most challenging setting, i.e., $n = 500$, $d = 30$, did the local constant smooth backfitting estimator outperform the local linear version. But even in that case the advantage is only by a small margin.

d=3						
	n=500			n=5000		
	MISE	Bias ²	Variance	MISE	Bias ²	Variance
LL-SBF	0.25	0.07	0.17	0.031	0.007	0.024
LC-SBF	0.30	0.05	0.25	0.051	0.011	0.041
LL-BF	43.14	0.69	42.46	0.779	0.041	0.737
LC-BF	1.44	0.48	0.96	0.077	0.020	0.058
d=10						
	n=500			n=5000		
	MISE	Bias ²	Variance	MISE	Bias ²	Variance
LL-SBF	0.22	0.05	0.17	0.020	0.005	0.015
LC-SBF	0.24	0.08	0.17	0.030	0.006	0.025
LL-BF	1118.80	10.88	1107.91	0.135	0.057	0.078
LC-BF	1.02	0.03	0.99	0.031	0.005	0.026
d=30						
	n=500			n=5000		
	MISE	Bias ²	Variance	MISE	Bias ²	Variance
LL-SBF	0.18	0.03	0.15	0.014	0.0007	0.0133
LC-SBF	0.16	0.05	0.10	0.029	0.0172	0.0114
LL-BF	NA	NA	NA	0.171	0.1494	0.0217
LC-BF	NA	NA	NA	0.033	0.0227	0.0105

Table 1: Simulation results comparing four different estimators: local constant smooth backfitting, local linear smooth backfitting, local constant backfitting, local linear backfitting. Values are calculated from 500 Monte Carlo simulations with MISE optimal bandwidth.

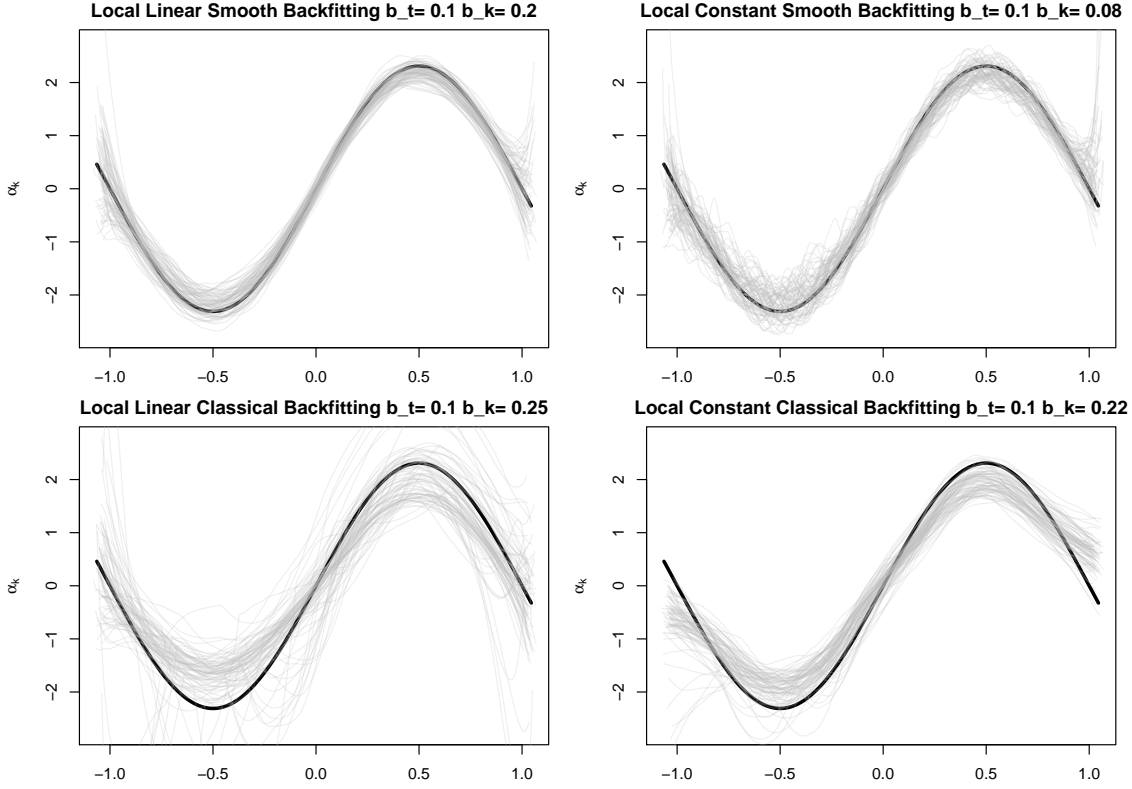


Figure 1: Simulation results comparing four different estimators: local constant smooth backfitting, local linear smooth backfitting, local constant backfitting, local linear backfitting. The grey lines are represent 100 Monte Carlo simulations with MISE optimal bandwidth estimating the true curve (black).

6. Data application: The TRACE study

The TRACE study group (see e.g. Jensen et al. [1997]) has collected information on more than 4000 consecutive patients with acute myocardial infarction (AMI) with the aim of studying the prognostic importance of ventricular fibrillation (vf) on mortality. We here consider a subset of these patients that are available in the `timereg` R package. We furthermore only consider those patients with more than 40 years of age, and only consider the first five years of follow-up time after the diagnosis. This results in $n = 1799$ observations. At entry, i.e., time of AMI occurrence, the patients had various risk factors recorded. Here, additionally to duration, i.e. time since AMI occurrence, we will consider age at AMI occurrence of the patient, a_i , and wall motion index (heart pumping effect based on ultrasound measurements where 2 is normal and 0 is worst [Scheike, 2009]), wmi_i . We will ignore additional binary covariates that have been recorded as our framework only covers continuous covariates. With that regard, this section should be seen as a simple illustration of our theoretical work rather than a serious attempt to answer a real-world question. In summary, we consider the model

$$\lambda_i(t) = Y_i(t)\{\alpha_0(t) + \alpha_2(a_i) + \alpha_3(wmi_i)\},$$

under the identifiability condition $\int \alpha_j(x_j)dx_j = 0$ for $j = 1, 2$. The initially estimated curve for α_0 can be seen in Figure 2. We find that the duration effect has two distinct periods with an increased risk in the beginning that flattens after approximately three months. This suggests that it might be beneficial to apply two different amounts of smoothing on those two periods. We therefore generate two different data sets from our original data set: The first data set covers the risk in the first three months (this can be achieved by censoring all patients who survived beyond three months) and the second data set covers

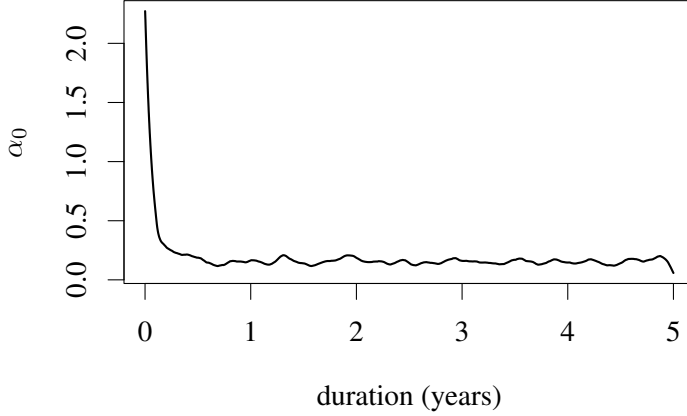


Figure 2: Local linear additive smooth backfitting fit of α_0 on the full data.

the risk conditional on surviving the first three months (i.e. one throws away all data where failure or censoring happened in the first three months). The results with our local linear estimator for the two different cohorts, i.e., those with ventricular fibrillation ($vf=1$) and those without ventricular fibrillation ($vf=0$) can be depicted in Figures 3 and 4.

The smoothing parameter was chosen manually: For the cohort with $vf = 0$ we have $n = 1655$ patients when considering the first three months and chose the bandwidths for (t, a, wmi) as $(0.1, 15, 0.8)$; for the dataset after surviving the first three months we have $n = 1482$ and chose a bandwidth of $(1, 15, 0.8)$. For the cohort with $vf = 1$ we have $n = 132$ patients for the first three months and chose a bandwidth of (t, a, wmi) as $(0.1, 20, 0.8)$; for the dataset after surviving the first three months we have $n = 75$ and chose a bandwidth of (t, a, wmi) as $(1, 20, 0.8)$. The dashed lines show a point-wise asymptotic 95% confidence interval based on Theorem 2. Note that it is hereby in particular assumed that (a) the bias can be neglected and (b) that the true underlying model is indeed additive. Therefore, the confidence intervals should be seen as rather optimistic. They nevertheless give an impression of the uncertainty under optimal conditions. Looking at Figures 3, we find that in the first three months $vf = 1$ leads to a significant increase in mortality risk. We also find that the risk increase is more severe for older people. Figure 4 does not provide evidence that $vf = 1$ leads to an increased risk after surviving the first three months. In the next section, we want to look at how confidence we can be with the model results.

6.1. Model robustness

6.1.1. CRPS score

We transform our estimated hazard function $\alpha = \alpha_0 + \alpha_1 + \alpha_2$ into a plug-in estimator of the survival function via the relationship $S(t) = \prod_{s \leq t} (1 - \alpha(s)ds)$. We then split our data randomly into an 80% training set and 20% test set. We train our model on the training set and evaluate the CPRS score [Avati et al., 2020] on the test set:

$$CRPS = n^{-1} \sum_{i=1}^n \int_0^{T_i} (1 - \hat{S}(s))^2 ds + \delta_i \int_{T_i}^{\infty} \hat{S}(s)^2 ds.$$

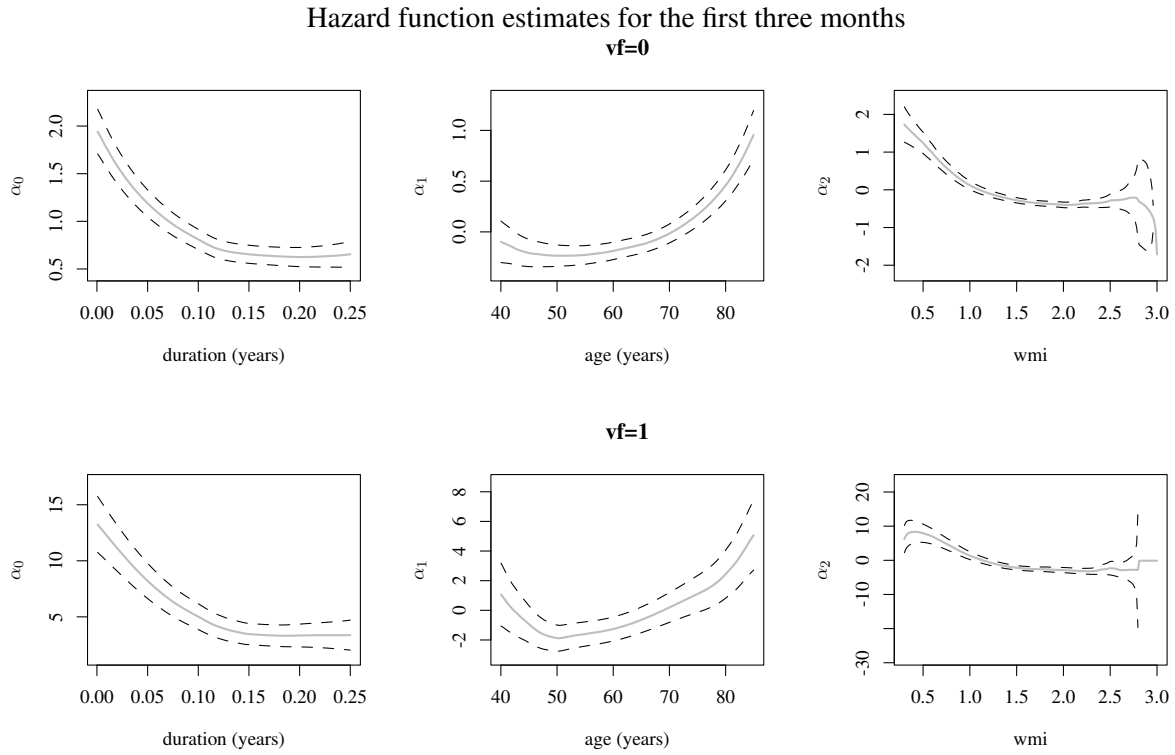


Figure 3: Local linear fit of $(\alpha_0, \alpha_1, \alpha_2)$ for the first three months for two different strata depending on the value of vf . Dashed line indicates asymptotic point-wise confidence interval.

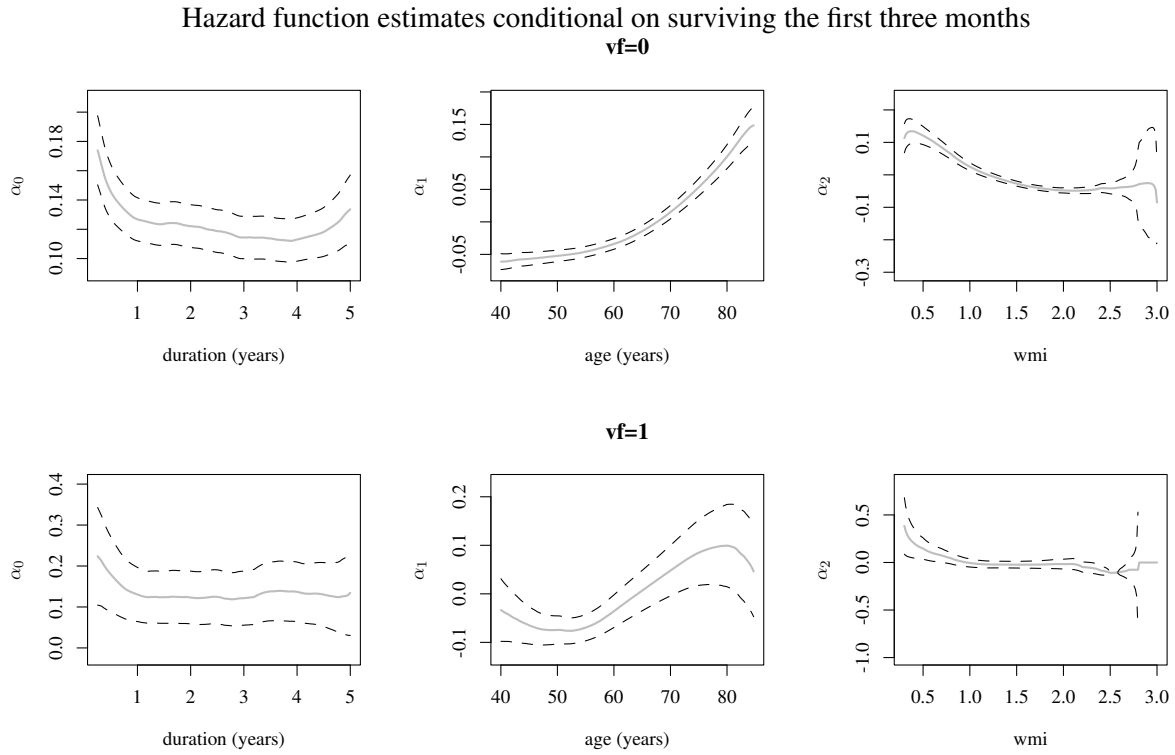


Figure 4: Local linear fit of $(\alpha_1, \alpha_2, \alpha_3)$ conditional on surviving the first three months for two different strata depending on the value of vf . Dashed line indicates asymptotic point-wise confidence interval.

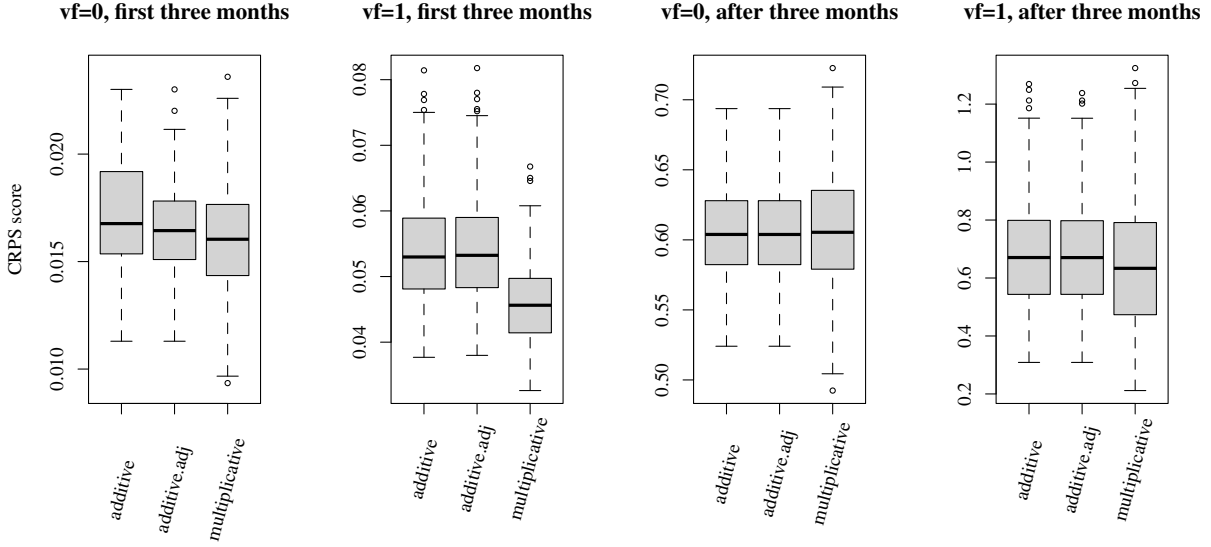


Figure 5: CPRS scores from 200 simulations of a 80/20 training-test-split. Boxplots are given for the four different data sets as described on top of the plots and each time for three different models: smooth backfitting additive model, smooth backfitting additive model using the adjusted survival estimates $\hat{S}^{\text{adj}}(s)$ and the smooth backfitting multiplicative model from Hiabu et al. [2021b].

Due to the additive structure, our survival prediction, while consistent can still be negative. We therefore consider a simple adjustment where we numerically calculate

$$\hat{S}^{\text{adj}}(s) = \prod_{s \leq t} (1 - \hat{\alpha}^{\text{adj}}(s) ds), \quad \hat{\alpha}^{\text{adj}}(s) = \max(\hat{\alpha}(s), 0).$$

Lastly, we compare our local linear additive fit with the local constant multiplicative smooth backfitting estimator from Hiabu et al. [2021a]. The results from 200 simulation runs can be seen in Figure 5. We have two main observations. Firstly, the model choice does not seem to have a big impact when considering survival conditional on surviving the first three months. Secondly, for survival during the first three months using the adjusted survival probability estimates improves the performance but even better performance can be achieved by using a multiplicative model. Nevertheless, remember that our smooth backfitting additive estimators has the desirable projection property that if the additive model assumption is violated the estimators converge to the closest additive fit; making the results therefore still interpretable. We now investigate this property in the next subsection.

6.1.2. Stability under model misspecification

We take the estimated multiplicative smooth backfitting model from the previous subsection, see also Figures 8 and 9 in the Appendix, as true model and investigate how in this case our additive estimator would look like. When generating the four data sets ($\text{vf} = 0, 1$; risk in the first three months, risk conditional on surviving the first three months), we keep the same number of samples as in the original data sets while sampling (a, wmi) with replacement from the original data sets. Afterwards, for each row, we draw a survival time from the multiplicative smooth backfitting model. The survival time is considered censored if it is greater than 0.25 when considering the first three months, and it is considered censored if it is greater than 5 when considering the period after the first three months.

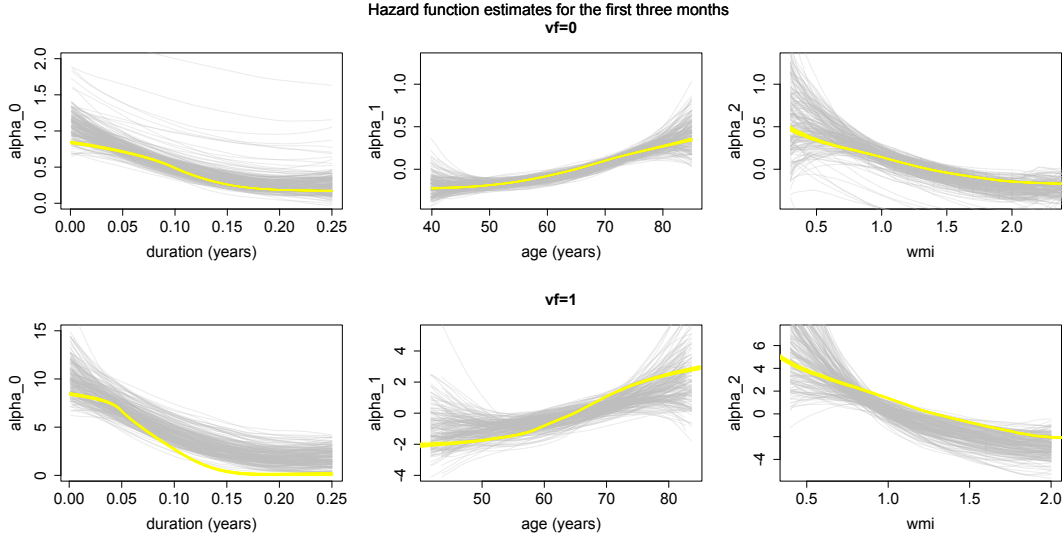


Figure 6: 200 simulations from a multiplicative hazard model, see Figure 8. Grey curves are fitted local linear smooth backfitting estimators. Yellow curves are approximately optimal additive fit derived from a smooth backfitting additive regression fit with the true hazard as response and an inflated sample size of 10,000.

We compare our additive smooth backfitting estimator to a somewhat optimal fit. Note that it is not clear how to analytically derive an optimal fit analytically or even numerically as it depends on the joint distribution of duration, age and wmi; which is not known. Therefore, we approximate the optimal fit by estimating an additive smooth backfitting regression function [Mammen et al., 1999a, Hiabu et al., 2023] based on 10,000 observations where the response is the known hazard. We consider 200 simulations and the fact that the regression estimator does not vary much gives us confidence that it is a good approximation of the optimal additive fit. The results are given in Figures 6 and 7. We find that our proposed estimators (grey lines) despite the limited sample sizes are reasonably close to the regression fit such that we can conclude that our approach is working reasonably well in estimating the optimal additive fit. Lastly, it should be noted that we also tried a classical backfitting approach with kernel smoothers with the result that the estimators for all components diverged in every simulation run and did not provide any result.

A. Appendix

A.1. Asymptotic theory for the local constant estimator

For the proof of Theorem 1, we apply the general theory for smooth backfitting estimators. We split the estimator into a stochastic part and a part consisting of its bias plus a function that vanishes. For counting processes martingales, these two parts are usually referred to as the variable and the stable part, respectively. One has to show three things: the convergence of the backfitting algorithm, asymptotic normality of the stochastic part and that the bias part vanishes asymptotically. In Mammen et al. [1999b], conditions for these three properties have been stated for a nonparametric regression setup. The main part of our proof is to verify these conditions under Assumptions A1–A5. For completeness we restate the modified conditions in our notation.

We also state propositions from Mammen et al. [1999b], adapted to our notation, which imply the properties we need if the following assumptions hold. The difference to Mammen et al. [1999b] is that we make use of martingale properties and counting process theory instead of the usual arguments for kernel density estimators.

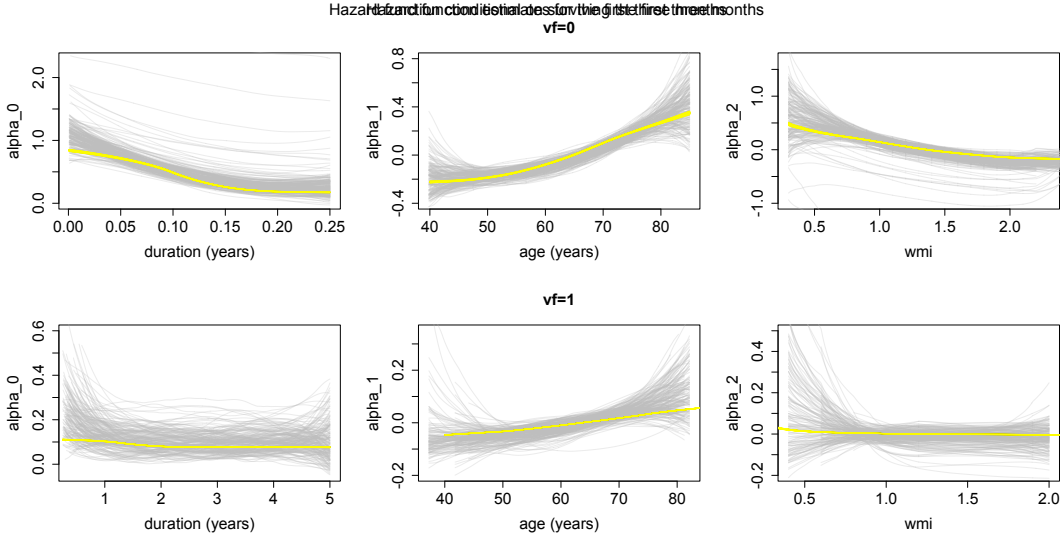


Figure 7: 200 simulations from a multiplicative hazard model, see Figure 9. Grey curves are fitted local linear additive smooth backfitting estimates. Yellow curves are approximately optimal additive fits derived from a smooth backfitting additive regression estimator with the true hazard as response and an inflated sample size of 10,000.

We start with assumptions about the marginal exposures and convergence of marginal exposure estimators. Note that we don't assume any particular definition of \hat{E}_j and $\hat{E}_{j,k}$, $j, k = 0, \dots, d$, for the following propositions.

B1 For all $j \neq k$ it holds

$$\int \frac{E_{j,k}(x_j, x_k)^2}{E_j(x_j)E_k(x_k)} dx_j dx_k < \infty.$$

B2 It holds

$$\begin{aligned} & \int \left[\frac{\hat{E}_j(x_j) - E_j(x_j)}{E_j(x_j)} \right]^2 E_j(x_j) dx_j = o_P(1), \\ & \int \left[\frac{\hat{E}_{j,k}(x_j, x_k)}{E_j(x_j)E_k(x_k)} - \frac{E_{j,k}(x_j, x_k)}{E_j(x_j)E_k(x_k)} \right]^2 E_j(x_j)E_k(x_k) dx_j dx_k = o_P(1), \\ & \int \left[\frac{\hat{E}_{j,k}(x_j, x_k)}{\hat{E}_j(x_j)E_k(x_k)} - \frac{E_{j,k}(x_j, x_k)}{E_j(x_j)E_k(x_k)} \right]^2 E_j(x_j)E_k(x_k) dx_j dx_k = o_P(1). \end{aligned}$$

Moreover, \hat{E}_j vanishes outside the support of E_j , $\hat{E}_{j,k}$ vanishes outside the support of $E_{j,k}$ and \hat{E} is symmetric, i.e. $\hat{E}_{j,k}(x_j, x_k) = \hat{E}_{k,j}(x_k, x_j)$.

We assume that the marginal pilot estimator and proportions of the marginal exposure estimators are somehow bounded in probability:

B3 There exists a constant C such that with probability tending to 1 for all j ,

$$\int \hat{\alpha}_j(x_j)^2 E_j(x_j) dx_j \leq C.$$

B4 For some finite intervals $S_j \subset \mathbb{R}$ that are contained in the support of E_j , $j = 1, \dots, d$, we suppose that there exists a finite constant C such that with probability tending to 1 for all $j \neq k$,

$$\sup_{x_j \in S_j} \int \frac{\hat{E}_{j,k}(x_j, x_k)}{E_j(x_j) \hat{E}_k(x_k)^2} dx_k \leq C.$$

We now introduce the notation $\hat{\alpha}_j = \hat{\alpha}_j^A + \hat{\alpha}_j^B$ for the one-dimensional smoother with

$$\hat{\alpha}_j^A = \hat{E}_j(x_j)^{-1} \frac{1}{n} \sum_{i=1}^n \int k_h(x_j, X_{ij}(s)) dM_i(s),$$

the variable part and

$$\hat{\alpha}_j^B = \hat{E}_j(x_j)^{-1} \frac{1}{n} \sum_{i=1}^n \int k_h(x_j, X_{ij}(s)) d\Lambda_i(s),$$

the stable part of $\hat{\alpha}_j$. Here, the compensator Λ_i of N_i is defined such that M_i is a martingale and $N_i = M_i + \Lambda_i$. The definition of M_i will be given later. Now we define the stochastic and stable components of the local constant smooth backfitting estimator, $\bar{\alpha}_{0,j}^s$, $\bar{\alpha}_j^s$, for $s \in \{A, B\}$, as the solution of

$$\bar{\alpha}_k^s(x_k) = \hat{\alpha}_k^s(x_k) - \hat{\alpha}_{0,k}^s - \sum_{j \neq k} \int_{\mathcal{X}_j} \bar{\alpha}_j^s(x_j) \left[\frac{\hat{E}_{j,k}(x_j, x_k)}{\hat{E}_k(x_k)} - \hat{E}_{j,[k+]}(x_j) \right] dx_j, \quad (35)$$

where $\hat{\alpha}_{0,k}^s = \int \hat{\alpha}_k^s(x_k) \hat{E}_k(x_k) dx_k / \int \hat{E}_k(x_k) dx_k$. Existence and uniqueness of $\hat{\alpha}_k^A, \hat{\alpha}_k^B$ is stated in Proposition 1 under the following assumptions. Assumption B6 assures converges of the variable part whereas B7 will be used for the structure of the bias part.

B5 There exists a constant C such that with probability tending to 1 for all j , it holds

$$\begin{aligned} \int \hat{\alpha}_j^A(x_j)^2 E_j(x_j) dx_j &\leq C, \\ \int \hat{\alpha}_j^B(x_j)^2 E_j(x_j) dx_j &\leq C. \end{aligned}$$

B6 We assume that there is a sequence $\Delta_n \rightarrow 0$ such that

$$\begin{aligned} \sup_{x_k \in S_k} \left| \int \frac{\hat{E}_{j,k}(x_j, x_k)}{\hat{E}_k(x_k)} \hat{\alpha}_j^A(x_j) dx_j \right| &= o_P(\Delta_n), \\ \left\| \int \frac{\hat{E}_{j,k}(x_j, x_k)}{\hat{E}_k(x_k)} \hat{\alpha}_j^A(x_j) dx_j \right\|_{2,k} &= o_P(\Delta_n), \end{aligned}$$

where $\|\cdot\|_{2,k}$ denotes norm defined via $\|g\|_{2,k} = \int g(u)^2 E_k(u) du$. The sets S_k have been introduced in Assumption B4.

B7 There exist deterministic functions $\mu_{n,j}$ such that

$$\sup_{x_j \in S_j} |\bar{\alpha}_j^B(x_j) - \mu_{n,j}(x_j)| = o_p(\Delta_n),$$

where S_k has been introduced in Assumption B4.

The following two propositions are results from Mammen et al. [1999b], adapted to our setting and notation. Under Assumptions B1–B3 and B5, Proposition 1 ensures that the backfitting algorithm converges and Propositions 2 and 3 give the asymptotic behavior of the backfitting estimator under Assumptions B1–B9.

Proposition 1 (Convergence of backfitting). *Under Assumptions B1–B3, with probability tending to 1, there exists a unique solution $\{\bar{\alpha}_j : j = 0, \dots, d\}$ to (19). Moreover, there exist constants $0 < \gamma < 1$ and $c > 0$ such that, with probability tending to 1, it holds:*

$$\int \left[\bar{\alpha}_j^{[r]}(x_j) - \bar{\alpha}_j(x_j) \right]^2 E_j(x_j) dx_j \leq c\gamma^{2r} \left(1 + \sum_{l=0}^d \int \left[\bar{\alpha}_l^{[0]}(x_l) \right]^2 E_l(x_l) dx_l \right),$$

for $j = 0, \dots, d$. The functions $\bar{\alpha}_l^{[0]}$ are the starting values of the backfitting algorithm. For $r > 0$ the functions $\bar{\alpha}_0^{[r]}, \dots, \bar{\alpha}_d^{[r]}$ are defined by equation (20).

Moreover, under the additional Assumption B5, with probability tending to 1, there exists a solution $\{\bar{\alpha}_j^s : j = 0, \dots, d\}$ of (35) that is unique for $s = A, B$, respectively.

Proposition 2 (Asymptotic behavior of stochastic part). *Suppose that Assumptions B1–B6 hold for a sequence Δ_n and intervals S_j , $j = 0, \dots, d$. Then it holds that*

$$\sup_{x_j \in S_j} |\bar{\alpha}_j^A(x_j) - [\hat{\alpha}_j^A(x_j) - \bar{\alpha}_{0,j}^A]| = o_P(\Delta_n).$$

Under the additional Assumption B7 it holds

$$\sup_{x_j \in S_j} |\bar{\alpha}_j^A(x_j) - [\hat{\alpha}_j^A(x_j) - \bar{\alpha}_{0,j}^A + \mu_{n,j}(x_j)]| = o_P(\Delta_n).$$

For the convergence of the bias term, we need the following.

B8 For all $j \neq k$, it holds

$$\sup_{x_j \in S_j} \int \left| \frac{\hat{E}_{j,k}(x_j, x_k)}{\hat{E}_j(x_j) \hat{E}_k(x_k)} - \frac{E_{j,k}(x_j, x_k)}{E_j(x_j) E_k(x_k)} \right| E_k(x_k) dx_k = o_p(1).$$

At last, Assumption B9 is about the structure of the bias term of the estimators.

B9 There exist deterministic functions $a_{n,0}(x_0), \dots, a_{n,d}(x_d)$ and constants $a_n^*, \gamma_{n,0}, \dots, \gamma_{n,d}$ and a function $\beta : \mathbb{R} \rightarrow \mathbb{R}$ (not depending on n), such that

$$\begin{aligned} \int a_{n,j}(x_j)^2 E_j(x_j) dx_j &< \infty, \\ \int \beta(x)^2 E(x) dx &< \infty, \\ \sup_{x_1 \in S_1, \dots, x_d \in S_d} |\beta(x)| &< \infty, \\ \gamma_{n,j} - \int a_{n,j}(x_j) \hat{E}_j(x_j) dx_j &= o_P(\Delta_n), \\ \sup_{x_j \in S_j} |\hat{\alpha}_j^B(x_j) - \hat{\mu}_{n,0} - \hat{\mu}_{n,j}(x_j)| &= o_P(\Delta_n), \\ \int |\hat{\alpha}_j^B(x_j) - \hat{\mu}_{n,0} - \hat{\mu}_{n,j}(x_j)|^2 E_j(x_j) dx_j &= o_P(\Delta_n^2), \end{aligned}$$

for random variables $\hat{\mu}_{n,0}$ and where

$$\hat{\mu}_{n,j}(x_j) = a_n^* + a_{n,j}(x_j) + \sum_{k \neq j} \int a_{n,k}(x_k) \frac{\hat{E}_{j,k}(x_j, x_k)}{\hat{E}_j(x_j)} dx_k + \Delta_n \int \beta(x) \frac{E(x)}{E_j(x_j)} dx_{-j}.$$

The following Proposition is taken from Mammen et al. [1999b] and we have adapted it to our notation. It implies in particular that the bias term of the smooth backfitting estimators equals the projections of the bias of the full-dimensional estimator of Linton et al. [2003].

Proposition 3 (Asymptotic behavior of bias part). *Under Assumptions B1–B6, B8, B9, for $j = 0, \dots, d$, it holds*

$$\sup_{x_j \in S_j} |\bar{\alpha}_j^B(x_j) - \mu_{n,j}(X_j)| = o_P(\Delta_n),$$

for $\mu_{n,j}(x_j) = a_{n,j}(x_j) - \gamma_{n,j} + \Delta_n \beta_j(x_j)$ with

$$(\beta_0, \beta_1, \dots, \beta_d) = \arg \min_{\mathcal{B}} \int [\beta(x) - \beta_0 - \beta_1(x_1) - \dots - \beta_d(x_d)]^2 E(x) dx,$$

and $\mathcal{B} = \{\tilde{\beta} = (\beta_0, \beta_1, \dots, \beta_d) : \int \beta_j(x_j) E_j(x_j) dx_j = 0; j = 0, \dots, d\}$. In particular, does Assumption B7 hold with this choice of $\mu_{n,j}(x_j)$.

With the next lemma we ensure that the constant α^* is estimated at parametric rate in the local constant setting. This standard result will also be needed in the proof of Theorem 1.

Lemma 1. *Let $\bar{\alpha}^* = (\sum_{i=1}^n \int dN_i(s)) / (\sum_{i=1}^n \int Y_i(s) ds)$ as defined in equation (14). Under the condition $\int \alpha_j(x_j) E_j(x_j) dx_j = 0$, for $j = 0, \dots, d$ together with Assumption A2, it holds*

$$n^{1/2}(\bar{\alpha}^* - \alpha^*) \rightarrow \mathcal{N}(0, \sigma_{\alpha^*}^2),$$

as $n \rightarrow \infty$ and for $\sigma_{\alpha^*}^2 = \alpha^*(1 - \alpha^*)$. This implies in particular $\bar{\alpha}^* - \alpha^* = O_p(n^{-1/2})$.

Proof. We first note that it holds $E_0(t) = \int E(x) dx_{-0} = \gamma(t)$ for $x = (t, z)$ and with γ from Assumption A2. Using $\frac{1}{n} \sum_{i=1}^n Y_i(s) = \gamma(s) + o_P(1)$ in the denominator and the usual martingale decomposition for counting processes in the numerator, we get

$$\begin{aligned} \mathbb{E}[n^{1/2} \bar{\alpha}^*] &= n^{1/2} \alpha^* + o(1), \\ \text{Var}(n^{1/2} \bar{\alpha}^*) &= \alpha^*(1 - \alpha^*) + o(1), \end{aligned}$$

because of the identification $\int \alpha_0(s) \gamma(s) ds = 0$. The terms $\mathbb{E}[\int \alpha_j(Z_{i,j}(s)) \gamma(s) ds]$ in the stable part of the martingale vanish because of $\gamma(t) = \int E(x) dx_{-0}$ and the identification criterion. The Central Limit Theorem for *i.i.d.* observations then yields the result. \square

Moreover, we will make use of the following counting process martingale central limit theorem, which is a direct application of Rebolledo's Theorem (Theorem II.5.1 in Andersen et al. [1993]). It is a multivariate extension of the central limit theorem for martingales in Ramlau-Hansen [1983].

Lemma 2 (Multivariate Ramlau-Hansen). *Let $\{M_i : i = 1, \dots, n\}$ be a sequence of *i.i.d.* martingales and let $g_{i,j}^{(n)}$ be predictable functions for $j = 1, \dots, d$. Furthermore, suppose it holds for $j, k = 1, \dots, d$,*

$$\sum_{i=1}^n \int g_{i,j}^{(n)}(s) g_{i,k}^{(n)}(s) d\langle M_i \rangle(s) \rightarrow \sigma_{j,k}^2, \quad (36)$$

$$\sum_{i=1}^n \int \left[g_{i,j}^{(n)}(s) \right]^2 I_{\{|g_{i,j}^{(n)}(s)| > \varepsilon\}} d\langle M_i \rangle(s) \rightarrow 0, \quad (37)$$

in probability for $n \rightarrow \infty$ with $\sigma_{j,k}^2 > 0$ and for every $\varepsilon > 0$. Then

$$\sum_{i=1}^n \begin{pmatrix} \int g_{i,1}^{(n)}(s) dM_i(s) \\ \vdots \\ \int g_{i,d}^{(n)}(s) dM_i(s) \end{pmatrix} \rightarrow \mathcal{N}(0, \Sigma),$$

in distribution for $n \rightarrow \infty$, where $\sigma_{j,k}^2$, $j, k = 1, \dots, d$ are the entries of the covariance matrix Σ .

To show Theorem 1 we apply Propositions 1–3 and Lemmas 1 and 2. According to the propositions it is sufficient to verify Assumptions B1–B9. In the proof of Theorem 1 we will show that our Assumptions A1–A5 imply Assumptions B1–B9 for the right choices of Δ_n , $a_{n,j}$, β , $\gamma_{n,j}$.

Proof of Theorem 1. In the following we show how Assumptions A1–A5 imply B1–B6, B8–B9 with our choice of marginal pilot estimators. Assumption B7 is established through Proposition 3 once the other assumptions are verified.

Without loss of generality, the proofs are done for $\mathcal{T} = R = 1$, i.e. for survival time and covariates with support $[0, 1]$ and we will show that Assumptions B1–B9 are satisfied on closed subsets $S_0 \subset (0, \mathcal{T})$ and $S_j \subset (0, R)$, $j = 1, \dots, d$.

We first note that Assumption B1 follows directly from A1.

For the remaining stochastic statements, we start with the derivation of convergence rates for the marginal exposure estimators. Moreover, we will show all statements for the rate $\Delta_n = h^2$. With $I_h = [2h, 1 - 2h]$, it holds for $j = 0, \dots, d$,

$$\sup_{x_j \in I_h} |\hat{E}_j(x_j) - E_j(x_j)| = O_P\left((\log n)^{1/2} n^{-2/5}\right), \quad (38)$$

$$\sup_{x_j, x_k \in I_h} |\hat{E}_{j,k}(x_j, x_k) - E_{j,k}(x_j, x_k)| = O_P\left((\log n)^{1/2} n^{-3/10}\right), \quad (39)$$

$$\sup_{0 \leq x_j \leq 1} |\hat{E}_j(x_j) - \int_0^1 k_h(x_j, u) du E_j(x_j)| = O_P\left(n^{-1/5}\right), \quad (40)$$

$$\sup_{0 \leq x_j, x_k \leq 1} |\hat{E}_{j,k}(x_j, x_k) - \int_0^1 k_h(x_j, u) du \int_0^1 k_h(x_k, v) dv E_{j,k}(x_j, x_k)| = O_P\left(n^{-1/5}\right). \quad (41)$$

Before proving equations (38)–(41) we emphasize that they imply in particular

$$\sup_{x_j \in [0,1]} |\hat{E}_j(x_j)| = O_P(1), \quad (42)$$

$$\sup_{x_j \in [0,1]} |\hat{E}_j(x_j)^{-1}| = O_P(1), \quad (43)$$

$$\sup_{x_j, x_k \in [0,1]} |\hat{E}_{j,k}(x_j, x_k)| = O_P(1). \quad (44)$$

Condition (38) follows with standard arguments (chaining, Bernstein inequality, c.f. Mammen et al. [1999b] for the regression case) from

$$\mathbb{E}[\hat{E}_j(x_j)] - E_j(x_j) = O\left(n^{-2/5}\right), \quad (45)$$

$$|\hat{E}_j(x_j)| \leq C_1 \quad \text{a.s.}, \quad (46)$$

$$|\hat{E}_j(u_1) - \hat{E}_j(u_2)| \leq C_2 |u_1 - u_2| n^m O_P(1), \quad (47)$$

$$\text{Var}(\hat{E}_j(x_j)) = O(n^{-4/5}), \quad (48)$$

for constants $0 < C_1, C_2 < \infty$, $m > 0$ and all $u_1 \neq u_2, x_j \in [0, 1]$. This can be seen with Taylor expansions and using the Lipschitz continuity of K . Condition (39)–(41) can be shown in the same

way. For (40) and (41) note that $\int_0^1 k_h(x_j, u)du$ corrects the kernel at the boundaries where it does not integrate to unity.

We now show (45)–(48). Condition (46) follows directly from A3 with K being bounded and the covariates having compact support. With usual kernel estimator arguments and a Taylor expansion of f_s around x_j we get

$$\mathbb{E}[\hat{E}_j(x_j)] - E_j(x_j) = o(h^2), \quad (49)$$

which implies condition (45) immediately. Condition (48) can be derived analogously. Eventually, the Lipschitz continuity of K in A3 yields (47).

Since the kernel k is cut off outside $[0, 1]$, Assumption B2 follows directly from (42)–(44).

For the remaining assumptions we split the marginal estimator $\hat{\alpha}_j(x_j)$ as described for B5 into the variable part

$$\hat{\alpha}_j^A(x_j) = \frac{\frac{1}{n} \sum_{i=1}^n \int k_h(x_j, X_{ij}(s)) dM_i(s)}{\hat{E}_j(x_j)},$$

and the stable part

$$\hat{\alpha}_j^B(x_j) = \frac{\frac{1}{n} \sum_{i=1}^n \int k_h(x_j, X_{ij}(s)) d\Lambda_i(s)}{\hat{E}_j(x_j)},$$

via $\hat{\alpha}_j(x_j) = \hat{\alpha}_j^A(x_j) + \hat{\alpha}_j^B(x_j)$. With the choice $\Lambda_i(t) = \int_0^t \lambda_i(s)ds$ for the intensity λ_i that was introduced in equation (4), we get that $M_i = N_i - \Lambda_i$ defines a unique square integrable martingale arising from the counting process N_i .

Next we derive the asymptotic behavior of $\hat{\alpha}_j^A(x_j)$ and $\hat{\alpha}_j^B(x_j)$ separately. With M_i being a martingale and $k_h(x_j, X_{ij}(s))$ being predictable, the integral $\int k_h(x_j, X_{ij}(s)) dM_i(s)$ is a martingale as well. Using the multivariate Ramlau-Hansen martingale central limit theorem in Lemma 2, we will show that $\hat{\alpha}_j^A(x_j)$ is asymptotically normally distributed whereas the difference between the stable part $\hat{\alpha}_j^B(x_j)$ and $\alpha_j(x_j)$ asymptotically behaves like the bias term $b_j(x_j)$.

For $x_j \in I_h$, we now show conditions (36) and (37) of Lemma 2 for $g_{ij}^{(n)}(s) = n^{-3/5} k_h(x_j - X_{ij}(s))$. Note that with Λ_i being the compensator of M_i , we get in particular $d\langle M_i \rangle(s) = d\Lambda_i(s) = \left[\alpha^* + \sum_{k=0}^d \alpha_k(X_{ik}(s)) \right] Y_i(s) ds$.

For cross-terms with $j \neq l$ in (36), it holds with this choice of $g_{ij}^{(n)}$ that

$$\begin{aligned} & \mathbb{E} \left[\sum_{i=1}^n \int g_{ij}^{(n)}(s) g_{il}^{(n)}(s) d\langle M_i \rangle(s) \right] \\ &= \mathbb{E} \left[\left(\frac{1}{n} n^{2/5} \right)^2 \sum_{i=1}^n \int k_h(x_j - X_{ij}(s)) k_h(x_l - X_{il}(s)) d\Lambda_i(s) \right] \\ &= n^{-1/5} \int \int k_h(x_j - u_j) k_h(x_l - u_l) \left[\alpha^* + \alpha_0(s) + \sum_{k=1}^d \alpha_k(u_k) \right] \\ & \quad \times \gamma(s) f_s(u_1, \dots, u_d) d(u_1, \dots, u_d) ds \\ &= O(h), \end{aligned} \quad (50)$$

because of the bounded support of the covariates and with the hazard rates being continuous. We write $f_s(u_1, \dots, u_d)$ for the conditional density of $(X_{i1}(s), \dots, X_{id}(s))$ at (u_1, \dots, u_d) given $Y_i(s) = 1$. Moreover, it can be shown easily with similar arguments that the variance of these terms satisfies

$$\text{Var} \left(\sum_{i=1}^n \int g_{ij}^{(n)}(s) g_{il}^{(n)}(s) d\langle M_i \rangle(s) \right) = O(h^6), \quad (51)$$

and hence $\sigma_{k,l}^2 = 0$ for $j \neq l$ is assured for (36). For the diagonal of the asymptotic covariance matrix $\tilde{\Sigma}$, we start with the following preliminary results. For $x_j \in I_h$ it holds

$$\begin{aligned}
& n^{4/5} \mathbb{E} \left[n^{-2} \sum_{i=1}^n \int k_h(x_j - X_{ij}(s))^2 \alpha_j(X_{ij}(s)) Y_i(s) ds \right] \\
&= n^{4/5} n^{-1} \int \int k_h(x_j - u)^2 \alpha_j(u) f_s(u) \gamma(s) du ds \\
&= n^{-1/5} h^{-1} \int \int k(v)^2 \alpha_j(x_j + vh) f_s(x_j + vh) \gamma(s) dv ds \\
&= (nh^5)^{-1/5} \int k(v)^2 \alpha_j(x_j) dv E_j(x_j) + o(1) \\
&= c_h^{-1} \int k(v)^2 dv \alpha_j(x_j) E_j(x_j) + o(1),
\end{aligned} \tag{52}$$

with usual kernel estimator arguments. Analogously, we get for $l \neq j$, that

$$\begin{aligned}
& n^{4/5} \mathbb{E} \left[n^{-2} \sum_{i=1}^n \int k_h(x_j - X_{il}(s))^2 \alpha_l(X_{il}(s)) Y_i(s) ds \right] \\
&= c_h^{-1} \int k(v)^2 dv \int \int \alpha_k(u_l) f_s(x_j, u_l) \gamma(s) du_l ds + o(1).
\end{aligned} \tag{53}$$

For the variance of the diagonal terms, one can derive

$$\text{Var} \left(\sum_{i=1}^n \int \left(g_{i,j}^{(n)}(s) \right)^2 d\langle M_i \rangle(s) \right) = O(h^5), \tag{54}$$

which yields the stochastic convergence of diagonal variance terms together with (52) and (53).

Summarizing, equations (50)–(54) imply condition (36) of Lemma 2 with $\sigma_{j,j}^2 = \tilde{\sigma}_j^2(x_j)$ for

$$\tilde{\sigma}_j^2(x_j) = c_h^{-1} \int k^2(v) dv \left(\alpha^* + \sum_{l \neq j} \int \int \alpha_k(u_l) f_s(x_j, u_l) \gamma(s) du_l ds + \alpha_j(x_j) E_j(x_j) \right),$$

and $\sigma_{j,k}^2 = 0$, $j \neq k$.

The Lindeberg condition (37) is satisfied under Assumption A3 since we assume bounded support for all covariates.

Hence, Lemma 2 implies

$$n^{2/5} \begin{pmatrix} \hat{\alpha}_0^A(x_0) \hat{E}_0(x_0) \\ \vdots \\ \hat{\alpha}_d^A(x_d) \hat{E}_d(x_d) \end{pmatrix} \rightarrow \mathcal{N}(0, \tilde{\Sigma}), \tag{55}$$

where $\tilde{\Sigma}$ is a diagonal matrix with the entries $\tilde{\sigma}_j^2(x_j)$, $j = 0, \dots, d$.

Equations (48) and (49) imply convergence in probability of $\hat{E}_j(x_j)$ to $E_j(x_j)$ at a fast enough rate and hence, we get

$$n^{2/5} \begin{pmatrix} \hat{\alpha}_0^A(x_0) \\ \vdots \\ \hat{\alpha}_d^A(x_d) \end{pmatrix} \rightarrow \mathcal{N}(0, \Sigma), \tag{56}$$

from (55) with Σ being a diagonal matrix with the entries $\sigma_j^2(x_j) = \tilde{\sigma}_j^2(x_j) E_j(x_j)^{-2}$, $j = 0, \dots, d$.

Note that condition (56), implies in particular $\text{Var}(\hat{\alpha}_j^A(x_j)) = O(n^{-4/5})$. Following the line of argumentation we used to prove (38) for $\hat{E}_j(x_j)$, this leads to

$$\sup_{x_j \in I_h} |\hat{\alpha}_j^A(x_j)| = O_P((\log n)^{1/2} n^{-2/5}). \quad (57)$$

Analogously, one can get a similar result at the boundary and thus

$$\sup_{x_j \in [0,1]} |\hat{\alpha}_j^A(x_j)| = O_P(1) \quad (58)$$

on the whole support.

For the stable part, we refer to Nielsen and Linton [1995] who have shown for

$$B_j(x_j) = \frac{1}{n} \sum_{i=1}^n \int k_h(x_j, X_{ij}(s)) d\Lambda_i(s)$$

that

$$\sup_{x_j \in [0,1]} |B_j(x_j) - \mathbb{E}[B_j(x_j)]| = o_P(1), \quad (59)$$

$$\sup_{x_j \in [0,1]} |\mathbb{E}[B_j(x_j)]| = o(1), \quad (60)$$

making use of the Lipschitz continuity of K from Assumption A3 and of Assumption A1. Together with (43), equations (59) and (60) imply

$$\sup_{x_j \in [0,1]} |\hat{\alpha}_j^B(x_j)| = O_P(1). \quad (61)$$

One can get Assumptions B3 and B5 immediately from (58) and (61). Assumptions B2, B4 and B8 follow from equations (38)–(41).

We illustrate the derivation of Assumption B6 for $x_j \in I_h$. First note that $\int E_{j,k}(x_j, x_k) (E_j(x_j))^{-1} k_h(x_j - X_{i,j}(s)) dx_j$ is a bounded function $g(h, x_k, X_{i,j}(s))$ of arguments h, x_k , and $X_{i,j}(s)$ and hence predictable. This leads to

$$\text{Var}\left(\int g(h, x_k, X_{i,j}(s)) dM_i(s)\right) = O(1),$$

due to M_i being a square integral martingale and a similar derivation to (50)–(54). Thus, it holds that

$$n^{1/2} \left(\frac{1}{n} \sum_{i=1}^n \int \int \frac{E_{j,k}(x_j, x_k)}{E_j(x_j)} k_h(x_j, X_{i,j}(s)) dx_j dM_i(s) \right)$$

is asymptotically normally distributed and in particular

$$\frac{1}{n} \sum_{i=1}^n \int \int \frac{E_{j,k}(x_j, x_k)}{E_j(x_j)} k_h(x_j, X_{i,j}(s)) dx_j dM_i(s) = O_P(n^{-1/2}).$$

Note that by integrating over x_k , we achieve the parametric rate $n^{1/2}$ making the usual rate h^{-1} vanish. Together with (38) and (39), the last equation yields

$$\int \frac{\hat{E}_{j,k}(x_j, x_k)}{\hat{E}_k(x_k)} \hat{\alpha}_j^A(x_j) dx_j$$

$$\begin{aligned}
&= \int \frac{E_{j,k}(x_j, x_k)}{E_k(x_k)} \hat{\alpha}_j^A(x_j) dx_j + O_P(n^{-3/10} n^{-2/5} \log n) \\
&= E_k(x_k)^{-1} \frac{1}{n} \sum_{i=1}^n \int \int \frac{E_{j,k}(x_j, x_k)}{E_j(x_j)} k_h(x_j, X_{i,j}(s)) dx_j dM_i(s) + O_P(n^{-3/10} n^{-2/5} \log n) \\
&= O_P(n^{-1/2}),
\end{aligned}$$

since (38) further implies $\hat{\alpha}_j^A(x_j) = E_j(x_j)^{-1} (x_j - X_{i,j}(s)) dM_i(s) + O_P(n^{-2/5} (\log n)^{1/2})$.

The last equation proves Assumption B6.

We prove Assumption B9 for the following choices for $j = 0, \dots, d$.

$$\begin{aligned}
a_n^* &= \alpha^*, \\
a_{n,j}(x_j) &= \alpha_j(x_j) + \alpha'_j(x_j) \int k_h(x_j, u)(u - x_j) \left[\int k_h(x_j, v) dv \right]^{-1} du, \\
\beta(x) &= \sum_{j=0}^d \left[\alpha'_j(x_j) \frac{\partial \log E(x)}{\partial x_j} + \frac{1}{2} \alpha''_j(x_j) \right] \int u^2 k(u) du, \\
\gamma_{n,j} &= 0.
\end{aligned}$$

The first three statement of B9 hold immediately with this choice of $a_{n,j}$ and Assumptions A1 and A3.

For the fourth statement it holds

$$\int a_{n,j}(x_j) \hat{E}_j(x_j) dx_j = \int \alpha_j(x_j) \hat{E}_j(x_j) dx_j + \int \alpha'_j(x_j) \hat{E}_j(x_j) \frac{\int k_h(x_j, u)(u - x_j)}{\int k_h(x_j, v) dv} dx_j, \quad (62)$$

and we investigate the two summands separately. For the first one it holds

$$\begin{aligned}
\int \alpha_j(x_j) \hat{E}_j(x_j) dx_j &= \frac{1}{n} \sum_{i=1}^n \int \int \alpha_j(x_j) k_h(x_j, X_{i,j}(s)) dx_j Y_i(s) ds \\
&= \frac{1}{n} \sum_{i=1}^n \int g_h(X_{i,j}(s)) Y_i(s) ds \\
&= \mathbb{E} \left[\int \alpha_j(x_j) \hat{E}_j(x_j) dx_j \right] + o_P(n^{-1/2}) \\
&= \int \int \int \alpha_j(x_j) k_h(x_j - u) \gamma(s) f_s(u) du ds dx_j + o_P(n^{-1/2}) \\
&= \int \int \alpha_j(x_j) k_h(x_j - u) E_j(u) du dx_j + o_P(n^{-1/2}) \\
&= \int \alpha_j(x_j) E_j(x_u) dx_j + o_P(n^{-1/2}),
\end{aligned}$$

since $\int g_h(X_{i,j}(s)) Y_i(s) ds$ are *i.i.d.* random variables with the definition $g_h(X_{i,j}(s)) = \int \alpha_j(x_j) k_h(x_j - X_{i,j}(s)) dx_j$ and the Central Limit Theorem applies as for B6. The last equation follows from a substitution, a Taylor expansion of E_j and the fact that k is a kernel of order one.

The second summand can be treated analogously yielding

$$\begin{aligned}
&\int \alpha'_j(x_j) \hat{E}_j(x_j) \frac{\int k_h(x_j, u)(u - x_j)}{\int k_h(x_j, v) dv} dx_j \\
&= \int \int \alpha'_j(x_j) k_h(x_j - u)(u - x_j) E_j(u) du dx_j + o_P(n^{-1/2}), \\
&= o_P(n^{-1/2}),
\end{aligned}$$

and hence in total

$$\int a_j(x_j) \hat{E}_j(x_j) dx_j = o_P(n^{-1/2}). \quad (63)$$

because of the identification $\int \alpha_j(x_j) E_j(x_u) dx_j = 0$. This verifies the fourth statement of B9 with $\gamma_{n,j} = 0$.

To prove B9, we start with two preliminary results:

$$\sup_{x_j \in I_h} |\hat{\alpha}_j^B(x_j) - \hat{\mu}_{n,j}(x_j)| = o_P(h^2), \quad (64)$$

$$\sup_{x_j \in I_h^c} |\hat{\alpha}_j^B(x_j) - \hat{\mu}_{n,j}(x_j)| = o_P(h). \quad (65)$$

Recall that by definition it holds

$$\begin{aligned} \hat{\alpha}_j^B(x_j) &= \frac{1}{n} \sum_{i=1}^n \int k_h(x_j - X_{ij}(s)) d\Lambda_i(s) \left(\hat{E}_j(x_j) \right)^{-1} \\ &= \frac{1}{n} \sum_{i=1}^n \int k_h(x_j - X_{ij}(s)) \left[\alpha^* + \sum_{l=0}^d \alpha_l(X_{il}(s)) \right] Y_i(s) ds \left(\hat{E}_j(x_j) \right)^{-1}, \end{aligned}$$

and

$$\begin{aligned} \hat{\mu}_{n,j}(x_j) &= a_{n,0} + a_{n,j}(x_j) + \sum_{k \neq j} \int a_{n,k}(x_k) \frac{\hat{E}_{j,k}(x_j, x_k)}{\hat{E}_j(x_j)} dx_k - \Delta_n \int \beta(x) \frac{E(x)}{E_j(x_j)} dx_{-j} \\ &= \alpha^* + \alpha_j(x_j) + \alpha'_j(x_j) \int k_h(x_j, u)(u - x_j) \left[\int k_h(x_j, v) dv \right]^{-1} du \\ &\quad + \sum_{k \neq j} \int \left(\alpha_k(x_k) + \alpha'_k(x_k) \int k_h(x_k, u)(u - x_k) \left[\int k_h(x_k, v) dv \right]^{-1} du \right) \\ &\quad \times \frac{\hat{E}_{j,k}(x_j, x_k)}{\hat{E}_j(x_j)} dx_k \\ &\quad + \Delta_n \int u^2 k(u) du \int \sum_{j=0}^d \left[\alpha'_j(x_j) \frac{\partial \log E(x)}{\partial x_j} + \frac{1}{2} \alpha''_j(x_j) \right] \frac{E(x)}{E_j(x_j)} dx_{-j}. \end{aligned}$$

Next, it holds for $j = 0, \dots, d$,

$$\begin{aligned} &\frac{1}{n} \sum_{i=1}^n \int k_h(x_j, X_{ij}(s)) \alpha_j(X_{ij}(s)) Y_i(s) ds \left(\hat{E}_j(x_j) \right)^{-1} \\ &= \alpha_j(x_j) + \alpha'_j(x_j) \int k_h(x_j, u)(u - x_j) du \left(\int k_h(x_j, u) du \right)^{-1} \\ &\quad + h^2 \int u^2 k(u) du \left[E'_j(x_j) \alpha'_j(x_j) + \frac{1}{2} E_j(x_j) \alpha''_j(x_j) \right] E_j(x_j)^{-1} + R_{n,j}(x_j), \end{aligned} \quad (66)$$

with $\sup_{x_j \in I_h} |R_{n,j}(x_j)| = o_p(h^2)$ and $\sup_{x_j \in [0,1] \setminus I_h} |R_{n,j}(x_j)| = O_p(h^2)$. Similarly, for $k \neq j$, we get

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n \int k_h(x_j, X_{ij}(s)) \alpha_k(X_{ik}(s)) Y_i(s) ds \left(\hat{E}_j(x_j) \right)^{-1} \\
&= \int \alpha_k(x_k) \frac{\hat{E}_{j,k}(x_j, x_k)}{\hat{E}_j(x_j)} dx_k \\
&+ \int \alpha'_k(x_k) \frac{\hat{E}_{j,k}(x_j, x_k)}{\hat{E}_j(x_j)} k_h(x_k, u) (u - x_k) du \left(\int k_h(x_j, u) du \right)^{-1} \\
&+ h^2 \int u^2 k(u) du \int \left[\frac{\partial E_{j,k}(x_j, x_k)}{\partial x_k} \alpha'_k(x_k) + \frac{1}{2} E_{j,k}(x_j, x_k) \alpha''_j(x_j) \right] E_j(x_j)^{-1} \\
&+ R_{n,j,j}(x_j),
\end{aligned} \tag{67}$$

with $\sup_{x_j \in I_h} |R_{n,j,k}(x_j)| = o_p(h^2)$ and $\sup_{x_j \in [0,1] \setminus I_h} |R_{n,j,k}(x_j)| = O_p(h^2)$. Equation (66) follows straightforward with a Taylor expansion of each α_j and E_j and for the derivation of (67) we refer to the proof of Theorem 4 in Mammen et al. [1999b], where the analogue is shown for the nonparametric regression case. Equations (66) and (67) imply (64) and (65) with above choices of $a_{n,j}$, β and $\gamma_{n,j}$. Eventually, together with (63), conditions (64) and (65) imply A9.

For the last statement of the theorem, we note that the constant component α^* in the conditional hazard can be estimated at a parametric rate $n^{-1/2}$ by $\bar{\alpha}^*$ due to Lemma 1. \square

A.2. Asymptotic theory for the local linear estimator

For the local linear estimator, we follow the same procedure as in Section A.1. We first introduce general assumptions as well as a set of results from Mammen et al. [1999b] which we will apply to prove Theorem 2. Then we verify the new assumptions under Assumptions A1–A5.

Let $E : \mathcal{X} \rightarrow [0, 1]$ be the exposure as defined earlier and let W be a (deterministic) positive definite $(d+1) \times (d+1)$ -matrix with elements $W_{r,s}$ such that $W_{0,0} = 1$. We set

$$M_j(x_j) = \begin{pmatrix} W_{0,0} & W_{j,0} \\ W_{j,0} & W_{j,j} \end{pmatrix} E_j(x_j), \tag{68}$$

$$S_{l,j}(x_l, x_j) = \begin{pmatrix} W_{0,0} & W_{l,0} \\ W_{j,0} & W_{l,j} \end{pmatrix} E_{l,j}(x_l, x_j). \tag{69}$$

These will later be the fixed but unknown matrices to which \hat{M}_j and \hat{S}_j , respectively, converge.

Now we make the following assumptions which are all of similar nature to B1–B9. Note that these are assumptions on $\hat{V}^j(x_j)$, $\hat{V}_j^j(x_j)$, $\hat{V}_j^j(x_j)$, $\hat{V}_{j,j}^j(x_j)$, $\hat{V}^{l,j}(x_l, x_j)$, $\hat{V}_l^{l,j}(x_l, x_j)$, $\hat{V}_j^{l,j}(x_l, x_j)$, $\hat{V}_{l,j}^{l,j}(x_l, x_j)$ and $\hat{\alpha}_j(x_j)$, $\hat{\alpha}^j(x_j)$, and all $x_j, x_l, j, l = 0, \dots, d$ and we don't assume any particular definition of these terms for the following propositions.

B1' For all $j \neq k$ it holds

$$\int \frac{E_{j,k}(x_j, x_k)^2}{E_j(x_j) E_k(x_k)} dx_j dx_k < \infty.$$

B2' For \hat{M}_j and $\hat{S}_{l,j}$ as in (27) and (28) it holds

$$\begin{aligned}
& \int \left[\frac{\hat{V}^j(x_j) - E_j(x_j)}{E_j(x_j)} \right]^2 E_j(x_j) dx_j = o_P(1), \\
& \int \left[\frac{\hat{V}^{j,k}(x_j, x_k)}{E_j(x_j) E_k(x_k)} - \frac{E_{j,k}(x_j, x_k)}{E_j(x_j) E_k(x_k)} \right]^2 E_j(x_j) E_k(x_k) dx_j dx_k = o_P(1),
\end{aligned}$$

$$\int \left[\hat{M}_j(x_j)^{-1} \hat{S}_{k,j}(x_k, x_j) - M_j(x_j)^{-1} S_{k,j}(x_k, x_j) \right]_{r,s}^2 E_j(x_j) E_k^{-1}(x_k) dx_j dx_k = o_P(1),$$

for $r, s = 1, 2$. Here $[A]_{r,s}$ denotes the element (r, s) of a matrix A . Moreover, \hat{M}_j vanishes outside the support of E_j , $\hat{S}_{j,k}$ vanishes outside the support of $E_{j,k}$ and \hat{S} is symmetric, i.e. $\hat{S}_{j,k}(x_j, x_k)^T = \hat{S}_{k,j}(x_k, x_j)$.

B3' There exists a constant C such that with probability tending to 1 for all j ,

$$\int \hat{\alpha}_j(x_j)^2 E_j(x_j) dx_j \leq C,$$

and

$$\int \hat{\alpha}^j(x_j)^2 E_j(x_j) dx_j \leq C.$$

B4' For some finite intervals $S_j \subset \mathbb{R}$ that are contained in the support of E_j , $j = 0, \dots, d$, we suppose that there exists a finite constant C such that with probability tending to 1 for all $j \neq k$,

$$\sup_{x_j \in S_j} \int \text{trace} \left[\hat{S}_{k,j}(x_k, x_j) \hat{M}_j(x_j)^{-2} \hat{S}_{k,j}(x_k, x_j) \right] E_k(x_k)^{-1} dx_k \leq C.$$

We now introduce the notation $\hat{\alpha}_j = \hat{\alpha}_j^A + \hat{\alpha}_j^B$ and $\hat{\alpha}^j = \hat{\alpha}^{j,A} + \hat{\alpha}^{j,B}$. Where $(\hat{\alpha}_j^A, \hat{\alpha}^{j,A})$ is the variable part and $(\hat{\alpha}_j^B, \hat{\alpha}^{j,B})$ is the stable part of the initialization $(\hat{\alpha}_j, \hat{\alpha}^j)$. The terms are given by

$$\begin{aligned} \hat{\alpha}_j^A(x_j) &= \left\{ (\hat{V}_j^j(x_j))^2 - \hat{V}_{j,j}^j(x_j) \hat{V}^j(x_j) \right\}^{-1} \frac{1}{n} \sum_{i=1}^n \int g_{i,j}(x_j) dM_i(s), \\ \hat{\alpha}^{j,A}(x_j) &= \left\{ (\hat{V}_j^j(x_j))^2 - \hat{V}_{j,j}^j(x_j) \hat{V}^j(x_j) \right\}^{-1} \frac{1}{n} \sum_{i=1}^n \int g_i^j(x_j) dM_i(s), \\ \hat{\alpha}_j^B(x_j) &= \left\{ (\hat{V}_j^j(x_j))^2 - \hat{V}_{j,j}^j(x_j) \hat{V}_{0,0}^j(x_j) \right\}^{-1} \frac{1}{n} \sum_{i=1}^n \int g_{i,j}(x_j) d\Lambda_i(s), \\ \hat{\alpha}^{j,B}(x_j) &= \left\{ (\hat{V}_j^j(x_j))^2 - \hat{V}_{j,j}^j(x_j) \hat{V}_{0,0}^j(x_j) \right\}^{-1} \frac{1}{n} \sum_{i=1}^n \int g_i^j(x_j) d\Lambda_i(s), \end{aligned}$$

with

$$\begin{aligned} g_{i,j}(x_j) &= \left[\hat{V}_j^j(x_j) \left(\frac{x_j - X_{ij}(s)}{h} \right) - \hat{V}_{j,j}^j(x_j) \right] k_h(x_j - X_{ij}(s)), \\ g_i^j(x_j) &= \left[\hat{V}_j^j(x_j) - \hat{V}^j(x_j) \left(\frac{x_j - X_{ij}(s)}{h} \right) \right] k_h(x_j - X_{ij}(s)). \end{aligned}$$

Equivalently, we can write

$$\begin{aligned} \begin{pmatrix} \hat{\alpha}_j^A(x_j) \\ \hat{\alpha}^{j,A}(x_j) \end{pmatrix} &= \frac{1}{n} \sum_{i=1}^n \int \begin{pmatrix} 1 \\ h^{-1}(x_j - X_{ij}(s)) \end{pmatrix} k_h(x_j, X_{ij}(s)) dM_i(s), \\ \begin{pmatrix} \hat{\alpha}_j^B(x_j) \\ \hat{\alpha}^{j,B}(x_j) \end{pmatrix} &= \frac{1}{n} \sum_{i=1}^n \int \hat{M}_j(x_j)^{-1} \begin{pmatrix} 1 \\ h^{-1}(x_j - X_{ij}(s)) \end{pmatrix} k_h(x_j, X_{ij}(s)) d\Lambda_i(s), \end{aligned}$$

As in Assumption B4, M_i is the martingale arising from N_i and Λ_i is its compensator. Later on, we will verify the following assumptions on $(\hat{\alpha}_j^A, \hat{\alpha}^{j,A})$ and $(\hat{\alpha}_j^B, \hat{\alpha}^{j,B})$. Moreover, for the whole estimator we define, for $s \in \{A, B\}$, $\tilde{\alpha}_{0,j}^s$, $\tilde{\alpha}_j^s$ and $\tilde{\alpha}^{j,s}$ as the solution of the equations

$$\hat{M}_j(x_j) \begin{pmatrix} \tilde{\alpha}_j^s(x_j) - \hat{\alpha}_j^s(x_j) \\ \tilde{\alpha}^{j,s}(x_j) - \hat{\alpha}^{j,s}(x_j) \end{pmatrix} = \tilde{\alpha}_{0,j}^s \begin{pmatrix} \hat{V}_j^j(x_j) \\ \hat{V}_j^j(x_j) \end{pmatrix} - \sum_{l \neq j} \int \hat{S}_{l,j}(x_l, x_j) \begin{pmatrix} \tilde{\alpha}_l^s(x_l) \\ \tilde{\alpha}^{l,s}(x_l) \end{pmatrix} dx_l, \quad (70)$$

$$\int \tilde{\alpha}_j^s(x_j) \hat{V}^j(x_j) dx_j = 0. \quad (71)$$

Existence and uniqueness of $\tilde{\alpha}_j^A, \tilde{\alpha}_j^B, \tilde{\alpha}^{j,A}, \tilde{\alpha}^{j,B}$ is stated in Proposition 4. We make the further assumptions

B5' There exists a constant C such that with probability tending to 1 for all j , it holds

$$\int \hat{\alpha}_j^s(x_j)^2 E_j(x_j) dx_j \leq C,$$

and

$$\int \hat{\alpha}^{j,s}(x_j)^2 E_j(x_j) dx_j \leq C,$$

for $s = A, B$.

B6' We assume that there is a sequence Δ_n such that

$$\begin{aligned} \sup_{x_k \in S_k} \left\| \int \hat{M}_k(x_k)^{-1} \hat{S}_{k,j}(x_k, x_j) \begin{pmatrix} \hat{\alpha}_j^A(x_j) \\ \hat{\alpha}_j^{j,A}(x_j) \end{pmatrix} dx_j \right\|_2 &= o_P(\Delta_n), \\ \left\| \int \hat{M}_k(x_k)^{-1} \hat{S}_{k,j}(x_k, x_j) \begin{pmatrix} \hat{\alpha}_j^A(x_j) \\ \hat{\alpha}_j^{j,A}(x_j) \end{pmatrix} dx_j \right\|_{M_k,2} &= o_P(\Delta_n), \end{aligned}$$

where $\|\cdot\|_2$ denotes the L_2 norm in \mathbb{R}^2 and where for functions $g : \mathbb{R} \rightarrow \mathbb{R}^2$ we define $\|g\|_{M_k,2}^2 = \int g(u) M_k(u) g(u) du$. The sets S_k have been introduced in Assumption B4'.

B7' There exist deterministic functions $\mu_{n,j}$ such that

$$\sup_{x_j \in S_j} |\tilde{\alpha}_j^B(x_j) - \mu_{n,j}(x_j)| = o_p(\Delta_n),$$

where S_k has been introduced in Assumption B4'.

The local linear equivalents to Propositions 1 and 2 are the following results from Mammen et al. [1999b], adapted to our setting. The following two propositions assure convergence of the backfitting algorithm and asymptotic normality of the stochastic part of the estimator under Assumptions B1'–B7'.

Proposition 4 (Convergence of backfitting). *Under Assumptions B1'–B3', with probability tending to 1, there exists a unique solution $\{\tilde{m}_{0,l}, \tilde{m}_l, \tilde{m}^l : l = 0, \dots, d\}$ to (26)–(28). Moreover, there exist constants $0 < \gamma < 1$ and $c > 0$ such that, with probability tending to 1, it holds:*

$$\begin{aligned} \int \left[\tilde{\alpha}_j^{[r]}(x_j) - \tilde{\alpha}_j(x_j) \right]^2 E_j(x_j) dx_j &\leq c\gamma^{2r}\Gamma, \\ \int \left[\tilde{\alpha}^{j,[r]}(x_j) - \tilde{\alpha}^j(x_j) \right]^2 E_j(x_j) dx_j &\leq c\gamma^{2r}\Gamma, \end{aligned}$$

where

$$\Gamma = 1 + \sum_{l=0}^d \int \left[\tilde{\alpha}_l^{[0]}(x_l) \right]^2 E_l(x_l) dx_l + \int \left[\tilde{\alpha}^{l,[0]}(x_l) \right]^2 E_l(x_l) dx_l.$$

The functions $\tilde{\alpha}_{0,l}^{[0]}, \tilde{\alpha}_l^{[0]}$ and $\tilde{\alpha}^{l,[0]}$ are the starting values of the backfitting algorithm. For $r > 0$ the functions $\tilde{\alpha}_l^{[r]}$ and $\tilde{\alpha}^{l,[r]}$ are defined by equations (29) and (30).

Moreover, under the additional Assumption B5', with probability tending to 1, there exists a solution $\{\tilde{\alpha}_0^s, \tilde{\alpha}_j^s, \tilde{\alpha}^{j,s} : j = 0, \dots, d\}$ of (70), (71) that is unique for $s = A, B$, respectively.

Proposition 5 (Asymptotic behavior of stochastic part). *Suppose that Assumptions B1'–B6' hold for a sequence Δ_n and intervals S_j , $j = 0, \dots, n$. Then it holds that*

$$\sup_{x_j \in S_j} |\tilde{\alpha}_j^A(x_j) - [\hat{\alpha}_j^A(x_j) - \tilde{\alpha}_{0,j}^A]| = o_P(\Delta_n).$$

Under the additional Assumption B7' it holds

$$\sup_{x_j \in S_j} |\tilde{\alpha}_j(x_j) - [\hat{\alpha}_j^A(x_j) - \tilde{\alpha}_{0,j}^A + \mu_{n,j}(x_j)]| = o_P(\Delta_n).$$

Before stating a result for the bias part, we assume the following.

B8' For all $j \neq k$, it holds

$$\sup_{x_j \in S_j} \int \left| \left[\hat{M}_j(x_j)^{-1} \hat{S}_{k,j}(s_k, x_j) - M_j^{-1}(x_j) S_{k,j}(x_k, x_j) \right]_{r,s} \right| E_k(x_k) dx_k = o_p(1),$$

for $r, s = 1, 2$.

B9' There exist deterministic functions $a_{n,0}(x_0), \dots, a_{n,d}(x_d), a_n^0(x_0), \dots, a_n^d(x_d)$ and constants $a_n^*, \gamma_{n,0}, \dots, \gamma_{n,d}$ such that

$$\begin{aligned} \int a_{n,j}(x_j)^2 E_j(x_j) dx_j &< \infty, \\ \int a_n^j(x_j)^2 E_j(x_j) dx_j &< \infty, \\ \gamma_{n,j} - \int a_{n,j}(x_j) \hat{V}^j(x_j) dx_j &= o_P(\Delta_n), \\ \sup_{x_j \in S_j} |\tilde{\alpha}_j^B(x_j) - \hat{\mu}_{n,0} - \hat{\mu}_{n,j}(x_j)| &= o_P(\Delta_n), \\ \int |\hat{\alpha}_j^B(x_j) - \hat{\mu}_{n,0} - \hat{\mu}_{n,j}(x_j)|^2 E_j(x_j) dx_j &= o_P(\Delta_n^2), \\ \sup_{x_j \in S_j} |\hat{\alpha}^{j,B}(x_j) - \hat{\mu}_{n,0} - \hat{\mu}_n^j(x_j)| &= o_P(\Delta_n), \\ \int |\hat{\alpha}^{j,B}(x_j) - \hat{\mu}_n^j(x_j)|^2 E_j(x_j) dx_j &= o_P(\Delta_n^2), \end{aligned}$$

for random variables $\hat{\mu}_{n,0}$ and where

$$\begin{pmatrix} \hat{\mu}_{n,j}(x_j) \\ \hat{\mu}_n^j(x_j) \end{pmatrix} = \begin{pmatrix} a_{n,0} + a_{n,j}(x_j) \\ a_n^j(x_j) \end{pmatrix} + \sum_{k \neq j} \int \hat{M}_j(x_j)^{-1} \hat{S}_{k,j}(x_k, x_j) \begin{pmatrix} a_{n,k}(x_k) \\ a_n^k(x_k) \end{pmatrix} dx_k.$$

The next proposition appears in Mammen et al. [1999b] with different notation for the nonparametric regression case. It assures convergence of the deterministic part of the estimator.

Proposition 6 (Asymptotic behavior of bias part). *Under Assumptions B1'–B6', B8', B9', it holds*

$$\begin{aligned} \sup_{x_j \in S_j} |\tilde{\alpha}_j^B(x_j) - \mu_{n,j}(X_j)| &= o_P(\Delta_n), \\ \sup_{x_j \in S_j} |\tilde{\alpha}^{j,B}(x_j) - \mu_n^j(X_j)| &= o_P(\Delta_n), \end{aligned}$$

for $\mu_{n,j}(x_j) = a_{n,j}(x_j) - \gamma_{n,j}$ and $\mu_n^j(x_j) = a_n^j(x_j)$. Assumption B7' holds with this choice of $\mu_{n,j}(x_j)$.

Proof of Theorem 2. To apply Propositions 4–6, we have to prove that Assumptions A1–A5 imply B1–B6, B8, B9. The proof is analogous to the proof of Theorem 1 and the assumptions can be shown in a similar way.

We now focus on the variance and bias part

$$\begin{aligned}\begin{pmatrix} \hat{\alpha}_j^A(x_j) \\ \hat{\alpha}_{j,A}^B(x_j) \end{pmatrix} &= \hat{M}_j(x_j)^{-1} \frac{1}{n} \sum_{i=1}^n \int \left(h^{-1}(x_j - X_{ij}(s)) \right) k_h(x_j, X_{ij}(s)) dM_i(s), \\ \begin{pmatrix} \hat{\alpha}_j^B(x_j) \\ \hat{\alpha}_{j,B}^A(x_j) \end{pmatrix} &= \hat{M}_j(x_j)^{-1} \frac{1}{n} \sum_{i=1}^n \int \left(h^{-1}(x_j - X_{ij}(s)) \right) k_h(x_j, X_{ij}(s)) d\Lambda_i(s).\end{aligned}$$

Analogously to (38)–(41), we show uniform convergence of $\hat{M}_j(x_j)$ and $\hat{S}_{l,j}(x_l, x_j)$ to $M_j(x_j)$ and $S_{l,j}(x_l, x_j)$, respectively, and then focus on

$$\frac{1}{n} \sum_{i=1}^n \int \left(h^{-1}(x_j - X_{ij}(s)) \right) k_h(x_j, X_{ij}(s)) dM_i(s)$$

for asymptotic normality and on

$$\frac{1}{n} \sum_{i=1}^n \int \left(h^{-1}(x_j - X_{ij}(s)) \right) k_h(x_j, X_{ij}(s)) d\Lambda_i(s)$$

for a bias term.

With M_i being the same martingale as in the proof of Theorem 1 occurring in the stochastic part, we get the same asymptotic variance σ_j^2 . Moreover, Assumptions A6–A9 can be verified with the choices

$$\begin{aligned}\Delta_n &= h^2, \\ a_n^* &= \alpha^*, \\ a_{n,j}(x_j) &= \alpha_j(x_j) + \frac{1}{2} h^2 \alpha_j''(x_j) \int u^2 k(u) du, \\ a_n^j(x_j) &= h \alpha_j'(x_j), \\ \beta(x) &= \sum_{j=1}^d \frac{1}{2} \int u^2 k(u) du \left[\alpha_j''(x_j) - \int \alpha_j''(x_j) E_j(x_j) dx_j \right], \\ \gamma_{n,j} &= \nu_{n,j} + \frac{h^2}{2} \int u^2 k(u) du \int \alpha_j''(x_j) E_j(x_j) dx_j, \\ \nu_{n,j} &= \int \int \alpha_j(x_j) k_h(x_j, u) E_j(u) du dx_j.\end{aligned}$$

□

A.3. Two-step smooth backfitting estimator

The interpretation as a projection motivates two different ways to compute the smooth backfitting hazard estimator. For the minimisation over all additive hazard functions, we can either minimize directly or we first minimize over the subspace of all (unstructured) local polynomial functions of degree p obtaining a solution $\hat{\alpha}_{pilot}$ from (8) which is a non-additive estimator and then minimize the integrated squared

errors between $\hat{\alpha}_{pilot}$ and all additive local polynomial functions of degree p :

$$\begin{aligned} \arg \min_{\substack{\alpha^* \in \mathbb{R}, \\ \alpha_j^{(l)}: \mathbb{R} \rightarrow \mathbb{R}, \\ j=0, \dots, d \\ l=0, \dots, p}} \sum_{i=1}^n \int \int \left\{ \hat{\alpha}_{pilot}(x) - \left[\alpha^* + \alpha_0(t) + \alpha_1(z_1) + \dots + \alpha_d(z_d) \right. \right. \\ \left. \left. + \alpha_0^{(p)}(x_0) \left(\frac{x_0 - X_{i0}(s)}{h} \right)^p + \dots + \alpha_d^{(p)}(x_d) \left(\frac{x_d - X_{id}(s)}{h} \right)^p \right] \right\}^2 \\ \times K_h(x - X_i(s)) Y_i(s) ds d\nu(x). \end{aligned} \quad (72)$$

We want to emphasize that the estimator we obtain via direct minimisation (9) or (10), respectively, and the one obtained through the two-step minimisation (72) are identical.

In the following, we want to illustrate how the estimator can be obtained from an unstructured hazard estimator. Although we don't make use of it, this representation enables us to derive the asymptotic theory for the final estimator making use of the known asymptotic behavior of the established unstructured local constant which is defined below. Moreover, the derivation is less technical and easier to follow and the implementation is more straightforward.

Let $\hat{\alpha}$ be the unstructured local constant pilot estimator, $\hat{\alpha}^{LC}$ defined in Section 4.3. Then, for a weighting w , the local constant smooth backfitting estimator $\bar{\alpha}$ can be equivalently defined as

$$\min_{\bar{\alpha}} \int_{\mathcal{X}} \left(\hat{\alpha}(x) - [\bar{\alpha}^* + \sum_{j=0}^d \bar{\alpha}_j(x_j)] \right)^2 w(x) dx.$$

Analogously, for $p = 1$ we get the local linear estimator $\hat{\alpha}^{LL}(x) = \hat{O}^{LL}(x) / \hat{E}^{LL}(x)$ for $x \in \mathcal{X}$ from equation (8), which is defined through

$$\begin{aligned} \hat{O}^{LL}(x) &= \frac{1}{n} \sum_{i=1}^n \int \{1 - (x - X_i(s))D(x)^{-1}c_1(x)\} K_h(x, X_i(s)) dN_i(s), \\ \hat{E}^{LL}(x) &= \frac{1}{n} \sum_{i=1}^n \int \{1 - (x - X_i(s))D(x)^{-1}c_1(x)\} K_h(x, X_i(s)) Y_i(s) ds, \end{aligned}$$

where $c_j(x) = n^{-1} \sum_{i=1}^n \int K_h(x, X_i(s))(x_j - X_{ij}(s))Y_i(s)ds$ and for the $(d+1) \times (d+1)$ -matrix $D(x) = [d_{jk}(x)]_{jk}$ with $d_{jk}(x) = \frac{1}{n} \sum_{i=1}^n \int K_h(x, X_i(s))(x_j - X_{ij}(s))(x_k - X_{ik}(s))Y_i(s)ds$.

Note that the matrix D is not necessarily regular for $d > 2$ and hence the existence of D^{-1} and the existence of $\hat{\alpha}^{LL}$ are not guaranteed for $d > 2$.

In contrast to the local linear estimator, the local constant estimator $\hat{\alpha}^{LC}$ is always well defined independent of the dimension d .

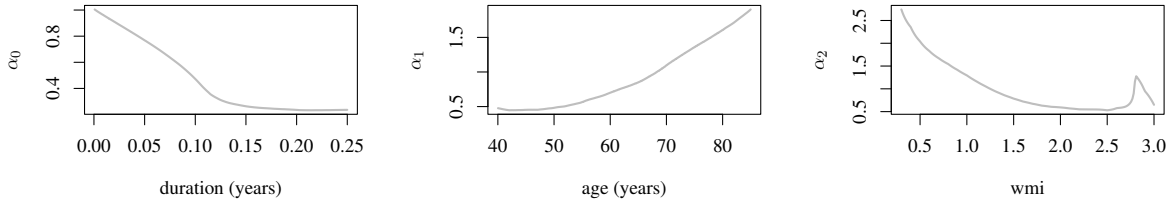
B. Fitted values from the multiplicative model

Here we show how the fitted values from the local constant multiplicative smooth backfitting model Hiabu et al. [2021a] look like.

Acknowledgment

This work was supported by the Deutsche Forschungsgemeinschaft (DFG) through the Research Training Group RTG 1953.

Hazard function estimates for the first three months
vf=0



vf=1

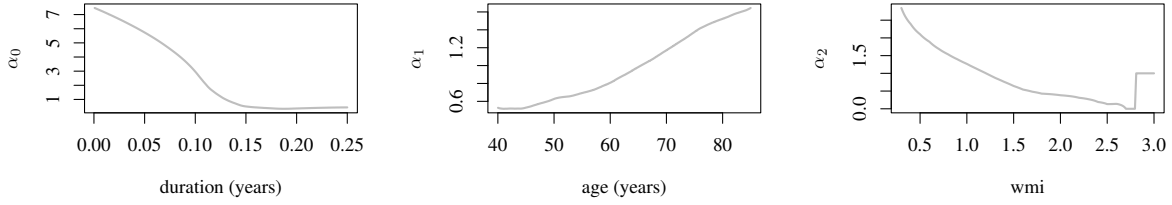
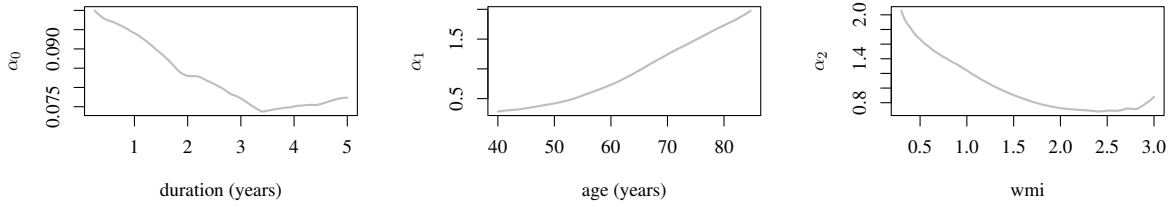


Figure 8: Local constant multiplicative smooth backfitting fit of $(\alpha_0, \alpha_1, \alpha_2)$ conditional on surviving the first three months for two different strata depending on the value of vf .

Hazard function estimates conditional on surviving the first three months
vf=0



vf=1

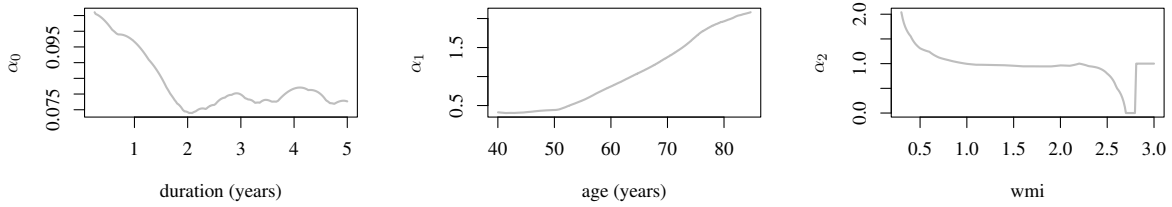


Figure 9: Local constant multiplicative smooth backfitting fit of $(\alpha_0, \alpha_1, \alpha_2)$ conditional on surviving the first three months for two different strata depending on the value of vf .

References

- O. O. Aalen. A model for nonparametric regression analysis of counting processes. In W. Klonecki, A. Kozek, and J. Rosiński, editors, *Mathematical Statistics and Probability Theory*, pages 1–25, New York, 1980. Springer New York.
- O. O. Aalen, O. Borgan, and H. Gjessing. *Survival and event history analysis: a process point of view*. Springer Science & Business Media, 2008.
- O. O. Aalen, M. Stensrud, V. Didelez, R. Daniel, K. Røysland, and S. Strohmaier. Time-dependent mediators in survival analysis: Modeling direct and indirect effects with the additive hazards model. *Biometrical Journal*, 2019. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/bimj.201800263>.
- P. Andersen, O. Borgan, R. Gill, and N. Keiding. *Statistical Models Based on Counting Processes*. Springer, New York, 1993.
- A. Avati, T. Duan, S. Zhou, K. Jung, N. H. Shah, and A. Y. Ng. Countdown regression: sharp and calibrated survival predictions. In *Uncertainty in Artificial Intelligence*, pages 145–155. PMLR, 2020.
- R. Beran. Nonparametric regression with randomly censored survival data. Technical report, Department of Statistics, University of California, Berkeley, 1981.
- N. E. Breslow and N. E. Day. Statistical methods in cancer research, vol. 2. *The Design and Analysis of Cohort Data*, Lyon, IARC, 1987.
- N. Bissantz, H. Dette, T. Hildebrandt, and K. Bissantz. Smooth backfitting in additive inverse regression. *Annals of the Institute of Statistical Mathematics*, 68(4):827–853, 2016.
- A. Buja, T. Hastie, and R. Tibshirani. Linear smoothers and additive models. *The Annals of Statistics*, 17(2):453–510, 1989a.
- A. Buja, T. Hastie, and R. Tibshirani. Linear smoothers and additive models. *The Annals of Statistics*, 17:453–510, 1989b.
- D. R. Cox. Regression models and life tables. *Journal of the Royal Statistical Society: Series B*, 34(2): 187–220, 1972.
- O. Dukes, T. Martinussen, E. J. Tchetgen Tchetgen, and S. Vansteelandt. On doubly robust estimation of the hazard difference. *Biometrics*, 75(1):100–109, 2019.
- K. Gregory, E. Mammen, and M. Wahl. Optimal estimation of sparse high-dimensional additive models. *The Annals of Statistics*, forthcoming, 2020.
- K. Han, B. U. Park, et al. Smooth backfitting for errors-in-variables additive models. *The Annals of Statistics*, 46:2216–2250, 2018.
- K. Han, H.-G. Müller, and B. U. Park. Additive functional regression for densities as responses. *Journal of the American Statistical Association*, 115:997–1010, 2020.
- M. Hiabu, E. Mammen, M. D. Martínez-Miranda, and J. P. Nielsen. Smooth backfitting of proportional hazards with multiplicative components. *Journal of the American Statistical Association*, 116(536): 1983–1993, 2021a.
- M. Hiabu, J. P. Nielsen, and T. H. Scheike. Nonsmooth backfitting for the excess risk additive regression model with two survival time scales. *Biometrika*, 108(2):491–506, 2021b.

- M. Hiabu, E. Mammen, and J. T. Meyer. Local linear smoothing in additive models as data projection. In D. Belomestny, C. Butucea, E. Mammen, E. Moulines, M. Reiß, and V. V. Ulyanov, editors, *Foundations of Modern Statistics*, pages 197–223, Cham, 2023. Springer International Publishing. ISBN 978-3-031-30114-8.
- L.-S. Huang and C.-H. Yu. Classical backfitting for smooth-backfitting additive models. *Journal of Computational and Graphical Statistics*, pages 1–22, 2019.
- F. W. Huffer and I. W. McKeague. Weighted least squares estimation for Aalen’s additive risk model. *Journal of the American Statistical Association*, 86(413):114–129, 1991.
- G. V. H. Jensen, C. Torp-Pedersen, P. Hildebrandt, L. Kober, F. Nielsen, T. Melchior, T. Joen, and P. Andersen. Does in-hospital ventricular fibrillation affect prognosis after myocardial infarction? *European heart journal*, 18(6):919–924, 1997.
- J. M. Jeon, B. U. Park, et al. Additive regression with Hilbertian responses. *The Annals of Statistics*, 48:2671–2697, 2020.
- Ø. Kravdal. The attractiveness of an additive hazard model: An example from medical demography. *European Journal of Population / Revue Européenne de Démographie*, 13(1):33–47, 1997.
- D. Y. Lin and Z. Ying. Semiparametric analysis of the additive risk model. *Biometrika*, 81(1):61–71, 1994.
- O. B. Linton, J. P. Nielsen, and S. Van de Geer. Estimating multiplicative and additive hazard functions by kernel methods. *The Annals of Statistics*, 31(1):464–492, 2003.
- E. Mammen and J. P. Nielsen. Generalised structured models. *Biometrika*, 90:551–566, 2003.
- E. Mammen and S. Sperlich. Additivity tests based on smooth backfitting. *Biometrika*, forthcoming, 2021.
- E. Mammen and K. Yu. Nonparametric estimation of noisy integral equations of the second kind. *Journal of the Korean Statistical Society*, 38:99–110, 2009.
- E. Mammen, O. Linton, and J. Nielsen. The existence and asymptotic properties of a backfitting projection algorithm under weak conditions. *The Annals of Statistics*, 27:1443–1490, 1999a.
- E. Mammen, O. B. Linton, and J. P. Nielsen. The existence and asymptotic properties of a backfitting projection algorithm under weak conditions. *The Annals of Statistics*, 27:1443–1490, 1999b.
- E. Mammen, J. S. Marron, B. A. Turlach, and M. P. Wand. A general framework for constrained smoothing. *Statistical Science*, 16:232–248, 2001.
- E. Mammen, B. U. Park, and M. Schienle. Additive models: Extensions and related models. In J. S. Racine, L. Su, and A. Ullah, editors, *The Oxford Handbook of Applied Nonparametric and Semiparametric Econometrics and Statistics*. Oxford Univ. Press, 2014.
- T. Martinussen and T. H. Scheike. A flexible additive multiplicative hazard model. *Biometrika*, 89(2): 283–298, 2002.
- T. Martinussen and T. H. Scheike. *Dynamic regression models for survival data*. Springer, New York, 2006.
- L. S. McDaniel, M. Yu, and R. Chappell. Analysis and design of clinical trials using additive hazards survival endpoints. *Statistics in Biopharmaceutical Research*, 11(3):274–282, 2019.

- I. W. McKeague. Asymptotic theory for weighted least squares estimators in aalen's additive risk model. In *Statistical Inference from Stochastic Processes: Proceedings of the AMS-IMS-SIAM Joint Summer Research Conference Held August 9–15, 1987, with Support from the National Science Foundation and the Army Research Office*, volume 80, pages 139–152. American Mathematical Society, 1988.
- I. W. McKeague and K. J. Utikal. Inference for a nonlinear counting process regression model. *The Annals of Statistics*, 18(3):1172–1187, 1990.
- J. P. Nielsen. Multiplicative bias correction in kernel hazard estimation. *Scandinavian Journal of Statistics*, 25(3):541–553, 1998.
- J. P. Nielsen and O. B. Linton. Kernel estimation in a non-parametric marker dependent hazard model. *The Annals of Statistics*, 23:1735–1748, 1995.
- J. P. Nielsen and S. Sperlich. Smooth backfitting in practice. *Journal of the Royal Statistical Society: Series B*, 67:43–61, 2005.
- J. P. Nielsen and C. Tanggaard. Boundary and bias correction in kernel hazard estimation. *Scandinavian Journal of Statistics*, 28:675–698, 2001.
- H. Ramlau-Hansen. Smoothing counting process intensities by means of kernel functions. *The Annals of Statistics*, 11:453–466, 1983.
- T. Scheike. Timereg package. *R package version*, 3(0):1–1, 2009.
- L. Spierdijk. Nonparametric conditional hazard rate estimation: a local linear approach. *Computational Statistics & Data Analysis*, 52(5):2419–2434, 2008.
- C. J. Stone. Optimal rates of convergence for nonparametric estimators. *The Annals of Statistics*, 8(6): 1348–1360, 1980.
- E. J. Tchetgen Tchetgen, S. Walter, S. Vansteelandt, T. Martinussen, and M. Glymour. Instrumental variable estimation in a survival context. *Epidemiology*, 26(3):402–410, 2015.
- I. Van Keilegom and N. Veraverbeke. Hazard rate estimation in nonparametric regression with censored data. *Annals of the Institute of Statistical Mathematics*, 53(4):730–745, 2001.
- K. Yu, B. U. Park, E. Mammen, et al. Smooth backfitting in generalized additive models. *The Annals of Statistics*, 36:228–260, 2008.