



City Research Online

City St George's, University of London

Citation: Dettoni, R., Marra, G. & Radice, R. (2020). Generalized Link-Based Additive Survival Models with Informative Censoring. *Journal of Computational and Graphical Statistics*, 29(3), pp. 503-512. doi: 10.1080/10618600.2020.1724544

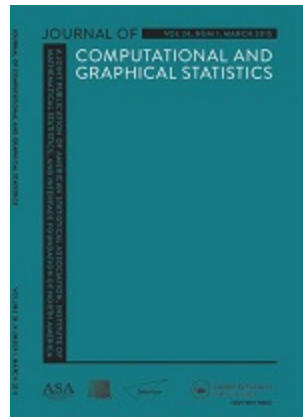
This is the accepted version of the paper.

This version of the publication may differ from the final published version. To cite this item please consult the publisher's version.

Permanent repository link: <https://openaccess.city.ac.uk/id/eprint/23594/>

Link to published version: <https://doi.org/10.1080/10618600.2020.1724544>

Copyright and Reuse: Copyright and Moral Rights remain with the author(s) and/or copyright holders. Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge, unless otherwise indicated, 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. For full details of reuse please refer to [City Research Online policy](#).



Generalized Link-Based Additive Survival Models with Informative Censoring

Journal:	<i>Journal of Computational and Graphical Statistics</i>
Manuscript ID	JCGS-19-163.R3
Manuscript Type:	Original Article
Keywords:	additive predictor, informative censoring, link-based survival model, penalised maximum likelihood, smoothing
Note: The following files were submitted by the author for peer review, but cannot be converted to PDF. You must view these files (e.g. movies) online.	
Supplemental Materials.zip	

Generalized Link-Based Additive Survival Models with Informative Censoring

Robinson Dettoni^{*†} Giampiero Marra[‡] Rosalba Radice[§]

2020-01-02

Abstract

Time to event data differ from other types of data because they are censored. Most of the related estimation techniques assume that the censoring mechanism is non-informative while in many applications it can actually be informative. The aim of this work is to introduce a class of flexible survival models which account for the information provided by the censoring times. The baseline functions are estimated non-parametrically by monotonic P-splines, whereas covariate effects are flexibly determined using additive predictors. Parameter estimation is reliably carried out within a penalised maximum likelihood framework with integrated automatic multiple smoothing parameter selection. We derive the \sqrt{n} -consistency and asymptotic normality of the non-informative and informative estimators, and shed light on the efficiency gains produced by the newly introduced informative estimator when compared to its non-informative counterpart. The finite sample properties of the estimators are investigated via a Monte Carlo simulation study which highlights the good empirical performance of the proposal. The modelling framework is illustrated on data about infants hospitalised for pneumonia. The models and methods discussed in the paper have been implemented in the R package GJRM to allow for transparent and reproducible research.

^{*}Department of Economics, Universidad de Santiago de Chile, Santiago 9170 022, Chile, robinson.dettoni@usach.cl

[†]Department of Statistical Science, University College London, Gower Street, London WC1E 6BT, UK, r.dettoni@ucl.ac.uk.

[‡]Department of Statistical Science, University College London, Gower Street, London WC1E 6BT, UK, giampiero.marra@ucl.ac.uk.

[§]Cass Business School, City, University of London, 106 Bunhill Row, EC1Y 8TZ London, UK, rosalba.radice@city.ac.uk.

Key Words: additive predictor, informative censoring, link-based survival model, penalised maximum likelihood, smoothing.

1 Introduction

Time to event data are different from other types of data because of censoring. This means that the response of interest, the time until a particular event occurs, can not be totally observed. As a result, models must be used to relate the observed and unobserved parts of the data since the recorded observations alone can not provide direct information on the event of interest. Most of the related estimation techniques assume that the censoring scheme is independent and non-informative conditional on covariates (e.g., Cox, 1972; Ma et al., 2014; Scheike & Zhang, 2003; Xue et al., 2018; Younes & Lachin, 1997). In many applications, however, these assumptions can at least be questioned (e.g., Chen, 2010; Huang & Zhang, 2008; Li & Peng, 2015; Lu & Zhang, 2012; Slud & Rubinstein, 1983; Wang et al., 2015; Xu et al., 2017, 2018; Zheng & Klein, 1995; Zeng et al., 2004).

Censoring is independent when the hazard rate of the event of interest for the censored observations is equal to the hazard rate for the uncensored ones, otherwise it is called dependent (Kalbfleisch & Prentice, 2002). If the event and censoring times are assumed to be dependent, then survival models accounting for this feature of the data face a problem of identification. In general, without additional assumptions, it is not possible to identify the survival distribution from the censored data alone or testing whether the censoring and survival mechanisms are independent (Cox, 1959; Tsiatis, 1975).

Censoring is informative when the censoring times, say T_2 , contain information on the parameters of the distribution of the event variable, say T_1 (Lagakos, 1979; Kalbfleisch & Prentice, 2002). In particular, let us write the hazard functions for the event and censored times as $h_{T_1}(t|\mathbf{x}_{T_1}; \boldsymbol{\theta}_{T_1})$ and $h_{T_2}(t|\mathbf{x}_{T_2}; \boldsymbol{\theta}_{T_2})$. If the vector of parameters $\boldsymbol{\theta}_{T_1}$ and $\boldsymbol{\theta}_{T_2}$ have components in common then censoring is informative. In this case, the observable data $(Y, \delta) = \{\min(T_1, T_2), I(T_1 < T_2)\}$, where I is the usual indicator function, provide sufficient information to identify the marginal survival functions of T_1 and T_2 (Kalbfleisch & Prentice, 2002).

1 Although dependent censoring is a well studied problem in the survival analysis and compet-
2 ing risk literature (e.g., Crowder, 2012; Emura & Chen, 2018), the specific literature analysing the
3 problem of informative censoring is scarce, even though ignoring it may have detrimental conse-
4 quences on inferential conclusions (e.g., Siannis et al., 2005; Lu & Zhang, 2012). In a seminal
5 work, Koziol & Green (1976) proposed an informative survival model where the hazard functions
6 of T_1 and T_2 satisfy $h_{T_2}(t) = p h_{T_1}(t)$, for some constant $0 < p < 1$. Since this model did not
7 incorporate covariates, it was further extended. For instance, Yuan (2005) introduced a semipara-
8 metric Cox model estimated via profile likelihood in which, for a given vector of covariates \mathbf{x} ,
9 $h_{T_2}(t|\mathbf{x}) = \varrho(t, \mathbf{x}; \theta) h_{T_1}(t|\mathbf{x})$, where ϱ is a function known up to a finite-dimensional parameter, θ .
10 The purpose of ϱ was to capture the possible information contained in the censoring times. Lu &
11 Zhang (2012) proposed a semi-parametric informative survival model where the baseline hazards
12 are estimated non-parametrically and the covariate effects parametrically. In their approach, the
13 hazard functions of T_1 and T_2 conditional on \mathbf{x} are modelled using $h_{T_v}(t|\mathbf{x}) = h_{0,T_v}(t) \exp(\mathbf{x}^\top \boldsymbol{\varphi}_v)$,
14 where $\mathbf{x}^\top \boldsymbol{\varphi}_v = \mathbf{x}_1^\top \boldsymbol{\vartheta}_0 + \mathbf{x}_2^\top \boldsymbol{\vartheta}_v$, for $v = 1, 2$.

15 In this article we deal with informative censoring. In particular, we develop a flexible, gen-
16 eral and tractable survival modelling framework where the baseline functions are estimated non-
17 parametrically via means of monotonic P-splines, covariate effects are flexibly determined using
18 additive predictors, and informative censoring is accounted for. Model fitting is based on an op-
19 timization scheme that allows for the reliable simultaneous penalized estimation of all model's
20 parameters as well as for stable and fast automatic multiple smoothing parameter selection. We
21 provide the \sqrt{n} -consistency and asymptotic normality of the non-informative and informative es-
22 timators, and show that the newly introduced informative estimator is more efficient than its non-
23 informative counterpart. A Monte Carlo simulation study highlights the merits of the proposal, and
24 the modelling framework is illustrated on data about infants hospitalised for pneumonia. The mod-
25 els and methods introduced in the article have been implemented in the R package GJRM (Marra
26 & Radice, 2019) to allow for transparent and reproducible research. To the best of our knowl-
27 edge, there are no alternative flexible survival models with informative censoring, nor respective
28 software implementations, of the type proposed here. Given that the assumption of absence of in-
29 formative censoring is often made for convenience, the proposed methodology is likely to appeal

the wider audience wishing to estimate possibly more realistic survival models or at least assess whether allowing for informative censoring can produce more plausible results.

The article is organized as follows. In the next section, the proposed model and its theoretical properties are discussed. In Section 3, the effectiveness of the proposed methodology is explored by means of a simulation study. In Section 4, the framework is illustrated on data about infants hospitalised for pneumonia. Section 5 concludes the paper with a discussion.

2 Methodology

In this paper, the case of right censored data is considered; the true event time is not always observed, in which case censoring (lower) times are observed. For individual i , where $i = 1, \dots, n$ and n represents the sample size, let T_{1i} and T_{2i} denote the true event and censoring times. Let also $\mathbf{z}_{\nu i}^T = (z_{\nu 1i}, \dots, z_{\nu K_{\nu}i})$ be a vector of baseline covariates of dimension K_{ν} , where \mathbf{z}^T stands for the transpose of a vector \mathbf{z} , $\nu = 1, 2$ and $\mathbf{z}_i^T = (\mathbf{z}_{1i}, \mathbf{z}_{2i})$. It is assumed that the (T_{1i}, \mathbf{z}_i) , for $i = 1, \dots, n$, are independently and identically distributed (*i.i.d.*). The censoring times, T_{2i} , are also assumed to be *i.i.d.* The distribution of T_2 depends on \mathbf{z} . In addition, we assume that T_{1i} and T_{2i} are conditionally independent given \mathbf{z}_i , and that T_{1i} is informatively right censored by T_{2i} through some covariates (Andersen & Keiding, 2006). We observe $(Y_i, \mathbf{z}_i, \delta_{1i})$, where $Y_i = \min\{T_{1i}, T_{2i}\}$ and $\delta_{1i} = I(T_{1i} \leq T_{2i})$. We also define $\delta_{2i} = [1 - \delta_{1i}]$. Finally, $\boldsymbol{\theta}$ is a generic vector of parameters.

2.1 Survival functions

The survival function of $T_{\nu i}$ taking values in $(0, 1)$, conditional on $\mathbf{z}_{\nu i}$ and $\boldsymbol{\theta}_{\nu}$, can be expressed as

$$P(T_{\nu i} > t_{\nu i} | \mathbf{z}_{\nu i}; \boldsymbol{\theta}_{\nu}) = S_{\nu}(t_{\nu i} | \mathbf{z}_{\nu i}; \boldsymbol{\theta}_{\nu}) = \mathcal{G}_{\nu}[\xi_{\nu i}(t_{\nu i}, \mathbf{z}_{\nu i}; \boldsymbol{\theta}_{\nu})], \quad (1)$$

where, for $\nu = 1, 2$, $\boldsymbol{\theta}_{\nu}$ and $\mathbf{z}_{\nu i}$ represent generic vectors of coefficients and covariates, respectively. The survival functions are modelled using generalised survival or link-based functions models (Younes & Lachin, 1997; Liu et al., 2018). That is, $S_{\nu}(t_{\nu i} | \mathbf{z}_{\nu i}; \boldsymbol{\theta}_{\nu})$ is defined as $\mathcal{G}_{\nu}[\xi_{\nu i}(t_{\nu i}, \mathbf{z}_{\nu i}; \boldsymbol{\theta}_{\nu})]$, where \mathcal{G}_{ν} is an inverse link function. The set up of the two ξ predictors is

discussed in the detail in the next section. As conveyed by the notation, ξ_{1i} and ξ_{2i} must include baseline functions of time. Different choices for function $\mathcal{G}_\nu[\xi_{\nu i}(t_{\nu i}, \mathbf{z}_{\nu i}; \boldsymbol{\theta}_\nu)]$ can be specified; some common examples are shown in Table 1 reported in Supplementary Material A. The cumulative hazard function, \mathcal{H}_ν , and the hazard function, h_ν , are given by

$$\begin{aligned}\mathcal{H}_\nu(t_{\nu i}|\mathbf{z}_{\nu i}; \boldsymbol{\theta}_\nu) &= -\log \mathcal{G}_\nu[\xi_{\nu i}(t_{\nu i}, \mathbf{z}_{\nu i}; \boldsymbol{\theta}_\nu)], \\ h_\nu(t_{\nu i}|\mathbf{z}_{\nu i}; \boldsymbol{\theta}_\nu) &= -\frac{\mathcal{G}'_\nu[\xi_{\nu i}(t_{\nu i}, \mathbf{z}_{\nu i}; \boldsymbol{\theta}_\nu)]}{\mathcal{G}_\nu[\xi_{\nu i}(t_{\nu i}, \mathbf{z}_{\nu i}; \boldsymbol{\theta}_\nu)]} \frac{\partial \xi_{\nu i}(t_{\nu i}, \mathbf{z}_{\nu i}; \boldsymbol{\theta}_\nu)}{\partial t_{\nu i}},\end{aligned}\quad (2)$$

where $\mathcal{G}'_\nu[\xi_{\nu i}(t_{\nu i}, \mathbf{z}_{\nu i}; \boldsymbol{\theta}_\nu)] = \partial \mathcal{G}_\nu[\xi_{\nu i}(t_{\nu i}, \mathbf{z}_{\nu i}; \boldsymbol{\theta}_\nu)] / \partial \xi_{\nu i}(t_{\nu i}, \mathbf{z}_{\nu i}; \boldsymbol{\theta}_\nu)$.

2.2 Additive predictors

This section provides some details on the set up of the two model's predictors for the cases of informative and non-informative censoring. Note that these must include baseline functions of time. To make the presentation simpler, the same design matrix is set up for the two additive predictors. Also, $t_{\nu i}$ can be treated like a covariate. The main advantages of using additive predictors are that various types of covariate effects can be dealt with and that such effects can be flexibly determined without making strong parametric a priori assumptions regarding their forms (e.g., Wood, 2017).

Let us consider a generic predictor $\xi_{\nu i} \in \mathbb{R}$ (where the dependence on the covariates and parameters is momentarily dropped), and the overall baseline covariate vector $\mathbf{x}_{\nu i}$, which contains $\mathbf{z}_{\nu i}$ and $t_{\nu i}$. The additive predictors for the censoring and event times can be defined generically as

$$\xi_{\nu i} = \alpha_{\nu 0} + \sum_{k_\nu=0}^{K_\nu} s_{\nu k_\nu}(\mathbf{x}_{\nu k_\nu i}), \quad i = 1, \dots, n. \quad (3)$$

In (3), $\alpha_{\nu 0} \in \mathbb{R}$ is an overall intercept, $\mathbf{x}_{\nu k_\nu i}$ denotes the k_ν^{th} sub-vector of the complete vector $\mathbf{x}_{\nu i}$ and the K_ν functions $s_{\nu k_\nu}(\mathbf{x}_{\nu k_\nu i})$ represent generic effects which are chosen according to the type of covariate(s) considered. Note that, in (3), k_ν starts from 0 since the summation also includes a smooth function of time.

If censoring is informative, some covariates in \mathbf{x}_{1i} must also appear in \mathbf{x}_{2i} . In particular, let us define the vectors of informative and non-informative covariates of dimensions Q and Q_ν as

$\mathbf{x}_i^0 \top = (x_{1i}^0, \dots, x_{Qi}^0)$ and $\mathbf{x}_{\nu i}^1 \top = (x_{\nu 1i}^1, \dots, x_{\nu Q_{\nu}i}^1)$, where $K_{\nu} = Q + Q_{\nu}$. Informative censoring implies that some components of $\sum_{k_1=1}^{K_1} s_{1k_1}(\mathbf{x}_{1k_1i})$ must appear in $\sum_{k_2=1}^{K_2} s_{2k_2}(\mathbf{x}_{2k_2i})$. Without loss of generality, we assume that the first Q components in $\sum_{k_1=1}^{K_1} s_{1k_1}(\mathbf{x}_{1k_1i})$ appear in $\sum_{k_2=1}^{K_2} s_{2k_2}(\mathbf{x}_{2k_2i})$. That is,

$$\begin{aligned} \sum_{k_1=1}^{K_1} s_{1k_1}(\mathbf{x}_{1k_1i}) &= \sum_{q_0=1}^Q s_q(\mathbf{x}_{qi}^0) + \sum_{q_1=1}^{Q_1} s_{1q_1}(\mathbf{x}_{1q_1i}^1) \\ \sum_{k_2=1}^{K_2} s_{2k_2}(\mathbf{x}_{2k_2i}) &= \sum_{q=1}^Q s_q(\mathbf{x}_{qi}^0) + \sum_{q_2=1}^{Q_2} s_{2q_2}(\mathbf{x}_{2q_2i}^1) \end{aligned} \quad (4)$$

Therefore, using (4), equation (3) becomes

$$\xi_{\nu i} = \alpha_{\nu 0} + \sum_{q=1}^Q s_q(\mathbf{x}_{qi}^0) + \sum_{q_{\nu}=0}^{Q_{\nu}} s_{\nu q_{\nu}}(\mathbf{x}_{\nu q_{\nu}i}^1), \quad (5)$$

where \mathbf{x}_{qi}^0 and $\mathbf{x}_{\nu q_{\nu}i}^1$ denote the informative and non-informative sub-vectors of the complete vectors \mathbf{x}_i^0 and $\mathbf{x}_{\nu i}^1$ respectively, and $s_{\nu q_{\nu}}(\mathbf{x}_{\nu q_{\nu}i}^1) = s_{\nu 0}(t_{\nu i})$ when $q_{\nu} = 0$.

In (5), the smooth functions are represented using the penalised regression spline approach (e.g., Wood, 2017). Specifically, each $s_{\nu q_{\nu}}(\mathbf{x}_{\nu q_{\nu}i}^1)$ can be approximated as a linear combination of $J_{\nu q_{\nu}}$ non-informative basis functions $\mathcal{Q}_{\nu q_{\nu} j_{\nu q_{\nu}}}(\mathbf{x}_{\nu q_{\nu}i}^1)$ and regression coefficients $\alpha_{\nu q_{\nu} j_{\nu q_{\nu}}} \in \mathbb{R}$. In a similar manner, each $s_q(\mathbf{x}_{qi}^0)$ can be approximated as a linear combination of J_q informative basis functions $\mathcal{Q}_{q j_q}(\mathbf{x}_{qi}^0)$ and regression coefficients $\alpha_{0 q j_q} \in \mathbb{R}$. More specifically, $s_q(\mathbf{x}_{qi}^0)$ and $s_{\nu q_{\nu}}(\mathbf{x}_{\nu q_{\nu}i}^1)$ are given by $s_q(\mathbf{x}_{qi}^0) = \sum_{j_q=1}^{J_q} \alpha_{0 q j_q} \mathcal{Q}_{q j_q}(\mathbf{x}_{qi}^0)$ and $s_{\nu q_{\nu}}(\mathbf{x}_{\nu q_{\nu}i}^1) = \sum_{j_{\nu q_{\nu}}=1}^{J_{\nu q_{\nu}}} \alpha_{\nu q_{\nu} j_{\nu q_{\nu}}} \mathcal{Q}_{\nu q_{\nu} j_{\nu q_{\nu}}}(\mathbf{x}_{\nu q_{\nu}i}^1)$, and therefore (5) can be written as

$$\xi_{\nu i} = \alpha_{\nu 0} + \sum_{q=1}^Q \mathcal{Q}_q(\mathbf{x}_{qi}^0) \top \alpha_{0q} + \sum_{q_{\nu}=0}^{Q_{\nu}} \mathcal{Q}_{\nu q_{\nu}}(\mathbf{x}_{\nu q_{\nu}i}^1) \top \alpha_{\nu q_{\nu}}, \quad (6)$$

where $\mathcal{Q}_q(\mathbf{x}_{qi}^0) \top \alpha_{0q} = \sum_{j_q=1}^{J_q} \alpha_{0 q j_q} \mathcal{Q}_{q j_q}(\mathbf{x}_{qi}^0)$, $\mathcal{Q}_{\nu q_{\nu}}(\mathbf{x}_{\nu q_{\nu}i}^1) \top \alpha_{\nu q_{\nu}} = \sum_{j_{\nu q_{\nu}}=1}^{J_{\nu q_{\nu}}} \alpha_{\nu q_{\nu} j_{\nu q_{\nu}}} \mathcal{Q}_{\nu q_{\nu} j_{\nu q_{\nu}}}(\mathbf{x}_{\nu q_{\nu}i}^1)$, $\mathcal{Q}_q(\mathbf{x}_{qi}^0) = \{\mathcal{Q}_{q1}(\mathbf{x}_{qi}^0), \dots, \mathcal{Q}_{qJ_q}(\mathbf{x}_{qi}^0)\} \top$, $\mathcal{Q}_{\nu q_{\nu}}(\mathbf{x}_{\nu q_{\nu}i}^1) = \{\mathcal{Q}_{\nu q_{\nu}1}(\mathbf{x}_{\nu q_{\nu}i}^1), \dots, \mathcal{Q}_{\nu q_{\nu}J_{\nu q_{\nu}}}(\mathbf{x}_{\nu q_{\nu}i}^1)\} \top$, $\alpha_{0q} = (\alpha_{0q1}, \dots, \alpha_{0qJ_q}) \top$ and $\alpha_{\nu q_{\nu}} = (\alpha_{\nu q_{\nu}1}, \dots, \alpha_{\nu q_{\nu}J_{\nu q_{\nu}}}) \top$. To write equation (6) in a more compact way, we define $\mathcal{Q}_i^0 \top \alpha_0 = \sum_{q=1}^Q \mathcal{Q}_q(\mathbf{x}_{qi}^0) \top \alpha_{0q}$ and $\mathcal{Q}_{\nu i}^1 \top \alpha_{\nu} = \sum_{q_{\nu}=0}^{Q_{\nu}} \mathcal{Q}_{\nu q_{\nu}}(\mathbf{x}_{\nu q_{\nu}i}^1) \top \alpha_{\nu q_{\nu}}$, where $\alpha_0 = (\alpha_{01}, \dots, \alpha_{0Q}) \top$, $\alpha_{\nu} = (\alpha_{\nu 0}, \alpha_{\nu 0}, \dots, \alpha_{\nu Q_{\nu}}) \top$, $\mathcal{Q}_i^0 = \{\mathcal{Q}_1(\mathbf{x}_{1i}^0) \top, \dots, \mathcal{Q}_Q(\mathbf{x}_{Qi}^0) \top\} \top$ and

$\mathbf{Q}_{\nu i}^1 = \{1, \mathbf{Q}_{\nu 0}(\mathbf{x}_{\nu 0 i}^1)^\top, \dots, \mathbf{Q}_{\nu Q_\nu}(\mathbf{x}_{\nu Q_\nu}^1)^\top\}^\top$. Therefore,

$$\xi_{\nu i} = \mathbf{Q}_i^{0\top} \boldsymbol{\alpha}_0 + \mathbf{Q}_{\nu i}^{1\top} \boldsymbol{\alpha}_\nu. \quad (7)$$

If $Q > 0$ then censoring is informative and $\sum_{q=1}^Q s_q(\mathbf{x}_{qi}^0)$ can be estimated using the information from both the censoring and event times. If $Q = 0$ (i.e., the components in $\sum_{k_1=1}^{K_1} s_{1k_1}(\mathbf{x}_{1k_1 i})$ and $\sum_{k_2=1}^{K_2} s_{2k_2}(\mathbf{x}_{2k_2 i})$ are assumed all distinct) then (6) reduces to the model with non-informative censoring and hence we would have

$$\xi_{\nu i} = \gamma_{\nu 0} + \sum_{k_\nu=0}^{K_\nu} \mathbf{Q}_{\nu k_\nu}(\mathbf{x}_{\nu k_\nu i})^\top \boldsymbol{\gamma}_{\nu k_\nu},$$

where $\mathbf{Q}_{\nu k_\nu}(\mathbf{x}_{\nu k_\nu i}) = \{\mathbf{Q}_{\nu k_\nu 1}(\mathbf{x}_{\nu k_\nu i}), \dots, \mathbf{Q}_{\nu k_\nu J_{\nu k_\nu}}(\mathbf{x}_{\nu k_\nu i})\}^\top$ and $\boldsymbol{\gamma}_{\nu k_\nu} = (\gamma_{\nu k_\nu 1}, \dots, \gamma_{\nu k_\nu J_{\nu k_\nu}})^\top$. Furthermore, if $\mathbf{Q}_{\nu i}^\top \boldsymbol{\gamma}_\nu = \sum_{k_\nu=0}^{K_\nu} \mathbf{Q}_{\nu k_\nu}(\mathbf{x}_{\nu k_\nu i})^\top \boldsymbol{\gamma}_{\nu k_\nu}$, $\boldsymbol{\gamma}_\nu = (\gamma_{\nu 0}, \gamma_{\nu 0}, \dots, \gamma_{\nu K_\nu})^\top$ and $\mathbf{Q}_{\nu i} = \{1, \mathbf{Q}_{\nu 0}(\mathbf{x}_{\nu 0 i})^\top, \dots, \mathbf{Q}_{\nu K_\nu}(\mathbf{x}_{\nu K_\nu}^1)^\top\}^\top$, we obtain

$$\xi_{\nu i} = \mathbf{Q}_{\nu i}^\top \boldsymbol{\gamma}_\nu. \quad (8)$$

Note that, for the case in which $Q = 0$, we have introduced the new parameter vector $\boldsymbol{\gamma}_\nu$ to stress the difference between the parameters of the informative and non-informative models. Some methods for determining the value of Q are discussed in Supplementary Material F.

The vectors of parameters $\boldsymbol{\alpha}_{0q}$ and $\boldsymbol{\alpha}_{\nu q_\nu}$ have associated quadratic penalties $\lambda_q \boldsymbol{\alpha}_{0q}^\top \mathcal{D}_q^0 \boldsymbol{\alpha}_{0q}$ and $\lambda_{\nu q_\nu} \boldsymbol{\alpha}_{\nu q_\nu}^\top \mathcal{D}_{\nu q_\nu}^1 \boldsymbol{\alpha}_{\nu q_\nu}$ used in fitting, whose role is to enforce specific properties on the respective functions, such as smoothness. It is important to note that \mathcal{D}_q^0 and $\mathcal{D}_{\nu q_\nu}^1$ only depend on the choice of the basis functions. Smoothing parameter $\lambda_{\nu k_\nu} \in [0, \infty)$ controls the trade-off between fit and smoothness, and plays a crucial role in determining the shape of $\hat{s}_{\nu k_\nu}(\mathbf{x}_{\nu k_\nu i})$. The overall penalty can be defined as $\boldsymbol{\alpha}_\nu^\top \mathcal{D}_\nu \boldsymbol{\alpha}_\nu$, where $\mathcal{D}_\nu = \text{diag}(\lambda_1 \mathcal{D}_1^0, \dots, \lambda_Q \mathcal{D}_Q^0, 0, \lambda_{\nu 0} \mathcal{D}_{\nu 0}^1, \dots, \lambda_{\nu Q_\nu} \mathcal{D}_{\nu Q_\nu}^1)$. Moreover, smooth functions are typically subject to centering (identifiability) constraints. The set up described above allows for several types of covariate effects such as linear, non-linear, spatial, random and functional effects, to name but a few. We refer the reader to Wood (2017) for the exact definitions of the spline bases and penalties of the above mentioned cases.

To give a concrete example, consider the informative additive model

$$g_\nu\{S_\nu(t_{\nu i}|\mathbf{z}_i^0, \mathbf{z}_{\nu i}^1)\} = g_\nu\{S_{\nu 0}(t_{\nu i})\} + \sum_{q=1}^Q \mathcal{Q}_q(\mathbf{z}_{qi}^0)^\top \boldsymbol{\alpha}_{0q} + \sum_{q_\nu=1}^{Q_\nu} \mathcal{Q}_{\nu q_\nu}(\mathbf{z}_{\nu q_\nu i}^1)^\top \boldsymbol{\alpha}_{\nu q_\nu}, \quad (9)$$

where $g_\nu : (0, 1) \rightarrow (-\infty, \infty)$ is a differentiable and invertible link function (see Table 1 in Supplementary Material A), $S_{\nu 0}(t_{\nu i})$ is a baseline survival function, and $g_\nu\{S_{\nu 0}(t_{\nu i})\}$ is represented using a smooth function of time, $s_{\nu 0}(t_{\nu i})$. When the log-log link is chosen, equation (9) yields the proportional hazards model

$$\log\{\mathcal{H}_\nu(t_{\nu i}|\mathbf{z}_i^0, \mathbf{z}_{\nu i}^1)\} = \log\{\mathcal{H}_{\nu 0}(t_{\nu i})\} + \sum_{q=1}^Q \mathcal{Q}_q(\mathbf{z}_{qi}^0)^\top \boldsymbol{\alpha}_{0q} + \sum_{q_\nu=1}^{Q_\nu} \mathcal{Q}_{\nu q_\nu}(\mathbf{z}_{\nu q_\nu i}^1)^\top \boldsymbol{\alpha}_{\nu q_\nu},$$

where $\mathcal{H}_\nu(t_{\nu i}|\mathbf{z}_i^0, \mathbf{z}_{\nu i}^1) = -\log\{S_\nu(t_{\nu i} | \mathbf{z}_i^0, \mathbf{z}_{\nu i}^1)\}$ and $\log\{\mathcal{H}_{\nu 0}(t_{\nu i})\} = -\log\{S_{\nu 0}(t_{\nu i})\}$ is the cumulative baseline hazard function. Analogously, equation (9) yields the proportional odds model when the -logit link is chosen.

The models considered in this paper are fundamentally parametric but flexible. It is worth noting that the more extensive use of parametric survival models in applications has been encouraged by Cox; see the discussion in Reid (1994). Moreover, as pointed out for instance by Hjort (1992), parametric approaches simplify somewhat model estimation and comparison, easily allow for the visualization of the estimated baseline hazard and survival functions, and allow to calculate several quantities of interest and their variances which would otherwise be difficult to obtain with a non-parametric approach. Another important advantage is that there is no necessity to use numerical integration methods to estimate the cumulative hazard function.

2.3 Estimation framework

The data consist of $\{Y_i, \delta_{1i}, \mathbf{z}_i\}$, where $Y_i = \min\{T_{1i}, T_{2i}\}$ and $\delta_{1i} = I(T_{1i} \leq T_{2i})$, for $i = 1, \dots, n$. Let $f(t_1, t_2|\mathbf{z})$ be the conditional joint distribution of (T_1, T_2) given \mathbf{z} . We can write $P(Y_i, \delta_{1i} = 1|\mathbf{z}_i) = \int_{y_i}^{\infty} f(y_i, t_2|\mathbf{z}_i)dt_2$ and $P(Y_i, \delta_{1i} = 0|\mathbf{z}_i) = \int_{y_i}^{\infty} f(t_1, y_i|\mathbf{z}_i)dt_1$. Therefore, the

conditional likelihood function of (Y_i, δ_{1i}) given \mathbf{z}_i , for all $i = 1, \dots, n$, is

$$\mathcal{L} = \prod_{i=1}^n \left\{ \int_{y_i}^{\infty} f(y_i, t_2 | \mathbf{z}_i) dt_2 \right\}^{\delta_{1i}} \left\{ \int_{y_i}^{\infty} f(t_1, y_i | \mathbf{z}_i) dt_1 \right\}^{\delta_{2i}}.$$

Below we provide the relevant details for the cases of informative and non-informative censoring, which highlight the differences between the two estimators and that are also required for the theoretical derivations in Section 2.4.

If it is assumed that T_{1i} and T_{2i} are conditionally independent given \mathbf{z}_i , then $\int_{y_i}^{\infty} f(y_i, t_2 | \mathbf{z}_i) dt_2 = f_1(y_i | \mathbf{z}_{1i}; \gamma_1) S_2(y_i | \mathbf{z}_{2i}; \gamma_2)$ and $\int_{y_i}^{\infty} f(t_1, y_i | \mathbf{z}_i) dt_1 = f_2(y_i | \mathbf{z}_{2i}; \gamma_2) S_1(y_i | \mathbf{z}_{1i}; \gamma_1)$ when censoring is non-informative. However, if censoring is informative γ_1 and γ_2 would have some components in common. Since it was assumed that the first Q components of γ_1 are the same as the first Q components of γ_2 , we have $\mathbf{Q}_{\nu i}^{\top} \gamma_{\nu} = \mathbf{Q}_i^{0\top} \alpha_0 + \mathbf{Q}_{\nu i}^{1\top} \alpha_{\nu}$. Using (1), (2), (7), (8), and $\xi_{\nu i}(\gamma_{\nu})$ and $\xi_{\nu i}(\alpha_0, \alpha_{\nu})$ as the shorthand notation for $\xi_{\nu i}(y_i, \mathbf{z}_{\nu i}; \gamma_{\nu})$ and $\xi_{\nu i}(y_i, \mathbf{z}_{\nu i}; \alpha_0, \alpha_{\nu})$, the non-informative and informative log-likelihood functions can be written, respectively, as

$$\begin{aligned} \ell(\boldsymbol{\gamma}) &= \sum_{i=1}^n \left\{ \log \mathcal{G}_1 [\xi_{1i}(\boldsymbol{\gamma}_1)] + \delta_{1i} \log \left\{ -\frac{\mathcal{G}'_1 [\xi_{1i}(\boldsymbol{\gamma}_1)]}{\mathcal{G}_1 [\xi_{1i}(\boldsymbol{\gamma}_1)]} \frac{\partial \xi_{1i}(\boldsymbol{\gamma}_1)}{\partial y_i} \right\} \right\} \\ &\quad + \sum_{i=1}^n \left\{ \log \mathcal{G}_2 [\xi_{2i}(\boldsymbol{\gamma}_2)] + \delta_{2i} \log \left\{ -\frac{\mathcal{G}'_2 [\xi_{2i}(\boldsymbol{\gamma}_2)]}{\mathcal{G}_2 [\xi_{2i}(\boldsymbol{\gamma}_2)]} \frac{\partial \xi_{2i}(\boldsymbol{\gamma}_2)}{\partial y_i} \right\} \right\}, \\ \ell(\boldsymbol{\alpha}) &= \sum_{i=1}^n \left\{ \log \mathcal{G}_1 [\xi_{1i}(\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_1)] + \delta_{1i} \log \left\{ -\frac{\mathcal{G}'_1 [\xi_{1i}(\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_1)]}{\mathcal{G}_1 [\xi_{1i}(\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_1)]} \frac{\partial \xi_{1i}(\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_1)}{\partial y_i} \right\} \right\} \\ &\quad + \sum_{i=1}^n \left\{ \log \mathcal{G}_2 [\xi_{2i}(\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_2)] + \delta_{2i} \log \left\{ -\frac{\mathcal{G}'_2 [\xi_{2i}(\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_2)]}{\mathcal{G}_2 [\xi_{2i}(\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_2)]} \frac{\partial \xi_{2i}(\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_2)}{\partial y_i} \right\} \right\}. \end{aligned} \tag{10}$$

To ensure that the hazard functions in (10) are positive, $[\partial \xi_{\nu i}(\boldsymbol{\theta}_{\nu}) / \partial y_i]$, for $\nu = (1, 2)$, must be positive. To this end, we model the time effects using B-splines with coefficients constrained such that the resulting smooth functions of time are monotonically increasing. In particular, we define $s_{\nu 0}(y_i) = \sum_{j_{\nu}=1}^{J_{\nu}} \vartheta_{\nu 0 j_{\nu}} \mathcal{M}_{\nu 0 j_{\nu}}(y_i)$, where the $\mathcal{M}_{\nu 0 j_{\nu}}(y_i)$ are B-spline basis functions of at least second order built over the interval $[a, b]$, based on equally spaced knots, and $\vartheta_{\nu 0 j_{\nu}}$ are spline coefficients. A sufficient condition for $[\partial s_{\nu 0}(y_i) / \partial y_i] \geq 0$ over $[a, b]$ is that $\vartheta_{\nu 0 j_{\nu}} \geq \vartheta_{\nu 0 j-1_{\nu 0}}, \forall j$ (Leitenstorfer & Tutz, 2006). Such condition can be imposed by re-parametrising the spline coefficient vector so that $\boldsymbol{\vartheta}_{\nu 0} = \boldsymbol{\Gamma}_{\nu 0} \tilde{\boldsymbol{\beta}}_{\nu 0}$, where $\boldsymbol{\beta}_{\nu 0}^{\top} = (\beta_{\nu 01}, \beta_{\nu 02}, \dots, \beta_{\nu 0 J_{\nu}})$, $\tilde{\boldsymbol{\beta}}_{\nu}^{\top} =$

$(\beta_{\nu 01}, \exp(\beta_{\nu 02}), \dots, \exp(\beta_{\nu 0J_{\nu 0}}))$ and $\Gamma_{\nu 0}[\kappa_{\nu 01}, \kappa_{\nu 02}] = 0$ if $\kappa_{\nu 01} < \kappa_{\nu 02}$, and $\Gamma_{\nu 0}[\kappa_{\nu 01}, \kappa_{\nu 02}] = 1$ if $\kappa_{\nu 01} \geq \kappa_{\nu 02}$. Following Pya & Wood (2015, Section 2.2.1), the penalty term is set up to penalise the squared differences between adjacent $\beta_{\nu 0j_{\nu 0}}$, starting from $\beta_{\nu 02}$, using $\mathcal{D}_{\nu 0} = \mathcal{D}_{\nu 0}^{\circ\top} \mathcal{D}_{\nu 0}^{\circ}$, where $\mathcal{D}_{\nu 0}^{\circ}$ is a $(J_{\nu 0} - 2) \times J_{\nu 0}$ matrix made up of zeros except that $\mathcal{D}_{\nu 0}^{\circ}[\kappa_{\nu 0}, \kappa_{\nu 0} + 1] = -\mathcal{D}_{\nu 0}^{\circ}[\kappa_{\nu 0}, \kappa_{\nu 0} + 2] = 1$ for $\kappa_{\nu 0} = 1, \dots, J_{\nu 0} - 2$. Therefore, the non-informative and informative additive predictors, that ensure positive hazard functions in (10), are

$$\begin{aligned}
 \xi_{\nu i} &= \gamma_{\nu 0} + \mathcal{Q}_{\nu 0}(y_i)^{\top} \Gamma_{\nu 0} \tilde{\gamma}_{\nu 0} + \sum_{k_{\nu}=1}^{K_{\nu}} \mathcal{Q}_{\nu k_{\nu}}(\mathbf{x}_{\nu k_{\nu} i})^{\top} \gamma_{\nu k_{\nu}}, \\
 \xi_{\nu i} &= \alpha_{\nu 0} + \mathcal{Q}_{\nu 0}(y_i)^{\top} \Gamma_{\nu 0} \tilde{\alpha}_{\nu 0} + \sum_{q=1}^Q \mathcal{Q}_q(\mathbf{x}_{qi}^0)^{\top} \alpha_{0q} + \sum_{q_{\nu}=1}^{Q_{\nu}} \mathcal{Q}_{\nu q_{\nu}}(\mathbf{x}_{\nu q_{\nu} i}^1)^{\top} \alpha_{\nu q_{\nu}}.
 \end{aligned} \tag{11}$$

Our model specification allows for a high degree of flexibility in modelling survival data. If an unpenalised estimation approach is employed to estimate $\boldsymbol{\gamma} = (\boldsymbol{\gamma}_1^{\top}, \boldsymbol{\gamma}_2^{\top})^{\top}$ and $\boldsymbol{\alpha} = (\boldsymbol{\alpha}_0^{\top}, \boldsymbol{\alpha}_1^{\top}, \boldsymbol{\alpha}_2^{\top})^{\top}$, then the resulting smooth function estimates are likely to be unduly wiggly (e.g., Wood, 2017). Therefore, to prevent over-fitting, the following functions are maximized

$$\ell_p(\boldsymbol{\gamma}) = \ell(\boldsymbol{\gamma}) - \frac{1}{2} \boldsymbol{\gamma}^{\top} \boldsymbol{S} \boldsymbol{\gamma}, \tag{12}$$

$$\ell_p(\boldsymbol{\alpha}) = \ell(\boldsymbol{\alpha}) - \frac{1}{2} \boldsymbol{\alpha}^{\top} \boldsymbol{S} \boldsymbol{\alpha}, \tag{13}$$

where $\ell_p(\boldsymbol{\gamma})$ and $\ell_p(\boldsymbol{\alpha})$ are the non-informative and informative penalized log-likelihoods. Moreover, $\boldsymbol{S} = \text{diag}(\boldsymbol{D}_1, \boldsymbol{D}_2)$, and \boldsymbol{D}_1 and \boldsymbol{D}_2 are overall penalties which contain $\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2$. The smoothing parameter vectors can be collected in the overall vector $\boldsymbol{\lambda} = (\boldsymbol{\lambda}_1^{\top}, \boldsymbol{\lambda}_2^{\top})^{\top}$. Estimation of the models' parameters and smoothing coefficients is achieved by using a stable and efficient trust region algorithm with integrated automatic multiple smoothing parameter selection (see Supplementary Material C for details). This required working with first and second order analytical derivatives which have been tediously derived as well as verified using numerical derivatives. Their structures are shown below. Note that these results were also required for the theoretical proofs presented in Section 2.4.

When censoring is non-informative, the gradient of (12) can be obtained as

$$\nabla_{\gamma} \ell_p(\gamma) = \nabla_{\gamma} \ell(\gamma) - \gamma \mathcal{S},$$

where $\nabla_{\gamma} \ell(\gamma) = (\nabla_{\gamma_1} \ell(\gamma)^\top, \nabla_{\gamma_2} \ell(\gamma)^\top)^\top$. The components of $\nabla_{\gamma} \ell(\gamma)$ can generically be calculated using the following expression

$$\nabla_{\gamma_{\nu k_\nu}} \ell(\gamma) = \begin{cases} \sum_{i=1}^n [\Delta_{\nu i} \mathcal{Q}_{\nu 0}^\Delta(y_i) + \Omega_{\nu i} \mathcal{Q}_{\nu 0}^{\Delta'}(y_i)] & \text{if } \gamma_{\nu k_\nu} = \gamma_{\nu 0}, \\ \sum_{i=1}^n [\Delta_{\nu i} \mathcal{Q}_{\nu k_\nu}(\mathbf{x}_{\nu k_\nu i})] & \text{otherwise.} \end{cases} \quad (14)$$

In (14), $\mathcal{Q}_{\nu 0}^\Delta(y_i)$ and $\mathcal{Q}_{\nu 0}^{\Delta'}(y_i)$ are design vectors. Furthermore, $\Omega_{\nu i} = \delta_{\nu i} \left(\frac{\partial \xi_{\nu i}}{\partial y_i} \right)^{-1}$ and $\Delta_{\nu i} = \left[\frac{\mathcal{G}'_\nu}{\mathcal{G}_\nu} + \delta_{\nu i} \left(\frac{\mathcal{G}''_\nu}{\mathcal{G}'_\nu} - \frac{\mathcal{G}'_\nu}{\mathcal{G}_\nu} \right) \right]$, for all $\nu = 1, 2$. The non-informative penalized Hessian can be calculated as

$$\nabla_{\gamma\gamma} \ell_p(\gamma) = \nabla_{\gamma\gamma} \ell(\gamma) - \mathcal{S},$$

where

$$\nabla_{\gamma\gamma} \ell(\gamma) = \begin{bmatrix} \nabla_{\gamma_1 \gamma_1} \ell(\gamma) & \mathbf{0} \\ \mathbf{0} & \nabla_{\gamma_2 \gamma_2} \ell(\gamma) \end{bmatrix}.$$

Further, the elements of $\nabla_{\gamma_\nu \gamma_\nu} \ell(\gamma)$ are calculated using

$$\begin{aligned} \nabla_{\gamma_{\nu k_\nu} \gamma_{\nu 0}} \ell(\gamma) &= \sum_{i=1}^n [\mathcal{Q}_{\nu k_\nu}(\mathbf{x}_{\nu k_\nu i}) \Phi_{\nu i} \mathcal{Q}_{\nu 0}^\Delta(y_i)^\top], \\ \nabla_{\gamma_{\nu 0} \gamma_{\nu s_\nu}} \ell(\gamma) &= \sum_{i=1}^n [\mathcal{Q}_{\nu 0}^\Delta(y_i) \Phi_{\nu i} \mathcal{Q}_{\nu s_\nu}(\mathbf{x}_{\nu s_\nu i})^\top], \\ \nabla_{\gamma_{\nu k_\nu} \gamma_{\nu s_\nu}} \ell(\gamma) &= \sum_{i=1}^n [\mathcal{Q}_{\nu k_\nu}(\mathbf{x}_{\nu k_\nu i}) \Phi_{\nu i} \mathcal{Q}_{\nu s_\nu}(\mathbf{x}_{\nu s_\nu i})^\top], \\ \nabla_{\gamma_{\nu 0} \gamma_{\nu 0}} \ell(\gamma) &= \sum_{i=1}^n [\mathcal{Q}_{\nu 0}^\Delta(y_i) \Phi_{\nu i} \mathcal{Q}_{\nu 0}^\Delta(y_i)^\top + \Delta_{\nu i} \mathcal{Q}_{\nu 0}^{\Delta\Delta}(y_i) - \mathcal{Q}_{\nu 0}^{\Delta'}(y_i) \Psi_{\nu i} \mathcal{Q}_{\nu 0}^{\Delta'}(y_i)^\top + \Omega_{\nu i} \mathcal{Q}_{\nu 0}^{\Delta\Delta'}(y_i)]. \end{aligned} \quad (15)$$

In these sub-matrices $\Phi_{\nu i} = \delta_{\nu i} \left(\frac{\mathcal{G}_{\nu}'''}{\mathcal{G}_{\nu}} - \frac{\mathcal{G}_{\nu}''^2}{\mathcal{G}_{\nu}'} - \frac{\mathcal{G}_{\nu}''}{\mathcal{G}_{\nu}} + \frac{\mathcal{G}_{\nu}'^2}{\mathcal{G}_{\nu}^2} \right)$ and $\Psi_{\nu i} = \left[\delta_{\nu i} \left(\frac{\partial \xi_{\nu i}}{\partial y_i} \right)^{-2} \right]$. In addition, $\mathcal{Q}_{\nu 0}^{\Delta \Delta}(y_i)$ and $\mathcal{Q}_{\nu 0}^{\Delta \Delta'}(y_i)$ are design diagonal matrices.

If the censoring is informative, the gradient of (13) can be calculated as

$$\nabla_{\alpha} \ell_p(\alpha) = \nabla_{\alpha} \ell(\alpha) - \alpha \mathcal{S},$$

where $\nabla_{\alpha} \ell(\alpha) = (\nabla_{\alpha_0} \ell(\alpha)^{\top}, \nabla_{\alpha_1} \ell(\alpha)^{\top}, \nabla_{\alpha_2} \ell(\alpha)^{\top})^{\top}$. To obtain $\nabla_{\alpha_0} \ell(\alpha)$ and $\nabla_{\alpha_{\nu}} \ell(\alpha)$, we use

$$\begin{aligned} \nabla_{\alpha_0} \ell(\alpha) &= \sum_{i=1}^n [\mathcal{Q}_i^0 (\Delta_{1i} + \Delta_{2i})], \\ \nabla_{\alpha_{\nu k_{\nu}}} \ell(\alpha) &= \begin{cases} \sum_{i=1}^n [\Delta_{\nu i} \mathcal{Q}_{\nu 0}^{i\Delta}(y_i) + \Omega_{\nu i} \mathcal{Q}_{\nu 0}^{i\Delta'}(y_i)] & \text{if } \alpha_{\nu k_{\nu}} = \alpha_{\nu 0}, \\ \sum_{i=1}^n [\Delta_{\nu i} \mathcal{Q}_{\nu k_{\nu}}(\mathbf{x}_{\nu k_{\nu} i})] & \text{otherwise,} \end{cases} \end{aligned} \quad (16)$$

where $\mathcal{Q}_{\nu 0}^{i\Delta}(y_i)$ and $\mathcal{Q}_{\nu 0}^{i\Delta'}(y_i)$ are design vectors. The informative penalized Hessian can be obtained as follow

$$\nabla_{\alpha \alpha} \ell_p(\alpha) = \nabla_{\alpha \alpha} \ell(\alpha) - \mathcal{S},$$

where

$$\nabla_{\alpha \alpha} \ell(\alpha) = \begin{bmatrix} \nabla_{\alpha_0 \alpha_0} \ell(\alpha) & \nabla_{\alpha_0 \alpha_1} \ell(\alpha) & \nabla_{\alpha_0 \alpha_2} \ell(\alpha) \\ \nabla_{\alpha_1 \alpha_0} \ell(\alpha) & \nabla_{\alpha_1 \alpha_1} \ell(\alpha) & \mathbf{0} \\ \nabla_{\alpha_2 \alpha_0} \ell(\alpha) & \mathbf{0} & \nabla_{\alpha_2 \alpha_2} \ell(\alpha) \end{bmatrix}.$$

Furthermore, $\nabla_{\alpha_0\alpha_0}\ell(\alpha)$ and the components of $\nabla_{\alpha_\nu\alpha_0}\ell(\alpha)$ and $\nabla_{\alpha_0\alpha_\nu}\ell(\alpha)$ are obtained using

$$\begin{aligned}\nabla_{\alpha_0\alpha_0}\ell(\alpha) &= \sum_{i=1}^n \left[\mathcal{Q}_i^0(\Phi_{1i} + \Phi_{2i})\mathcal{Q}_i^{0\top} \right], \\ \nabla_{\alpha_0\alpha_{\nu q_\nu}}\ell(\alpha) &= \begin{cases} \sum_{i=1}^n \left[\mathcal{Q}_i^0\Phi_{\nu i}\mathcal{Q}_{\nu 0}^{\prime\Delta}(y_i)^\top \right] & \text{if } \alpha_{\nu q_\nu} = \alpha_{\nu 0}, \\ \sum_{i=1}^n \left[\mathcal{Q}_i^0\Phi_{\nu i}\mathcal{Q}_{\nu q_\nu}(\mathbf{x}_{\nu q_\nu i}^1)^\top \right] & \text{otherwise,} \end{cases} \\ \nabla_{\alpha_{\nu q_\nu}\alpha_0}\ell(\alpha) &= \begin{cases} \sum_{i=1}^n \left[\mathcal{Q}_{\nu 0}^{\prime\Delta}(y_i)\Phi_{\nu i}\mathcal{Q}_i^{0\top} \right] & \text{if } \alpha_{\nu q_\nu} = \alpha_{\nu 0}, \\ \sum_{i=1}^n \left[\mathcal{Q}_{\nu q_\nu}(\mathbf{x}_{\nu q_\nu i}^1)\Phi_{\nu i}\mathcal{Q}_i^{0\top} \right] & \text{otherwise.} \end{cases}\end{aligned}\quad (17)$$

Finally, the elements of $\nabla_{\alpha_\nu\alpha_\nu}\ell(\alpha)$ are calculated using

$$\begin{aligned}\nabla_{\alpha_{\nu q_\nu}\alpha_{\nu 0}}\ell(\alpha) &= \sum_{i=1}^n \left[\mathcal{Q}_{\nu q_\nu}(\mathbf{x}_{\nu q_\nu i}^1)\Phi_{\nu i}\mathcal{Q}_{\nu 0}^{\prime\Delta}(y_i)^\top \right], \\ \nabla_{\alpha_{\nu 0}\alpha_{\nu q_\nu}}\ell(\alpha) &= \sum_{i=1}^n \left[\mathcal{Q}_{\nu 0}^{\prime\Delta}(y_i)\Phi_{\nu i}\mathcal{Q}_{\nu q_\nu}(\mathbf{x}_{\nu q_\nu i}^1)^\top \right], \\ \nabla_{\alpha_{\nu q_\nu}\alpha_{\nu r_\nu}}\ell(\alpha) &= \sum_{i=1}^n \left[\mathcal{Q}_{\nu q_\nu}(\mathbf{x}_{\nu q_\nu i}^1)\Phi_{\nu i}\mathcal{Q}_{\nu r_\nu}(\mathbf{x}_{\nu r_\nu i}^1)^\top \right], \\ \nabla_{\alpha_{\nu 0}\alpha_{\nu 0}}\ell(\alpha) &= \sum_{i=1}^n \left[\mathcal{Q}_{\nu 0}^{\prime\Delta}(y_i)\Phi_{\nu i}\mathcal{Q}_{\nu 0}^{\prime\Delta}(y_i)^\top + \Delta_{\nu i}\mathcal{Q}_{\nu 0}^{\prime\Delta\Delta}(y_i) - \mathcal{Q}_{\nu 0}^{\prime\Delta'}(y_i)\Psi_{\nu i}\mathcal{Q}_{\nu 0}^{\prime\Delta'}(y_i)^\top + \Omega_{\nu i}\mathcal{Q}_{\nu 0}^{\prime\Delta\Delta'}(y_i) \right].\end{aligned}\quad (18)$$

As before, $\mathcal{Q}_{\nu 0}^{\prime\Delta\Delta}(y_i)$ and $\mathcal{Q}_{\nu 0}^{\prime\Delta\Delta'}(y_i)$ represent design diagonal matrices.

The derivations of the results reported here as well as some algorithmic details are given in Supplementary Materials B and C.

Remark 1. The scores and Hessian components described in this section have been implemented in a modular way, hence no substantial programming work will be required to incorporate link functions not considered in this article. Furthermore, quantities such as those defined in (14), (15), (16), (17) and (18), are needed for the theoretical proofs of the next section.

2.4 Theoretical properties

In this section, we derive the \sqrt{n} consistency and asymptotic normality of the non-informative and informative estimators, and shed light on the efficiency gains produced by the newly introduced informative estimator when compared to its non-informative counterpart. As far as the number of basis functions is concerned, we use the fixed-knot asymptotic framework since it is closer to practical statistical modelling (e.g., Vatter & Chavez-Demoulin, 2015, and references therein). In what follows, we define $\hat{S}_{\nu 0}(\hat{\boldsymbol{\theta}}_{\nu 0}) = \mathcal{G}_{\nu 0}[s(\hat{\boldsymbol{\theta}}_{\nu 0})]$ as the short notation for $\hat{S}_{\nu 0}(y_i, \hat{\boldsymbol{\theta}}_{\nu 0}) = \mathcal{G}_{\nu 0}[s(y_i, \hat{\boldsymbol{\theta}}_{\nu 0})]$ and $\boldsymbol{\theta}^*$ as the true vector of parameters.

The informative penalized maximum log-likelihood estimator (IPMLE) can be defined as

$$\hat{\boldsymbol{\alpha}} = \operatorname{argmax}_{\boldsymbol{\alpha} \in \Theta} \ell_p(\boldsymbol{\alpha}),$$

and the non-informative counterpart (NPMLE) as

$$\hat{\boldsymbol{\gamma}} = \operatorname{argmax}_{\boldsymbol{\gamma} \in \Theta} \ell_p(\boldsymbol{\gamma}).$$

Theorem 1 (Asymptotic properties of the IPMLE estimator). If the set of Assumptions 1 and 2 in Supplementary Material D hold then

- (i) the informative penalized maximum log-likelihood estimator $\hat{\boldsymbol{\alpha}}$ exists, is \sqrt{n} -consistent and

$$\sqrt{n}(\hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}^*) \xrightarrow{d} \mathcal{N}\{\mathbf{0}, [\mathcal{I}(\boldsymbol{\alpha}^*)]^{-1}\},$$

where $\mathcal{I}(\boldsymbol{\alpha}^*) = \mathbb{E}[-\nabla_{\boldsymbol{\alpha}\boldsymbol{\alpha}} \ell(\mathbf{w}; \boldsymbol{\alpha}^*)]$.

- (ii) $\hat{S}_{10}(\hat{\boldsymbol{\alpha}}_{10})$ is asymptotically independent of $\hat{S}_{20}(\hat{\boldsymbol{\alpha}}_{20})$ and

$$\sqrt{n}[\hat{S}_{\nu 0}(\hat{\boldsymbol{\alpha}}_{\nu 0}) - S_{\nu 0}(\boldsymbol{\alpha}_{\nu 0}^*)] \xrightarrow{d} \mathcal{N}\{\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\alpha}_{\nu 0}^*}\}, \nu = 1, 2,$$

where $\boldsymbol{\Sigma}_{\boldsymbol{\alpha}_{\nu 0}^*} = \mathcal{G}'_{\nu 0}[s(\boldsymbol{\alpha}_{\nu 0}^*)] \nabla_{\boldsymbol{\alpha}_{\nu 0}} s(\boldsymbol{\alpha}_{\nu 0}^*) [\mathcal{I}(\boldsymbol{\alpha}_{\nu 0}^*)]^{-1} \nabla_{\boldsymbol{\alpha}_{\nu 0}} s(\boldsymbol{\alpha}_{\nu 0}^*)^\top \mathcal{G}'_{\nu 0}[s(\boldsymbol{\alpha}_{\nu 0}^*)]$ and $\mathcal{I}(\boldsymbol{\alpha}_{\nu 0}^*) = \mathbb{E}[-\nabla_{\boldsymbol{\alpha}_{\nu 0} \boldsymbol{\alpha}_{\nu 0}} \ell(\mathbf{w}; \boldsymbol{\alpha}_{\nu 0}^*)]$.

Theorem 2 (Asymptotic properties of the NPMLE estimator). If the set of Assumptions **1** and **2** in Supplementary Material D hold then

- (i) the non-informative penalized maximum log-likelihood estimator $\hat{\gamma}$ exists, is \sqrt{n} -consistent and

$$\sqrt{n}(\hat{\gamma} - \gamma^*) \xrightarrow{d} \mathcal{N}\{\mathbf{0}, [\mathcal{I}(\gamma^*)]^{-1}\},$$

where $\mathcal{I}(\gamma^*) = \mathbb{E}[-\nabla_{\gamma\gamma}\ell(\mathbf{w}; \gamma^*)]$.

- (ii) $\hat{S}_{10}(\hat{\gamma}_{10})$ is asymptotically independent of $\hat{S}_{20}(\hat{\gamma}_{20})$ and

$$\sqrt{n}[\hat{S}_{\nu 0}(\hat{\gamma}_{\nu 0}) - S_{\nu 0}(\gamma_{\nu 0}^*)] \xrightarrow{d} \mathcal{N}\{\mathbf{0}, \Sigma_{\gamma_{\nu 0}^*}\}, \nu = 1, 2,$$

where $\Sigma_{\gamma_{\nu 0}^*} = \mathcal{G}'_{\nu 0}[s(\gamma_{\nu 0}^*)] \nabla_{\gamma_{\nu 0}} s(\gamma_{\nu 0}^*) [\mathcal{I}(\gamma_{\nu 0}^*)]^{-1} \nabla_{\gamma_{\nu 0}} s(\gamma_{\nu 0}^*)^\top \mathcal{G}'_{\nu 0}[s(\gamma_{\nu 0}^*)]$ and $\mathcal{I}(\gamma_{\nu 0}^*) = \mathbb{E}[-\nabla_{\gamma_{\nu 0}\gamma_{\nu 0}}\ell(\mathbf{w}; \gamma_{\nu 0}^*)]$.

Theorem 3 (Efficiency of the IPMLE estimator). For $\nu = 1, 2$, let $\gamma_\nu = (\gamma_\nu^l, \gamma_\nu^{nl})^\top$ be the informative and non-informative parameters of the non-informative model, respectively. Under the set of Assumptions **1** and **2** in Supplementary Material D, and if we further assume that $\gamma_{\nu 0}^{nl} = \alpha_{\nu 0}$, then

$$\mathcal{ACov}(\hat{\alpha}_0) < \mathcal{ACov}(\hat{\gamma}_\nu^l),$$

$$\mathcal{ACov}(\hat{\alpha}_\nu) < \mathcal{ACov}(\hat{\gamma}_\nu^{nl}),$$

where $\mathcal{ACov}(\hat{\alpha}_0) = \Sigma_{\alpha_0^*}$, $\mathcal{ACov}(\hat{\alpha}_\nu) = \Sigma_{\alpha_\nu^*}$, $\mathcal{ACov}(\hat{\gamma}_\nu^l) = \Sigma_{\gamma_{\nu 0}^{*l}}$, and $\mathcal{ACov}(\hat{\gamma}_\nu^{nl}) = \Sigma_{\gamma_{\nu 0}^{*nl}}$ represent the asymptotic covariance matrices of $\hat{\alpha}_0$, $\hat{\alpha}_\nu$, $\hat{\gamma}_\nu^l$ and $\hat{\gamma}_\nu^{nl}$ respectively.

The proofs of Theorems 1, 2 and 3 are given in Supplementary Material D.

Remark 2. The fact that the informative and non-informative survival functions are orthogonal (part (ii) of Theorems 1 and 2) suggests that the estimation algorithm will yield more accurate parameter vector updates throughout the iterations (e.g., Nocedal & Wright, 2006). Moreover, Theorem 3 shows that under informative censoring it is possible to estimate the model's coefficients more efficiently since more information is exploited by the informative model.

Remark 3. As far as the construction of confidence intervals and p-values are concerned, for practical purposes it is convenient to adapt to the current context the results discussed in Marra et al. (2017). Supplementary Material E provides more details on this.

3 Simulation study

This section provides evidence on the empirical effectiveness of the proposed methodology in recovering true linear effects, non linear effects and baseline functions under informative censoring for three Data Generating Processes (DGPs). The performance of the informative penalized maximum log-likelihood estimator against that of its non-informative counterpart was also examined.

- (i) DGP1 (z_{1i} non-informative, z_{2i} informative and censoring rate of about 78%). Event times, T_{1i} , were generated from a proportional hazard model, while censored times, T_{2i} , were generated from a proportional odd model. These, defined on the survival function scale, are given by

$$\begin{aligned} & \log [-\log \{S_{10}(t_{1i})\}] + \alpha_{01} + \alpha_{11}z_{1i} + s_{11}(z_{2i}), \\ & \log \left[\frac{\{1 - S_{20}(t_{2i})\}}{S_{20}(t_{2i})} \right] + \alpha_{02} + \alpha_{12}z_{1i} + s_{12}(z_{2i}), \end{aligned} \quad (19)$$

where $S_{10}(t_{1i}) = 0.72 \exp(-0.4t_{1i}^{2.4}) + 0.28 \exp(-0.1t_{1i}^{1.0})$ and $S_{20}(t_{2i}) = 0.99 \exp(-0.1t_{2i}^{2.2}) + 0.01 \exp(-0.4t_{2i}^{1.1})$ (Crowther & Lambert, 2013). Covariate z_{1i} was generated using a binomial distribution and z_{2i} using a uniform distribution. As for the smooth functions, we used $s_{11}(z_{2i}) = s_{12}(z_{2i}) = -0.2 \exp(3.2z_{2i})$, whereas the parametric coefficients were: $\alpha_{01} = 0.25$, $\alpha_{02} = 0.85$, $\alpha_{11} = -2.0$ and $\alpha_{12} = 1.8$.

Sample sizes were set to 500, 1000 and 4000, and the number of replicates to 1000. Replicates in which the models did not converge were discarded and replaced with additional ones. The models were fitted using `gamLSS()` in GJRM by employing the proportional hazard link ("PH") for the event times and the proportional odd link ("PO") for the censoring times (see Supplementary Material A for some software details). The smooth components of z_2 were represented using penalized low rank thin plate splines with second order penalty

and 10 bases (the default in GJRM), and the smooths of times using monotonic penalized B-splines with penalty defined in Section 2.3 and 10 bases. Note that smooth terms of explanatory variables can also be represented using different spline definitions (see Supplementary Material A). In the case of one-dimensional smooth functions, all definitions lead to virtually the same result as long as the amount of smoothing is selected in a data-driven manner (e.g., Wood, 2017). For each replicate, curve estimates were constructed using 200 equally spaced fixed values in the $(0, 8)$ range for the baseline functions and $(0, 1)$ otherwise.

Results: Regarding the estimates for α_{11} (the parameter of the non-informative covariate), Figure 4 (in Supplementary Material G) and Table 1 show that overall the mean estimates for the IPMLE and NPMLE are very close to the respective true values and improve as the sample size increases, and that the variability of the estimates decreases as the sample size grows large.

As for the smooth effect of the informative covariate, Figures 6 and 7 (in Supplementary Material G), and Table 1 show that overall the true functions are recovered well by the proposed estimation methods and that the results improve in terms of bias and efficiency as the sample size increases. However, the IPMLE is more efficient than the NPMLE for all sample sizes examined in the simulation study; for example, for $n = 500, 1000$ the RMSE for the NPMLE is more than twice as large as the IPMLE. Some gains in efficiency are also observed for the baseline functions.

- (ii) DGP2 (z_{1i} informative, z_{2i} informative and censoring rate of about 74%). As for DGP1, T_{1i} and T_{2i} were generated using the model defined in (19). However, in this case, the baseline survival functions were defined as $S_{10}(t_{1i}) = 0.75 \exp(-0.4t_{1i}^{2.4}) + 0.25 \exp(-0.1t_{1i}^{1.0})$ and $S_{20}(t_{2i}) = 0.99 \exp(-0.1t_{2i}^{2.2}) + 0.01 \exp(-0.4t_{2i}^{1.1})$. The informative covariates, z_{1i} and z_{2i} , were generated using binomial and uniform distributions, respectively. Finally, $s_{11}(z_{2i}) = s_{12}(z_{2i}) = -0.2 \exp(3.2z_{2i})$, $\alpha_{01} = 0.25$, $\alpha_{02} = 0.85$ and $\alpha_{11} = \alpha_{12} = -1.5$.

Results: Similarly to DGP1, Figures 5, 8 and 9 and Table 3 (in Supplementary Material G) show that overall the mean estimates for the two estimators are very close to the respective true values and improve as the sample size increases. The variability of the estimates also

(a) Informative Penalized Maximum Log-likelihood Estimator (IPMLE)						
	Bias			RMSE		
	n = 500	n = 1000	n = 4000	n = 500	n = 1000	n = 4000
α_{11}	-0.047	-0.013	-0.001	0.369	0.239	0.118
s_{11}	0.036	0.028	0.013	0.161	0.114	0.061
h_{10}	0.095	0.069	0.034	0.336	0.245	0.104
S_{10}	0.027	0.024	0.018	0.071	0.054	0.033

(b) Non-informative Penalized Maximum Log-likelihood Estimator (NPMLE)						
	Bias			RMSE		
	n = 500	n = 1000	n = 4000	n = 500	n = 1000	n = 4000
α_{11}	-0.079	-0.015	-0.005	0.360	0.245	0.116
s_{11}	0.085	0.069	0.046	0.383	0.206	0.118
h_{10}	0.120	0.070	0.034	0.427	0.292	0.121
S_{10}	0.034	0.025	0.017	0.086	0.068	0.039

Table 1: Bias and root mean squared error (RMSE) for the IPMLE and NPMLE obtained by applying the `gamLSS()` to informative survival data simulated according to DGP1 characterised by a censoring rate of about 78%. Bias and RMSE for the smooth terms are calculated, respectively, as $n_s^{-1} \sum_{i=1}^{n_s} |\bar{s}_i - s_i|$ and $n_s^{-1} \sum_{i=1}^{n_s} \sqrt{n_{rep}^{-1} \sum_{rep=1}^{n_{rep}} (\hat{s}_{rep,i} - s_i)^2}$, where $\bar{s}_i = n_{rep}^{-1} \sum_{rep=1}^{n_{rep}} \hat{s}_{rep,i}$, n_s is the number of equally spaced fixed values in the (0, 8) or (0, 1) range, and n_{rep} is the number of simulation replicates. In this case, $n_s = 200$ and $n_{rep} = 1000$. The bias for the smooth terms is based on absolute differences in order to avoid compensating effects when taking the sum.

decreases as the sample size grows large. However, the IPMLE is significantly more efficient than the NPMLE for all cases considered.

Computing times for the proposed approach were on average 8 seconds for $n = 4000$ and around 5 seconds for $n = 1000, 500$. A third DGP with a different smooth function for z_{2i} and with a censoring rate of about 47% was explored (see Supplementary Material G). This DGP suggested the perhaps expected result that the gain in efficiency of the IPMLE tends not to be too significant when a mild censoring rate is considered. Finally, for the above DGPs, we explored the ability of information criteria such as the Akaike information criterion (AIC) and the Bayesian information criterion (BIC), defined in Supplementary Material F, to select the correct model. When doing this, we also considered the informative estimator with incorrectly chosen set of informative covariates (e.g., for DGP1, in estimation, z_1 was assumed to be informative instead of z_2). For all sample sizes and cases considered both AIC and BIC always chose the correct model.

4 Empirical illustration

The modelling framework is illustrated using the data employed by Lu & Zhang (2012), where the aim was to assess how several factors affect the contraction of pneumonia in infants in the presence of informative censoring. According to the World Health Organization (WHO), pneumonia accounted for 16% of all deaths of children under five years old in 2015. The data set consists of 3470 annual personal interviews conducted for the National Longitudinal Survey of Youth from 1979 through 1986 (NLSY, 1995). The response variable, Y_i , is the age, in months, at which the infant was hospitalised for pneumonia, and 97.9% of this variable is right censored.

The covariates considered in the modeling were age of the mother in years (`mothage`), urban environment (urban = 1, rural = 0), region (1 = north-east, 2 = north central, 3 = south, 4 = west), poverty (1 = yes, 0 = no), whether the infant had a normal birth weight as defined by weighting at least 5.5 pounds (`wmonth` = 1 if yes and 0 otherwise), race (1 = white, 2 = black, 3 = other), education (years of school of mother), month the child started to be on solid food (`sfmonth`), average number of cigarettes smoked per week during pregnancy (`smoke` = 0, 1 or 2) and alcohol used by mother during pregnancy (0, 1, 2), where the higher the number the higher the frequency of alcohol consumption. To capture the effect of housing crowding (since pneumonia is a communicable disease), number of siblings of the child (`nsibs`) was considered and grouped in three categories (0 for infants without siblings, 1 for infants with one to three siblings, and 2 for more than three siblings).

To assess whether the censoring mechanism was informative, we employed the AIC, BIC, and K -Fold Cross validation (Υ^{KCV}) with $K = 20$ (decreasing or increasing this value did not alter the conclusions); see Supplementary Material F for their definitions. Since several combinations of covariates and link functions had to be considered, a number of models were tried out and the final models selected using the above mentioned criteria. Table 2 in Supplementary Material F shows the results for the chosen models and supports the presence of informative censoring through the alcohol and region variables (Model 3). Table 2 and Figure 1 present the results for Model 3 and Model 1 (the latter neglects informative censoring).

Main findings: From a quick overall look at Table 2, the results exhibit a smaller estimation uncertainty for the informative model. Analysing the table in more detail, the coefficients of `wmonth`,

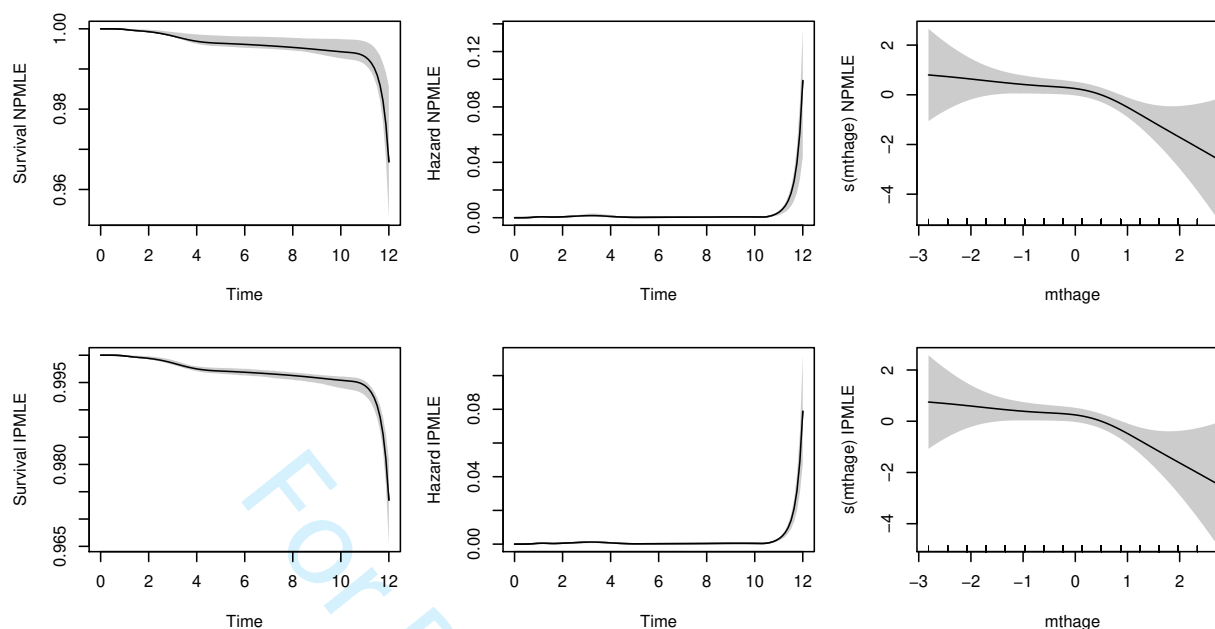


Figure 1: Smooth function estimates and their corresponding 95% intervals for Model 1 (non-informative model) and Model 3 (informative model) obtained by applying `gamLSS()` in GJRM to pneumonia data. The intervals have been obtained using the approach described in Supplementary Material E.

`nsibs1`, `nsibs2` are statistically significant for both models. For instance, the expected hazard for infants with one to three siblings is 2 times that for infants without siblings. Similarly, the expected hazard is 6.4 times higher in infants with more than 3 siblings as compared to infants with no siblings. The parameters of categories `alcohol1` and `region4` of the respective variables are statistically significant at the 10% level for the informative model and are not significant for the non-informative model. The implication of this result is that using the non-informative model the variables `alcohol` and `region` would most likely be removed from the model, hence missing out on some potentially important behavioral and geographical patterns. The table also shows that the smooth functions estimates for $s(u)$ and $s(\text{mthage})$ are statistically significant for both models, whereas Figure 1 displays their estimated functional forms along with the survival and hazard curves. The plots show, for instance, that, after a certain point, the hazard to contract pneumonia decreases with mother's age. The survival and hazard curves are very similar across the two models with the main difference that the informative approach yields considerably less variable estimates. Our results are consistent with those of Lu & Zhang (2012) who found that the censoring mechanism is informative in this dataset, and that the informative model provides a better fit as compared to its non-informative counterpart.

(a) Model 1 (NPMLE)				
Linear Covariates	Estimate	Standart Error	Z-value	P-value
intercept	-71.28	45.52	-1.566	0.117
alcohol1	0.364	0.310	1.174	0.240
alcohol2	-0.130	0.336	-0.386	0.700
nsibs1	0.696	0.258	2.695	0.007 **
nsibs2	1.833	0.761	2.408	0.016 *
region2	-0.004	0.343	-0.012	0.991
region3	-0.489	0.343	-1.426	0.154
region4	-0.698	0.438	-1.595	0.111
wmonth	-0.767	0.293	-2.617	0.009 **
Smooth Variables	EDF	Ref. DF	Chi-square	P-value
s (u)	7.776	8.640	101.94	<2e-16 ***
s (mthage)	2.503	3.171	10.41	0.019 *
(b) Model 3 (IPMLE)				
Linear Covariates	Estimate	Standart Error	Z-value	P-value
intercept	-71.50	45.51	-1.571	0.116
alcohol1	0.086	0.046	1.859	0.063 ·
alcohol2	0.022	0.046	0.472	0.637
nsibs1	0.687	0.257	2.670	0.008 **
nsibs2	1.860	0.760	2.448	0.014 *
region2	-0.063	0.056	-1.135	0.256
region3	-0.017	0.052	-0.325	0.745
region4	-0.107	0.059	-1.814	0.070 ·
wmonth	-0.761	0.291	-2.616	0.009 **
Smooth Variables	EDF	Ref. DF	Chi-square	P-value
s (u)	7.776	8.640	101.59	<2e-16 ***
s (mthage)	2.466	3.127	9.501	0.026 *

Table 2: Estimation results of the non-informative and informative models (Models 1 and 3, respectively, in Table 5 in Supplementary Material F) applied to pneumonia data. The models were fitted using `gamlss()` in GJRM by employing the "PH-PH" link functions combination. Furthermore, EDF and Ref. DF refer to the effective degrees of freedom and reference degrees of freedom of the smooths. More details can be founded in Supplementary Materials C and E.

5 Discussion

In this article, we have introduced generalized link-based additive survival models with informative censoring and their potential illustrated using simulated and real data. The proofs of the \sqrt{n} -consistency and asymptotic normality of the non-informative and informative estimators have been provided. Further, we showed that the newly introduced informative estimator is more efficient than its non-informative counterpart.

Important features of the modelling framework are that: the survival models can account for informative censoring; the baseline functions are estimated non-parametrically via means of monotonic P-splines, which allows one to obtain coherent estimated survival functions; covariate effects are flexibly determined using additive predictors; the optimization scheme allows for the reliable simultaneous penalized estimation of all model's parameters as well as for stable and fast automatic multiple smoothing parameter selection; the models can be easily utilized using the freely available GJRM R package which allows for several modelling choice as well as for transparent and reproducible research. Given that the assumption of absence of informative censoring is often made for mathematical convenience as well as lack of software, the proposed methodology is likely to appeal researchers in various fields wishing to estimate possibly more realistic survival models.

Future research will focus on extending the proposed informative model to include time varying covariates along with the incorporation of left and interval censored responses, and on the construction of efficient schemes for selecting automatically the set of informative covariates. We will also look into the case of dependent censoring as well as alternative estimation approaches such as sieve maximum likelihood.

References

- Andersen, P. K. & Keiding, N. (2006). *Survival and event history analysis*. Wiley Chichester.
- Chen, Y.-H. (2010). Semiparametric marginal regression analysis for dependent competing risks under an assumed copula. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(2), 235–251.

- 1 Cox, D. (1972). Regression models and life tables (with discussion). *Journal of the Royal Statis-*
2 *tical Society, Series B*, 34(2), 187–220.
3
4
5 Cox, D. R. (1959). The analysis of exponentially distributed life-times with two types of failure.
6 *Journal of the Royal Statistical Society: Series B (Methodological)*, 21(2), 411–421.
7
8
9
10 Crowder, M. J. (2012). *Multivariate survival analysis and competing risks*. Chapman and
11 Hall/CRC.
12
13
14
15 Crowther, M. J. & Lambert, P. C. (2013). Simulating biologically plausible complex survival data.
16 *Statistics in medicine*, 32(23), 4118–4134.
17
18
19
20 Emura, T. & Chen, Y.-H. (2018). *Analysis of Survival Data with Dependent Censoring: Copula-*
21 *Based Approaches*. Springer.
22
23
24
25 Hjort, N. L. (1992). On inference in parametric survival data models. *International Statistical*
26 *Review/Revue Internationale de Statistique*, 60(3), 355–387.
27
28
29
30 Huang, X. & Zhang, N. (2008). Regression survival analysis with an assumed copula for depen-
31 dent censoring: a sensitivity analysis approach. *Biometrics*, 64(4), 1090–1099.
32
33
34
35 Kalbfleisch, J. D. & Prentice, R. L. (2002). *The Statistical Analysis of Failure Time Data*. 2nd Ed.
36 Hoboken, Wiley.
37
38
39
40 Koziol, J. A. & Green, S. B. (1976). A cramér-von mises statistic for randomly censored data.
41 *Biometrika*, 63(3), 465–474.
42
43
44
45 Lagakos, S. (1979). General right censoring and its impact on the analysis of survival data. *Bio-*
46 *metrics*, 35(1), 139–156.
47
48
49
50 Leitenstorfer, F. & Tutz, G. (2006). Generalized monotonic regression based on b-splines with an
51 application to air pollution data. *Biostatistics*, 8(3), 654–673.
52
53
54
55 Li, R. & Peng, L. (2015). Quantile regression adjusting for dependent censoring from semicom-
56 peting risks. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 77(1),
57 107–130.
58
59
60

- 1 Liu, X.-R., Pawitan, Y., & Clements, M. (2018). Parametric and penalized generalized survival
2 models. *Statistical methods in medical research*, 27(5), 1531–1546.
3
4
5 Lu, Z. & Zhang, W. (2012). Semiparametric likelihood estimation in survival models with infor-
6 mative censoring. *Journal of Multivariate Analysis*, 106, 187–211.
7
8
9
10 Ma, J., Heritier, S., & L \hat{o} , S. N. (2014). On the maximum penalized likelihood approach for
11 proportional hazard models with right censored survival data. *Computational Statistics & Data*
12 *Analysis*, 74, 142–156.
13
14
15
16
17 Marra, G. & Radice, R. (2019). *GJRM: Generalised Joint Regression Modelling*. R package
18 version 0.2-1.
19
20
21
22 Marra, G., Radice, R., B \ddot{a} rnighausen, T., Wood, S. N., & McGovern, M. E. (2017). A simultaneous
23 equation approach to estimating hiv prevalence with nonignorable missing responses. *Journal*
24 *of the American Statistical Association*, 112(518), 484–496.
25
26
27
28
29 NLSY (1995). National Longitudinal Survey of Youth Handbook. The Ohio State University.
30
31
32
33 Nocedal, J. & Wright, S. (2006). *Numerical optimization, series in operations research and finan-*
34 *cial engineering*. Springer, New York, USA, 2006.
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

- 1 Tsiatis, A. (1975). A nonidentifiability aspect of the problem of competing risks. *Proceedings of*
2 *the National Academy of Sciences*, 72(1), 20–22.
3
4
5 Vatter, T. & Chavez-Demoulin, V. (2015). Generalized additive models for conditional dependence
6 structures. *Journal of Multivariate Analysis*, 141, 147–167.
7
8
9
10 Wang, A., Chandra, K., Xu, R., & Sun, J. (2015). The identifiability of dependent competing risks
11 models induced by bivariate frailty models. *Scandinavian Journal of Statistics*, 42(2), 427–437.
12
13
14
15 Wood, S. N. (2017). *Generalized Additive Models: An Introduction With R*. Second Edition,
16 Chapman & Hall/CRC, London.
17
18
19
20 Xu, J., Ma, J., Connors, M. H., & Brodaty, H. (2018). Proportional hazard model estimation under
21 dependent censoring using copulas and penalized likelihood. *Statistics in medicine*, 37(14),
22 2238–2251.
23
24
25
26
27 Xu, J., Ma, J., & Prvan, T. (2017). Non parametric hazard estimation with dependent censoring
28 using penalized likelihood and an assumed copula. *Communications in Statistics-Theory and*
29 *Methods*, 46(22), 11383–11403.
30
31
32
33
34 Xue, X., Xie, X., & Strickler, H. D. (2018). A censored quantile regression approach for the
35 analysis of time to event data. *Statistical methods in medical research*, 27(3), 955–965.
36
37
38
39 Younes, N. & Lachin, J. (1997). Link-based models for survival data with interval and continuous
40 time censoring. *Biometrics*, 53(4), 1199–1211.
41
42
43
44 Yuan, M. (2005). Semiparametric censorship model with covariates. *Test*, 14(2), 489–514.
45
46
47
48 Zeng, D. et al. (2004). Estimating marginal survival function by adjusting for dependent censoring
49 using many covariates. *The Annals of Statistics*, 32(4), 1533–1555.
50
51
52
53 Zheng, M. & Klein, J. P. (1995). Estimates of marginal survival for dependent competing risks
54 based on an assumed copula. *Biometrika*, 82(1), 127–138.
55
56
57
58
59
60

Supplementary Material: "Generalized Link-Based Additive Survival Models with Informative Censoring"

Supplementary Material A: Software

The models proposed in this article can be employed via the `gamlss()` function in the R package GJRM (Marra & Radice, 2019). As an example, consider the following call

```
f1 <- list(u ~ s(u, bs = "mpi") + z1 + s(z2), u ~ s(u, bs = "mpi") + z1 + s(z2))
M1 <- gamlss(f1, data = data, surv = TRUE, margin = "PH", margin2 = "PH"
  cens = delta, informative = "yes", inform.cov = c("z1"))
```

where `f1` is a list containing the two additive predictors of the informative model, and `s(u, bs = "mpi")` represents the monotonic P-spline function which models a transformation of the baseline survival function. As for `s(z2)`, the default is `bs = "tp"` (penalized low rank thin plate spline) with `k = 10` (number of basis functions) and `m = 2` (order of derivatives). However, argument `bs` can also be set to, for example, `cr` (penalized cubic regression spline), `ps` (P-spline) and `mrf` (Markov random field), to name but a few. In the `gamlss` function, `surv = TRUE` indicates that a survival model is fitted. The arguments `margin = "PH"` and `margin2 = "PH"` specify the link functions for the survival and censoring times, respectively. Table 1 shows the possible choices for the links that have been implemented for this article. In this example, we specify the proportional hazard link ("PH") for the two equations. Argument `cens = delta` is a binary censoring indicator; this variable has to be equal to 1 if the event occurred and 0 otherwise. Finally, `informative = "yes"` indicates that we are fitting a survival model with informative censoring, and `inform.cov = c("z1")` specifies the set of informative covariates.

Model	Link $g(S)$	Inverse link $g^{-1}(\xi) = G(\xi)$	$G'(\xi)$
Prop.hazards ("PH")	$\log \{-\log(S)\}$	$\exp \{-\exp(\xi)\}$	$-G'(\xi) \exp(\xi)$
Prop.odds ("PO")	$-\log \left(\frac{S}{1-S} \right)$	$\frac{\exp(-\xi)}{1+\exp(-\xi)}$	$-G^2(\xi) \exp(-\xi)$
Probit ("probit")	$-\Phi^{-1}(S)$	$\Phi(-\xi)$	$-\phi(-\xi)$

Table 1: Link functions implemented in GJRM. Φ and ϕ are the cumulative distribution and density functions of a univariate standard normal distribution. Alternative links can be implemented. The first two functions can be called log-log and -logit links, respectively.

Supplementary Material B: Scores and Hessians

In this section, the detailed derivations of the informative and non-informative Scores and Hessians are presented.

B.1. Informative and Non-informative Scores

If censoring is informative then γ_1 and γ_2 would have some components in common. Because the first Q components of γ_1 are the same as the first Q components of γ_2 , we have

$$\mathcal{Q}_{\nu i}^\top \gamma_\nu = \mathcal{Q}_i^{0\top} \alpha_0 + \mathcal{Q}_{\nu i}^{1\top} \alpha_\nu.$$

Therefore, defining $\alpha = (\alpha_0^\top, \alpha_1^\top, \alpha_2^\top)^\top$, the informative penalized log-likelihood function can be written as

$$\ell_p(\alpha) = \ell(\alpha) - \frac{1}{2} \alpha^\top \mathcal{S} \alpha, \quad (1)$$

where $\ell(\alpha)$ is defined as

$$\begin{aligned} \ell(\alpha) = & \sum_{i=1}^n \left\{ \log \mathcal{G}_1 [\xi_{1i}(\alpha_0, \alpha_1)] + \delta_{1i} \log \left\{ -\frac{\mathcal{G}'_1 [\xi_{1i}(\alpha_0, \alpha_1)]}{\mathcal{G}_1 [\xi_{1i}(\alpha_0, \alpha_1)]} \frac{\partial \xi_{1i}(\alpha_0, \alpha_1)}{\partial y_i} \right\} \right\} \\ & + \sum_{i=1}^n \left\{ \log \mathcal{G}_2 [\xi_{2i}(\alpha_0, \alpha_2)] + \delta_{2i} \log \left\{ -\frac{\mathcal{G}'_2 [\xi_{2i}(\alpha_0, \alpha_2)]}{\mathcal{G}_2 [\xi_{2i}(\alpha_0, \alpha_2)]} \frac{\partial \xi_{2i}(\alpha_0, \alpha_2)}{\partial y_i} \right\} \right\}. \end{aligned}$$

The gradient of equation (1) can be calculated as

$$\nabla_\alpha \ell_p(\alpha) = \nabla_\alpha \ell(\alpha) - \alpha \mathcal{S},$$

where $\nabla_{\alpha} \ell(\alpha) = (\nabla_{\alpha_0} \ell(\alpha)^{\top}, \nabla_{\alpha_1} \ell(\alpha)^{\top}, \nabla_{\alpha_2} \ell(\alpha)^{\top})^{\top}$. where $\nabla_{\alpha_0} \ell(\alpha)$, $\nabla_{\alpha_1} \ell(\alpha)$ and $\nabla_{\alpha_2} \ell(\alpha)$ can be obtained as $\frac{\partial \ell(\alpha)}{\partial \alpha_0} = \left[\frac{\partial \ell(\alpha)}{\partial \alpha_{011}} \dots \frac{\partial \ell(\alpha)}{\partial \alpha_{0Q_1 J_1 Q_1}} \right]^{\top}$, $\frac{\partial \ell(\alpha)}{\partial \alpha_1} = \left[\frac{\partial \ell(\alpha)}{\partial \alpha_{111}} \dots \frac{\partial \ell(\alpha)}{\partial \alpha_{1Q_1 J_1 Q_1}} \right]^{\top}$ and $\frac{\partial \ell(\alpha)}{\partial \alpha_2} = \left[\frac{\partial \ell(\alpha)}{\partial \alpha_{21}} \dots \frac{\partial \ell(\alpha)}{\partial \alpha_{2Q_2 J_2 Q_2}} \right]^{\top}$. In particular, the scalar derivatives of $\nabla_{\alpha_0} \ell(\alpha)$, $\nabla_{\alpha_1} \ell(\alpha)$ and $\nabla_{\alpha_2} \ell(\alpha)$ can be calculated as

$$\begin{aligned} \frac{\partial \ell(\alpha)}{\partial \alpha_{0j}} &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \right\} + \sum_{i=1}^n \delta_{1i} \left\{ \left[-\frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial y_i} \right]^{-1} \left[-\frac{\mathcal{G}''_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial y_i} + \frac{\mathcal{G}'_1{}^2}{\mathcal{G}_1^2} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial y_i} - \frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \alpha_{0j}} \right] \right\} \\ &+ \sum_{i=1}^n \left\{ \frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \right\} + \sum_{i=1}^n \delta_{2i} \left\{ \left[-\frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial y_i} \right]^{-1} \left[-\frac{\mathcal{G}''_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial y_i} + \frac{\mathcal{G}'_2{}^2}{\mathcal{G}_2^2} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial y_i} - \frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \alpha_{0j}} \right] \right\} \\ &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} + \delta_{1i} \left[\frac{\mathcal{G}''_1}{\mathcal{G}'_1} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} - \frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} + \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \alpha_{0j}} \left(\frac{\partial \xi_{1i}}{\partial y_i} \right)^{-1} \right] \right\} \\ &+ \sum_{i=1}^n \left\{ \frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} + \delta_{2i} \left[\frac{\mathcal{G}''_2}{\mathcal{G}'_2} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} - \frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} + \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \alpha_{0j}} \left(\frac{\partial \xi_{2i}}{\partial y_i} \right)^{-1} \right] \right\} \\ &= \sum_{i=1}^n \left\{ \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \left[\frac{\mathcal{G}'_1}{\mathcal{G}_1} + \delta_{1i} \left(\frac{\mathcal{G}''_1}{\mathcal{G}'_1} - \frac{\mathcal{G}'_1}{\mathcal{G}_1} \right) \right] + \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \left[\frac{\mathcal{G}'_2}{\mathcal{G}_2} + \delta_{2i} \left(\frac{\mathcal{G}''_2}{\mathcal{G}'_2} - \frac{\mathcal{G}'_2}{\mathcal{G}_2} \right) \right] \right\} \\ &= \sum_{i=1}^n \left\{ \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \Delta_1 + \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \Delta_2 \right\}, \end{aligned} \tag{2}$$

$$\begin{aligned} \frac{\partial \ell(\alpha)}{\partial \alpha_{1j}} &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} \right\} + \sum_{i=1}^n \delta_{1i} \left\{ \left[-\frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial y_i} \right]^{-1} \left[-\frac{\mathcal{G}''_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} \frac{\partial \xi_{1i}}{\partial y_i} + \frac{\mathcal{G}'_1{}^2}{\mathcal{G}_1^2} \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} \frac{\partial \xi_{1i}}{\partial y_i} - \frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \alpha_{1j}} \right] \right\} \\ &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} + \delta_{1i} \left[\frac{\mathcal{G}''_1}{\mathcal{G}'_1} \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} - \frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} + \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \alpha_{1j}} \left(\frac{\partial \xi_{1i}}{\partial y_i} \right)^{-1} \right] \right\} \\ &= \sum_{i=1}^n \left\{ \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} \left[\frac{\mathcal{G}'_1}{\mathcal{G}_1} + \delta_{1i} \left(\frac{\mathcal{G}''_1}{\mathcal{G}'_1} - \frac{\mathcal{G}'_1}{\mathcal{G}_1} \right) \right] + \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \alpha_{1j}} \delta_{1i} \left(\frac{\partial \xi_{1i}}{\partial y_i} \right)^{-1} \right\} \\ &= \sum_{i=1}^n \left\{ \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} \Delta_1 + \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \alpha_{1j}} \Omega_1 \right\}, \end{aligned} \tag{3}$$

$$\begin{aligned}
\frac{\partial \ell(\boldsymbol{\alpha})}{\partial \alpha_{2j}} &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} \right\} + \sum_{i=1}^n \delta_{2i} \left\{ \left[-\frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial y_i} \right]^{-1} \left[-\frac{\mathcal{G}''_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} \frac{\partial \xi_{2i}}{\partial y_i} + \frac{\mathcal{G}'_2{}^2}{\mathcal{G}_2^2} \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} \frac{\partial \xi_{2i}}{\partial y_i} - \frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \alpha_{2j}} \right] \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} + \delta_{2i} \left[\frac{\mathcal{G}''_2}{\mathcal{G}'_2} \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} - \frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} + \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \alpha_{2j}} \left(\frac{\partial \xi_{2i}}{\partial y_i} \right)^{-1} \right] \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} \left[\frac{\mathcal{G}'_2}{\mathcal{G}_2} + \delta_{2i} \left(\frac{\mathcal{G}''_2}{\mathcal{G}'_2} - \frac{\mathcal{G}'_2}{\mathcal{G}_2} \right) \right] + \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \alpha_{2j}} \delta_{2i} \left(\frac{\partial \xi_{2i}}{\partial y_i} \right)^{-1} \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} \Delta_2 + \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \alpha_{2j}} \Omega_2 \right\},
\end{aligned} \tag{4}$$

where $\xi_{\nu i} = \xi_{\nu i}(\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_\nu)$, $\Delta_\nu = \left[\frac{\mathcal{G}'_\nu}{\mathcal{G}_\nu} + \delta_{\nu i} \left(\frac{\mathcal{G}''_\nu}{\mathcal{G}'_\nu} - \frac{\mathcal{G}'_\nu}{\mathcal{G}_\nu} \right) \right]$ and $\Omega_\nu = \delta_{\nu i} \left(\frac{\partial \xi_{\nu i}}{\partial y_i} \right)^{-1}$. The last terms of equations (2), (3) and (4) allow to express $\nabla_{\boldsymbol{\alpha}_0} \ell(\boldsymbol{\alpha})$, $\nabla_{\boldsymbol{\alpha}_1} \ell(\boldsymbol{\alpha})$ and $\nabla_{\boldsymbol{\alpha}_2} \ell(\boldsymbol{\alpha})$ as follow

$$\begin{aligned}
\nabla_{\boldsymbol{\alpha}_0} \ell(\boldsymbol{\alpha}) &= \sum_{i=1}^n \left[\Delta_1 \frac{\partial \xi_{1i}}{\partial \boldsymbol{\alpha}_0} + \Delta_2 \frac{\partial \xi_{2i}}{\partial \boldsymbol{\alpha}_0} \right], \\
\nabla_{\boldsymbol{\alpha}_1} \ell(\boldsymbol{\alpha}) &= \sum_{i=1}^n \left[\Delta_1 \frac{\partial \xi_{1i}}{\partial \boldsymbol{\alpha}_1} + \Omega_1 \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \boldsymbol{\alpha}_1} \right], \\
\nabla_{\boldsymbol{\alpha}_2} \ell(\boldsymbol{\alpha}) &= \sum_{i=1}^n \left[\Delta_2 \frac{\partial \xi_{2i}}{\partial \boldsymbol{\alpha}_2} + \Omega_2 \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \boldsymbol{\alpha}_2} \right],
\end{aligned}$$

where, for all $i = 1, \dots, n$ and $\nu = 1, 2$, $\frac{\partial \xi_{\nu i}}{\partial \boldsymbol{\alpha}_0} = \left[\frac{\partial \xi_{\nu i}}{\partial \alpha_{011}} \dots \frac{\partial \xi_{\nu i}}{\partial \alpha_{0QJ_Q}} \right]^\top$, $\frac{\partial \xi_{\nu i}}{\partial \boldsymbol{\alpha}_\nu} = \left[\frac{\partial \xi_{\nu i}}{\partial \alpha_{\nu 11}} \dots \frac{\partial \xi_{\nu i}}{\partial \alpha_{\nu Q_\nu J_{\nu Q_\nu}} \right]^\top$ and $\frac{\partial^2 \xi_{\nu i}}{\partial y_i \partial \boldsymbol{\alpha}_\nu} = \left[\frac{\partial^2 \xi_{\nu i}}{\partial y_i \partial \alpha_{\nu 11}} \dots \frac{\partial^2 \xi_{\nu i}}{\partial y_i \partial \alpha_{\nu Q_\nu J_{\nu Q_\nu}} \right]^\top$. These expressions can be calculated using the design vectors defined in Section 2.2 as

$$\begin{aligned}
\frac{\partial \xi_{\nu i}}{\partial \boldsymbol{\alpha}_0} &= (\boldsymbol{Q}_1(\mathbf{x}_{1i}^0)^\top, \dots, \boldsymbol{Q}_Q(\mathbf{x}_{Qi}^0)^\top)^\top = \boldsymbol{Q}_i^0, \\
\frac{\partial \xi_{\nu i}}{\partial y_i} &= \lim_{\varepsilon \rightarrow 0} \left\{ \frac{\boldsymbol{Q}_{\nu 0}(y_i + \varepsilon) - \boldsymbol{Q}_{\nu 0}(y_i - \varepsilon)}{2\varepsilon} \right\}^\top \boldsymbol{\Gamma}_{\nu 0} \tilde{\boldsymbol{\alpha}}_{\nu 0} = \boldsymbol{Q}'_{\nu 0}(y_i)^\top \boldsymbol{\Gamma}_{\nu 0} \tilde{\boldsymbol{\alpha}}_{\nu 0}, \\
\frac{\partial \xi_{\nu i}}{\partial \boldsymbol{\alpha}_{\nu q_\nu}} &= \begin{cases} \boldsymbol{Q}'_{\nu 0}{}^{\Delta}(y_i) & \text{if } \boldsymbol{\alpha}_{\nu q_\nu} = \boldsymbol{\alpha}_{\nu 0} \\ \boldsymbol{Q}_{\nu q_\nu}(\mathbf{x}_{\nu q_\nu i}^1) & \text{otherwise,} \end{cases} \\
\frac{\partial^2 \xi_{\nu i}}{\partial y_i \partial \boldsymbol{\alpha}_{\nu q_\nu}} &= \begin{cases} \boldsymbol{Q}'_{\nu 0}{}^{\Delta'}(y_i) & \text{if } \boldsymbol{\alpha}_{\nu q_\nu} = \boldsymbol{\alpha}_{\nu 0} \\ \mathbf{0} & \text{otherwise,} \end{cases}
\end{aligned}$$

where $\mathcal{Q}'_{\nu 0}(y_i)$ can be conveniently obtained using a finite-difference method. Moreover, we define the design vectors $\mathcal{Q}'_{\nu 0 \Delta}(y_i)$ and $\mathcal{Q}'_{\nu 0 \Delta'}(y_i)$ as

$$\mathcal{Q}'_{\nu 0 \Delta}(y_i) = \begin{bmatrix} \sum_{j_{\nu 0}=1}^{J_{\nu 0}} \mathcal{Q}_{\nu 0 j_{\nu 0}}(y_i) \\ [\sum_{j_{\nu 0}=2}^{J_{\nu 0}} \mathcal{Q}_{\nu 0 j_{\nu 0}}(y_i)] \exp(\alpha_{\nu 02}) \\ [\sum_{j_{\nu 0}=3}^{J_{\nu 0}} \mathcal{Q}_{\nu 0 j_{\nu 0}}(y_i)] \exp(\alpha_{\nu 03}) \\ \vdots \\ \mathcal{Q}_{\nu 0 J_{\nu 0}}(y_i) \exp(\alpha_{\nu 0 J_{\nu 0}}) \end{bmatrix} \quad \mathcal{Q}'_{\nu 0 \Delta'}(y_i) = \begin{bmatrix} \sum_{j_{\nu 0}=1}^{J_{\nu 0}} \mathcal{Q}'_{\nu 0 j_{\nu 0}}(y_i) \\ [\sum_{j_{\nu 0}=2}^{J_{\nu 0}} \mathcal{Q}'_{\nu 0 j_{\nu 0}}(y_i)] \exp(\alpha_{\nu 02}) \\ [\sum_{j_{\nu 0}=3}^{J_{\nu 0}} \mathcal{Q}'_{\nu 0 j_{\nu 0}}(y_i)] \exp(\alpha_{\nu 03}) \\ \vdots \\ \mathcal{Q}'_{\nu 0 J_{\nu 0}}(y_i) \exp(\alpha_{\nu 0 J_{\nu 0}}) \end{bmatrix}.$$

On the other hand, when censoring is non-informative the penalized log-likelihood function is

$$\ell_p(\gamma) = \ell(\gamma) - \frac{1}{2} \gamma^\top \mathcal{S} \gamma, \quad (5)$$

where $\ell(\gamma)$ can be written as

$$\begin{aligned} \ell(\gamma) = & \sum_{i=1}^n \left\{ \log \mathcal{G}_1 [\xi_{1i}(\gamma_1)] + \delta_{1i} \log \left\{ -\frac{\mathcal{G}'_1 [\xi_{1i}(\gamma_1)]}{\mathcal{G}_1 [\xi_{1i}(\gamma_1)]} \frac{\partial \xi_{1i}(\gamma_1)}{\partial y_i} \right\} \right\} \\ & + \sum_{i=1}^n \left\{ \log \mathcal{G}_2 [\xi_{2i}(\gamma_2)] + \delta_{2i} \log \left\{ -\frac{\mathcal{G}'_2 [\xi_{2i}(\gamma_2)]}{\mathcal{G}_2 [\xi_{2i}(\gamma_2)]} \frac{\partial \xi_{2i}(\gamma_2)}{\partial y_i} \right\} \right\}. \end{aligned}$$

The gradient of (5) can be calculated as

$$\nabla_{\gamma} \ell_p(\gamma) = \nabla_{\gamma} \ell(\gamma) - \gamma \mathcal{S},$$

where $\nabla_{\gamma} \ell(\gamma) = (\nabla_{\gamma_1} \ell(\gamma)^\top, \nabla_{\gamma_2} \ell(\gamma)^\top)^\top$. In addition, $\nabla_{\gamma_1} \ell(\gamma)$ and $\nabla_{\gamma_2} \ell(\gamma)$ can be calculated as $\frac{\partial \ell(\gamma)}{\partial \gamma_1} = \left[\frac{\partial \ell(\gamma)}{\partial \gamma_{111}} \dots \frac{\partial \ell(\gamma)}{\partial \gamma_{1K_1 J_{1K_1}}} \right]^\top$ and $\frac{\partial \ell(\gamma)}{\partial \gamma_2} = \left[\frac{\partial \ell(\gamma)}{\partial \gamma_{211}} \dots \frac{\partial \ell(\gamma)}{\partial \gamma_{2K_2 J_{2K_2}}} \right]^\top$. Furthermore,

the scalar derivatives of $\nabla_{\gamma_1} \ell(\gamma)$ and $\nabla_{\gamma_2} \ell(\gamma)$ can be obtained as

$$\begin{aligned}
 \frac{\partial \ell(\gamma)}{\partial \gamma_{1j}} &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \gamma_{1j}} \right\} + \sum_{i=1}^n \delta_{1i} \left\{ \left[-\frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial y_i} \right]^{-1} \left[-\frac{\mathcal{G}''_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \gamma_{1j}} \frac{\partial \xi_{1i}}{\partial y_i} + \frac{\mathcal{G}_1'^2}{\mathcal{G}_1^2} \frac{\partial \xi_{1i}}{\partial \gamma_{1j}} \frac{\partial \xi_{1i}}{\partial y_i} - \frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \gamma_{1j}} \right] \right\} \\
 &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \gamma_{1j}} + \delta_{1i} \left[\frac{\mathcal{G}''_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \gamma_{1j}} - \frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \gamma_{1j}} + \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \gamma_{1j}} \left(\frac{\partial \xi_{1i}}{\partial y_i} \right)^{-1} \right] \right\} \\
 &= \sum_{i=1}^n \left\{ \frac{\partial \xi_{1i}}{\partial \gamma_{1j}} \left[\frac{\mathcal{G}'_1}{\mathcal{G}_1} + \delta_{1i} \left(\frac{\mathcal{G}''_1}{\mathcal{G}_1} - \frac{\mathcal{G}'_1}{\mathcal{G}_1} \right) \right] + \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \gamma_{1j}} \delta_{1i} \left(\frac{\partial \xi_{1i}}{\partial y_i} \right)^{-1} \right\} \\
 &= \sum_{i=1}^n \left\{ \frac{\partial \xi_{1i}}{\partial \gamma_{1j}} \Delta_1 + \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \gamma_{1j}} \Omega_1 \right\},
 \end{aligned} \tag{6}$$

$$\begin{aligned}
 \frac{\partial \ell(\gamma)}{\partial \gamma_{2j}} &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \gamma_{2j}} \right\} + \sum_{i=1}^n \delta_{2i} \left\{ \left[-\frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial y_i} \right]^{-1} \left[-\frac{\mathcal{G}''_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \gamma_{2j}} \frac{\partial \xi_{2i}}{\partial y_i} + \frac{\mathcal{G}_2'^2}{\mathcal{G}_2^2} \frac{\partial \xi_{2i}}{\partial \gamma_{2j}} \frac{\partial \xi_{2i}}{\partial y_i} - \frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \gamma_{2j}} \right] \right\} \\
 &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \gamma_{2j}} + \delta_{2i} \left[\frac{\mathcal{G}''_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \gamma_{2j}} - \frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \gamma_{2j}} + \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \gamma_{2j}} \left(\frac{\partial \xi_{2i}}{\partial y_i} \right)^{-1} \right] \right\} \\
 &= \sum_{i=1}^n \left\{ \frac{\partial \xi_{2i}}{\partial \gamma_{2j}} \left[\frac{\mathcal{G}'_2}{\mathcal{G}_2} + \delta_{2i} \left(\frac{\mathcal{G}''_2}{\mathcal{G}_2} - \frac{\mathcal{G}'_2}{\mathcal{G}_2} \right) \right] + \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \gamma_{2j}} \delta_{2i} \left(\frac{\partial \xi_{2i}}{\partial y_i} \right)^{-1} \right\} \\
 &= \sum_{i=1}^n \left\{ \frac{\partial \xi_{2i}}{\partial \gamma_{2j}} \Delta_2 + \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \gamma_{2j}} \Omega_2 \right\},
 \end{aligned} \tag{7}$$

where $\xi_{\nu i} = \xi_{\nu i}(\gamma_\nu)$. The last terms of equations (6) and (7) allow $\nabla_{\gamma_1} \ell(\gamma)$ and $\nabla_{\gamma_2} \ell(\gamma)$ to be expressed as

$$\begin{aligned}
 \nabla_{\gamma_1} \ell(\gamma) &= \sum_{i=1}^n \left[\Delta_1 \frac{\partial \xi_{1i}}{\partial \gamma_1} + \Omega_1 \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \gamma_1} \right] \\
 \nabla_{\gamma_2} \ell(\gamma) &= \sum_{i=1}^n \left[\Delta_2 \frac{\partial \xi_{2i}}{\partial \gamma_2} + \Omega_2 \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \gamma_2} \right],
 \end{aligned}$$

where $\frac{\partial \xi_{\nu i}}{\partial \gamma_\nu} = \left[\frac{\partial \xi_{\nu i}}{\partial \gamma_{\nu 11}} \dots \frac{\partial \xi_{\nu i}}{\partial \gamma_{\nu K_\nu J_\nu K_\nu}} \right]^\top$ and $\frac{\partial^2 \xi_{\nu i}}{\partial y_i \partial \gamma_\nu} = \left[\frac{\partial^2 \xi_{\nu i}}{\partial y_i \partial \gamma_{\nu 11}} \dots \frac{\partial^2 \xi_{\nu i}}{\partial y_i \partial \gamma_{\nu K_\nu J_\nu K_\nu}} \right]^\top$ for all $i = 1, \dots, n$ and $\nu = 1, 2$. Furthermore, $\frac{\partial \xi_{\nu i}(\gamma_\nu)}{\partial y_i}$, can be generically calculated using

$$\frac{\partial \xi_{\nu i}(\gamma_\nu)}{\partial y_i} = \lim_{\varepsilon \rightarrow 0} \left\{ \frac{\mathcal{Q}_{\nu 0}(y_i + \varepsilon) - \mathcal{Q}_{\nu 0}(y_i - \varepsilon)}{2\varepsilon} \right\}^\top \Gamma_{\nu 0} \hat{\gamma}_{\nu 0} = \mathcal{Q}'_{\nu 0}(y_i)^\top \Gamma_{\nu 0} \hat{\gamma}_{\nu 0},$$

where $\mathcal{Q}'_{\nu 0}(y_i)$ can also be calculated using a finite-difference method. The design vectors for $\frac{\partial \xi_{\nu i}(\gamma_\nu)}{\partial \gamma_\nu}$ and $\frac{\partial^2 \xi_{\nu i}(\gamma_\nu)}{\partial y_i \partial \gamma_\nu}$ can be obtained using

$$\frac{\partial \xi_{\nu i}(\gamma_\nu)}{\partial \gamma_{\nu k_\nu}} = \begin{cases} \mathcal{Q}_{\nu 0}^\Delta(y_i) & \text{if } \gamma_{\nu k_\nu} = \gamma_{\nu 0} \\ \mathcal{Q}_{\nu k_\nu}(\mathbf{x}_{\nu k_\nu i}) & \text{otherwise,} \end{cases}$$

$$\frac{\partial^2 \xi_{\nu i}(\gamma_\nu)}{\partial y_i \partial \gamma_{\nu k_\nu}} = \begin{cases} \mathcal{Q}_{\nu 0}^{\Delta'}(y_i) & \text{if } \gamma_{\nu k_\nu} = \gamma_{\nu 0} \\ 0 & \text{otherwise.} \end{cases}$$

Finally, we have that

$$\mathcal{Q}_{\nu 0}^\Delta(y_i) = \begin{bmatrix} \sum_{j_{\nu 0}=1}^{J_{\nu 0}} \mathcal{Q}_{\nu 0 j_{\nu 0}}(y_i) \\ [\sum_{j_{\nu 0}=2}^{J_{\nu 0}} \mathcal{Q}_{\nu 0 j_{\nu 0}}(y_i)] \exp(\gamma_{\nu 02}) \\ [\sum_{j_{\nu 0}=3}^{J_{\nu 0}} \mathcal{Q}_{\nu 0 j_{\nu 0}}(y_i)] \exp(\gamma_{\nu 03}) \\ \vdots \\ \mathcal{Q}_{\nu 0 J_{\nu 0}}(y_i) \exp(\gamma_{\nu 0 J_{\nu 0}}) \end{bmatrix} \quad \mathcal{Q}_{\nu 0}^{\Delta'}(y_i) = \begin{bmatrix} \sum_{j_{\nu 0}=1}^{J_{\nu 0}} \mathcal{Q}'_{\nu 0 j_{\nu 0}}(y_i) \\ [\sum_{j_{\nu 0}=2}^{J_{\nu 0}} \mathcal{Q}'_{\nu 0 j_{\nu 0}}(y_i)] \exp(\gamma_{\nu 02}) \\ [\sum_{j_{\nu 0}=3}^{J_{\nu 0}} \mathcal{Q}'_{\nu 0 j_{\nu 0}}(y_i)] \exp(\gamma_{\nu 03}) \\ \vdots \\ \mathcal{Q}'_{\nu 0 J_{\nu 0}}(y_i) \exp(\gamma_{\nu 0 J_{\nu 0}}) \end{bmatrix}.$$

B.2. Informative and Non-informative Hessians

The informative penalized Hessian can be obtained as

$$\nabla_{\alpha\alpha} \ell_p(\alpha) = \nabla_{\alpha\alpha} \ell(\alpha) - \mathcal{S},$$

where $\nabla_{\alpha\alpha} \ell(\alpha)$ is

$$\nabla_{\alpha\alpha} \ell(\alpha) = \begin{bmatrix} \nabla_{\alpha_0 \alpha_0} \ell(\alpha) & \nabla_{\alpha_0 \alpha_1} \ell(\alpha) & \nabla_{\alpha_0 \alpha_2} \ell(\alpha) \\ \nabla_{\alpha_1 \alpha_0} \ell(\alpha) & \nabla_{\alpha_1 \alpha_1} \ell(\alpha) & \nabla_{\alpha_1 \alpha_2} \ell(\alpha) \\ \nabla_{\alpha_2 \alpha_0} \ell(\alpha) & \nabla_{\alpha_2 \alpha_1} \ell(\alpha) & \nabla_{\alpha_2 \alpha_2} \ell(\alpha) \end{bmatrix}. \quad (8)$$

In addition, $\nabla_{\alpha_v \alpha_\kappa} \ell(\alpha) = \frac{\partial^2 \ell(\alpha)}{\partial \alpha_v \partial \alpha_\kappa^\top}$, for all $v = 0, 1, 2$ and $\kappa = 0, 1, 2$. This expression is calculated using

$$\nabla_{\alpha_v \alpha_\kappa} \ell(\alpha) = \begin{bmatrix} \frac{\partial^2 \ell(\alpha)}{\partial \alpha_{v11} \partial \alpha_{\kappa 11}} & \cdots & \frac{\partial^2 \ell(\alpha)}{\partial \alpha_{v11} \partial \alpha_{\kappa Q_\kappa J_\kappa Q_\kappa}} \\ \cdots & \ddots & \cdots \\ \frac{\partial^2 \ell(\alpha)}{\partial \alpha_{v Q_v J_v Q_v} \partial \alpha_{\kappa 11}} & \cdots & \frac{\partial^2 \ell(\alpha)}{\partial \alpha_{v Q_v J_v Q_v} \partial \alpha_{\kappa Q_\kappa J_\kappa Q_\kappa}} \end{bmatrix}.$$

Since α_1 appears only in $\xi_{1i}(\alpha_0, \alpha_1)$ and α_2 only in $\xi_{2i}(\alpha_0, \alpha_2)$, then $\nabla_{\alpha_1 \alpha_2} \ell(\alpha) = \nabla_{\alpha_2 \alpha_1} \ell(\alpha) = \mathbf{0}$. Hence, (8) can be written as

$$\nabla_{\alpha \alpha} \ell(\alpha) = \begin{bmatrix} \nabla_{\alpha_0 \alpha_0} \ell(\alpha) & \nabla_{\alpha_0 \alpha_1} \ell(\alpha) & \nabla_{\alpha_0 \alpha_2} \ell(\alpha) \\ \nabla_{\alpha_1 \alpha_0} \ell(\alpha) & \nabla_{\alpha_1 \alpha_1} \ell(\alpha) & \mathbf{0} \\ \nabla_{\alpha_2 \alpha_0} \ell(\alpha) & \mathbf{0} & \nabla_{\alpha_2 \alpha_2} \ell(\alpha) \end{bmatrix}. \quad (9)$$

In equation (9), the scalar derivatives of $\nabla_{\alpha_0 \alpha_0} \ell(\alpha)$, $\nabla_{\alpha_1 \alpha_0} \ell(\alpha)$, $\nabla_{\alpha_0 \alpha_2} \ell(\alpha)$, $\nabla_{\alpha_1 \alpha_1} \ell(\alpha)$ and $\nabla_{\alpha_2 \alpha_2} \ell(\alpha)$, can be calculated as

$$\begin{aligned} \frac{\partial^2 \ell(\alpha)}{\partial \alpha_{0j} \partial \alpha_{0k}} &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}_1''}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{0k}} - \frac{\mathcal{G}_1'^2}{\mathcal{G}_1^2} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{0k}} + \frac{\mathcal{G}_1'''}{\mathcal{G}_1} \delta_{1i} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{0k}} - \frac{\mathcal{G}_1''^2}{\mathcal{G}_1^2} \delta_{1i} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{0k}} \right. \\ &\quad \left. - \frac{\mathcal{G}_1''}{\mathcal{G}_1} \delta_{1i} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{0k}} + \frac{\mathcal{G}_1'^2}{\mathcal{G}_1^2} \delta_{1i} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{0k}} + \frac{\mathcal{G}_1'}{\mathcal{G}_1} \frac{\partial^2 \xi_{1i}}{\partial \alpha_{0j} \partial \alpha_{0k}} + \frac{\mathcal{G}_1''}{\mathcal{G}_1} \delta_{1i} \frac{\partial^2 \xi_{1i}}{\partial \alpha_{0j} \partial \alpha_{0k}} - \frac{\mathcal{G}_1'}{\mathcal{G}_1} \delta_{1i} \frac{\partial^2 \xi_{1i}}{\partial \alpha_{0j} \partial \alpha_{0k}} \right\} \\ &\quad + \sum_{i=1}^n \left\{ \frac{\mathcal{G}_2''}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{0k}} - \frac{\mathcal{G}_2'^2}{\mathcal{G}_2^2} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{0k}} + \frac{\mathcal{G}_2'''}{\mathcal{G}_2} \delta_{2i} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{0k}} - \frac{\mathcal{G}_2''^2}{\mathcal{G}_2^2} \delta_{2i} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{0k}} \right. \\ &\quad \left. - \frac{\mathcal{G}_2''}{\mathcal{G}_2} \delta_{2i} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{0k}} + \frac{\mathcal{G}_2'^2}{\mathcal{G}_2^2} \delta_{2i} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{0k}} + \frac{\mathcal{G}_2'}{\mathcal{G}_2} \frac{\partial^2 \xi_{2i}}{\partial \alpha_{0j} \partial \alpha_{0k}} + \frac{\mathcal{G}_2''}{\mathcal{G}_2} \delta_{2i} \frac{\partial^2 \xi_{2i}}{\partial \alpha_{0j} \partial \alpha_{0k}} - \frac{\mathcal{G}_2'}{\mathcal{G}_2} \delta_{2i} \frac{\partial^2 \xi_{2i}}{\partial \alpha_{0j} \partial \alpha_{0k}} \right\} \\ &= \sum_{i=1}^n \left\{ \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{0k}} \left[\left(\frac{\mathcal{G}_1''}{\mathcal{G}_1} - \frac{\mathcal{G}_1'^2}{\mathcal{G}_1^2} \right) + \delta_{1i} \left(\frac{\mathcal{G}_1'''}{\mathcal{G}_1} - \frac{\mathcal{G}_1''^2}{\mathcal{G}_1^2} - \frac{\mathcal{G}_1''}{\mathcal{G}_1} + \frac{\mathcal{G}_1'^2}{\mathcal{G}_1^2} \right) \right] \right. \\ &\quad \left. + \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{0k}} \left[\left(\frac{\mathcal{G}_2''}{\mathcal{G}_2} - \frac{\mathcal{G}_2'^2}{\mathcal{G}_2^2} \right) + \delta_{2i} \left(\frac{\mathcal{G}_2'''}{\mathcal{G}_2} - \frac{\mathcal{G}_2''^2}{\mathcal{G}_2^2} - \frac{\mathcal{G}_2''}{\mathcal{G}_2} + \frac{\mathcal{G}_2'^2}{\mathcal{G}_2^2} \right) \right] \right\} \\ &= \sum_{i=1}^n \left\{ \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{0k}} \Phi_1 + \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{0k}} \Phi_2 \right\}, \end{aligned}$$

(10)

$$\begin{aligned}
\frac{\partial^2 \ell(\boldsymbol{\alpha})}{\partial \alpha_{0j} \partial \alpha_{1k}} &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}_1''}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} - \frac{\mathcal{G}_1'^2}{\mathcal{G}_1^2} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} + \frac{\mathcal{G}_1'''}{\mathcal{G}_1} \delta_{1i} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} - \frac{\mathcal{G}_1''^2}{\mathcal{G}_1^2} \delta_{1i} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} \right. \\
&\quad \left. - \frac{\mathcal{G}_1''}{\mathcal{G}_1} \delta_{1i} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} + \frac{\mathcal{G}_1'^2}{\mathcal{G}_1^2} \delta_{1i} \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} + \frac{\mathcal{G}_1'}{\mathcal{G}_1} \frac{\partial^2 \xi_{1i}}{\partial \alpha_{0j} \partial \alpha_{1k}} + \frac{\mathcal{G}_1''}{\mathcal{G}_1} \delta_{1i} \frac{\partial^2 \xi_{1i}}{\partial \alpha_{0j} \partial \alpha_{1k}} - \frac{\mathcal{G}_1'}{\mathcal{G}_1} \delta_{1i} \frac{\partial^2 \xi_{1i}}{\partial \alpha_{0j} \partial \alpha_{1k}} \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} \left[\left(\frac{\mathcal{G}_1''}{\mathcal{G}_1} - \frac{\mathcal{G}_1'^2}{\mathcal{G}_1^2} \right) + \delta_{1i} \left(\frac{\mathcal{G}_1'''}{\mathcal{G}_1} - \frac{\mathcal{G}_1''^2}{\mathcal{G}_1^2} - \frac{\mathcal{G}_1''}{\mathcal{G}_1} + \frac{\mathcal{G}_1'}{\mathcal{G}_1^2} \right) \right] \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\partial \xi_{1i}}{\partial \alpha_{0j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} \Phi_1 \right\},
\end{aligned} \tag{11}$$

$$\begin{aligned}
\frac{\partial^2 \ell(\boldsymbol{\alpha})}{\partial \alpha_{0j} \partial \alpha_{2k}} &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}_2''}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} - \frac{\mathcal{G}_2'^2}{\mathcal{G}_2^2} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} + \frac{\mathcal{G}_2'''}{\mathcal{G}_2} \delta_{2i} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} - \frac{\mathcal{G}_2''^2}{\mathcal{G}_2^2} \delta_{2i} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} \right. \\
&\quad \left. - \frac{\mathcal{G}_2''}{\mathcal{G}_2} \delta_{2i} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} + \frac{\mathcal{G}_2'^2}{\mathcal{G}_2^2} \delta_{2i} \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} + \frac{\mathcal{G}_2'}{\mathcal{G}_2} \frac{\partial^2 \xi_{2i}}{\partial \alpha_{0j} \partial \alpha_{2k}} + \frac{\mathcal{G}_2''}{\mathcal{G}_2} \delta_{2i} \frac{\partial^2 \xi_{2i}}{\partial \alpha_{0j} \partial \alpha_{2k}} - \frac{\mathcal{G}_2'}{\mathcal{G}_2} \delta_{2i} \frac{\partial^2 \xi_{2i}}{\partial \alpha_{0j} \partial \alpha_{2k}} \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} \left[\left(\frac{\mathcal{G}_2''}{\mathcal{G}_2} - \frac{\mathcal{G}_2'^2}{\mathcal{G}_2^2} \right) + \delta_{2i} \left(\frac{\mathcal{G}_2'''}{\mathcal{G}_2} - \frac{\mathcal{G}_2''^2}{\mathcal{G}_2^2} - \frac{\mathcal{G}_2''}{\mathcal{G}_2} + \frac{\mathcal{G}_2'}{\mathcal{G}_2^2} \right) \right] \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\partial \xi_{2i}}{\partial \alpha_{0j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} \Phi_2 \right\},
\end{aligned} \tag{12}$$

$$\begin{aligned}
\frac{\partial^2 \ell(\boldsymbol{\alpha})}{\partial \alpha_{1j} \partial \alpha_{1k}} &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}_1''}{\mathcal{G}_1} \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} - \frac{\mathcal{G}_1'^2}{\mathcal{G}_1^2} \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} + \frac{\mathcal{G}_1'''}{\mathcal{G}_1} \delta_{1i} \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} - \frac{\mathcal{G}_1''^2}{\mathcal{G}_1^2} \delta_{1i} \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} \right. \\
&- \frac{\mathcal{G}_1''}{\mathcal{G}_1} \delta_{1i} \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} + \frac{\mathcal{G}_1'^2}{\mathcal{G}_1^2} \delta_{1i} \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} + \frac{\mathcal{G}_1'}{\mathcal{G}_1} \frac{\partial^2 \xi_{1i}}{\partial \alpha_{1j} \partial \alpha_{1k}} + \frac{\mathcal{G}_1''}{\mathcal{G}_1} \delta_{1i} \frac{\partial^2 \xi_{1i}}{\partial \alpha_{1j} \partial \alpha_{1k}} - \frac{\mathcal{G}_1'}{\mathcal{G}_1} \delta_{1i} \frac{\partial^2 \xi_{1i}}{\partial \alpha_{1j} \partial \alpha_{1k}} \\
&+ \left. \frac{\partial^3 \xi_{1i}}{\partial y_i \partial \alpha_{1j} \partial \alpha_{1k}} \delta_{1i} \left(\frac{\partial \xi_{1i}}{\partial y_i} \right)^{-1} - \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \alpha_{1k}} \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \alpha_{1j}} \delta_{1i} \left(\frac{\partial \xi_{1i}}{\partial y_i} \right)^{-2} \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} \left[\left(\frac{\mathcal{G}_1''}{\mathcal{G}_1} - \frac{\mathcal{G}_1'^2}{\mathcal{G}_1^2} \right) + \delta_{1i} \left(\frac{\mathcal{G}_1'''}{\mathcal{G}_1} - \frac{\mathcal{G}_1''^2}{\mathcal{G}_1^2} - \frac{\mathcal{G}_1''}{\mathcal{G}_1} + \frac{\mathcal{G}_1'^2}{\mathcal{G}_1^2} \right) \right] \right. \\
&+ \frac{\partial^2 \xi_{1i}}{\partial \alpha_{1j} \partial \alpha_{1k}} \left[\frac{\mathcal{G}_1'}{\mathcal{G}_1} + \delta_{1i} \left(\frac{\mathcal{G}_1''}{\mathcal{G}_1} - \frac{\mathcal{G}_1'}{\mathcal{G}_1} \right) \right] - \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \alpha_{1k}} \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \alpha_{1j}} \left[\delta_{1i} \left(\frac{\partial \xi_{1i}}{\partial y_i} \right)^{-2} \right] \\
&+ \left. \frac{\partial^3 \xi_{1i}}{\partial y_i \partial \alpha_{1j} \partial \alpha_{1k}} \left[\delta_{1i} \left(\frac{\partial \xi_{1i}}{\partial y_i} \right)^{-1} \right] \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\partial \xi_{1i}}{\partial \alpha_{1j}} \frac{\partial \xi_{1i}}{\partial \alpha_{1k}} \Phi_1 + \frac{\partial^2 \xi_{1i}}{\partial \alpha_{1j} \partial \alpha_{1k}} \Delta_1 - \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \alpha_{1k}} \frac{\partial^2 \xi_{1i}}{\partial y_i \partial \alpha_{1j}} \Psi_1 + \frac{\partial^3 \xi_{1i}}{\partial y_i \partial \alpha_{1j} \partial \alpha_{1k}} \Omega_1 \right\}, \tag{13}
\end{aligned}$$

$$\begin{aligned}
\frac{\partial^2 \ell(\boldsymbol{\alpha})}{\partial \alpha_{2j} \partial \alpha_{2k}} &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}_2''}{\mathcal{G}_2} \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} - \frac{\mathcal{G}_2'^2}{\mathcal{G}_2^2} \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} + \frac{\mathcal{G}_2'''}{\mathcal{G}_2} \delta_{2i} \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} - \frac{\mathcal{G}_2''^2}{\mathcal{G}_2^2} \delta_{2i} \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} \right. \\
&- \frac{\mathcal{G}_2''}{\mathcal{G}_2} \delta_{2i} \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} + \frac{\mathcal{G}_2'^2}{\mathcal{G}_2^2} \delta_{2i} \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} + \frac{\mathcal{G}_2'}{\mathcal{G}_2} \frac{\partial^2 \xi_{2i}}{\partial \alpha_{2j} \partial \alpha_{2k}} + \frac{\mathcal{G}_2''}{\mathcal{G}_2} \delta_{2i} \frac{\partial^2 \xi_{2i}}{\partial \alpha_{2j} \partial \alpha_{2k}} - \frac{\mathcal{G}_2'}{\mathcal{G}_2} \delta_{2i} \frac{\partial^2 \xi_{2i}}{\partial \alpha_{2j} \partial \alpha_{2k}} \\
&+ \left. \frac{\partial^3 \xi_{2i}}{\partial y_i \partial \alpha_{2j} \partial \alpha_{2k}} \delta_{2i} \left(\frac{\partial \xi_{2i}}{\partial y_i} \right)^{-1} - \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \alpha_{2k}} \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \alpha_{2j}} \delta_{2i} \left(\frac{\partial \xi_{2i}}{\partial y_i} \right)^{-2} \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} \left[\left(\frac{\mathcal{G}_2''}{\mathcal{G}_2} - \frac{\mathcal{G}_2'^2}{\mathcal{G}_2^2} \right) + \delta_{2i} \left(\frac{\mathcal{G}_2'''}{\mathcal{G}_2} - \frac{\mathcal{G}_2''^2}{\mathcal{G}_2^2} - \frac{\mathcal{G}_2''}{\mathcal{G}_2} + \frac{\mathcal{G}_2'^2}{\mathcal{G}_2^2} \right) \right] \right. \\
&+ \frac{\partial^2 \xi_{2i}}{\partial \alpha_{2j} \partial \alpha_{2k}} \left[\frac{\mathcal{G}_2'}{\mathcal{G}_2} + \delta_{2i} \left(\frac{\mathcal{G}_2''}{\mathcal{G}_2} - \frac{\mathcal{G}_2'}{\mathcal{G}_2} \right) \right] - \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \alpha_{2k}} \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \alpha_{2j}} \left[\delta_{2i} \left(\frac{\partial \xi_{2i}}{\partial y_i} \right)^{-2} \right] \\
&+ \left. \frac{\partial^3 \xi_{2i}}{\partial y_i \partial \alpha_{2j} \partial \alpha_{2k}} \left[\delta_{2i} \left(\frac{\partial \xi_{2i}}{\partial y_i} \right)^{-1} \right] \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\partial \xi_{2i}}{\partial \alpha_{2j}} \frac{\partial \xi_{2i}}{\partial \alpha_{2k}} \Phi_2 + \frac{\partial^2 \xi_{2i}}{\partial \alpha_{2j} \partial \alpha_{2k}} \Delta_2 - \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \alpha_{2k}} \frac{\partial^2 \xi_{2i}}{\partial y_i \partial \alpha_{2j}} \Psi_2 + \frac{\partial^3 \xi_{2i}}{\partial y_i \partial \alpha_{2j} \partial \alpha_{2k}} \Omega_2 \right\}, \tag{14}
\end{aligned}$$

where $\Phi_\nu = \delta_{\nu i} \left(\frac{\mathcal{G}_\nu'''}{\mathcal{G}_\nu} - \frac{\mathcal{G}_\nu''^2}{\mathcal{G}_\nu^2} - \frac{\mathcal{G}_\nu''}{\mathcal{G}_\nu} + \frac{\mathcal{G}_\nu'^2}{\mathcal{G}_\nu^2} \right)$ and $\Psi_\nu = \left[\delta_{\nu i} \left(\frac{\partial \xi_{\nu i}}{\partial y_i} \right)^{-2} \right]$. Collecting the last terms of (10), (11), (12), (13) and (14), we obtain

$$\frac{\partial^2 \ell(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}_0 \partial \boldsymbol{\alpha}_0^\top} = \sum_{i=1}^n \left\{ \Phi_1 \frac{\partial \xi_{1i}}{\partial \boldsymbol{\alpha}_0} \left[\frac{\partial \xi_{1i}}{\partial \boldsymbol{\alpha}_0} \right]^\top + \Phi_2 \frac{\partial \xi_{2i}}{\partial \boldsymbol{\alpha}_0} \left[\frac{\partial \xi_{2i}}{\partial \boldsymbol{\alpha}_0} \right]^\top \right\},$$

$$\frac{\partial^2 \ell(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}_0 \partial \boldsymbol{\alpha}_\nu^\top} = \sum_{i=1}^n \left\{ \Phi_\nu \frac{\partial \xi_{\nu i}}{\partial \boldsymbol{\alpha}_0} \left[\frac{\partial \xi_{\nu i}}{\partial \boldsymbol{\alpha}_0} \right]^\top \right\},$$

$$\frac{\partial^2 \ell(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}_\nu \partial \boldsymbol{\alpha}_\nu^\top} = \sum_{i=1}^n \left\{ \Phi_\nu \frac{\partial \xi_{\nu i}}{\partial \boldsymbol{\alpha}_\nu} \left[\frac{\partial \xi_{\nu i}}{\partial \boldsymbol{\alpha}_\nu} \right]^\top + \Delta_\nu \frac{\partial^2 \xi_{\nu i}}{\partial \boldsymbol{\alpha}_\nu \partial \boldsymbol{\alpha}_\nu^\top} - \Psi_\nu \frac{\partial^2 \xi_{\nu i}}{\partial y_i \partial \boldsymbol{\alpha}_\nu} \left[\frac{\partial^2 \xi_{\nu i}}{\partial y_i \partial \boldsymbol{\alpha}_\nu} \right]^\top + \Omega_\nu \frac{\partial^3 \xi_{\nu i}}{\partial y_i \partial \boldsymbol{\alpha}_\nu \partial \boldsymbol{\alpha}_\nu^\top} \right\},$$

where

$$\frac{\partial^2 \xi_{\nu i}}{\partial \boldsymbol{\alpha}_\nu \partial \boldsymbol{\alpha}_\nu^\top} = \begin{bmatrix} \frac{\partial^2 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial \alpha_{\nu 11} \partial \alpha_{\nu 11}} & \cdots & \frac{\partial^2 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial \alpha_{\nu 11} \partial \alpha_{\nu Q_\nu J_\nu Q_\nu}} \\ \cdots & \ddots & \cdots \\ \frac{\partial^2 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial \alpha_{\nu Q_\nu J_\nu Q_\nu} \partial \alpha_{\nu 11}} & \cdots & \frac{\partial^2 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial \alpha_{\nu Q_\nu J_\nu Q_\nu} \partial \alpha_{\nu Q_\nu J_\nu Q_\nu}} \end{bmatrix},$$

$$\frac{\partial^3 \xi_{\nu i}}{\partial y_i \partial \boldsymbol{\alpha}_\nu \partial \boldsymbol{\alpha}_\nu^\top} = \begin{bmatrix} \frac{\partial^3 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial y_i \partial \alpha_{\nu 11} \partial \alpha_{\nu 11}} & \cdots & \frac{\partial^3 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial y_i \partial \alpha_{\nu 11} \partial \alpha_{\nu Q_\nu J_\nu Q_\nu}} \\ \cdots & \ddots & \cdots \\ \frac{\partial^3 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial y_i \partial \alpha_{\nu Q_\nu J_\nu Q_\nu} \partial \alpha_{\nu 11}} & \cdots & \frac{\partial^3 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial y_i \partial \alpha_{\nu Q_\nu J_\nu Q_\nu} \partial \alpha_{\nu Q_\nu J_\nu Q_\nu}} \end{bmatrix}.$$

In particular, the design sub-matrices of $\frac{\partial^2 \xi_{\nu i}}{\partial \boldsymbol{\alpha}_\nu \partial \boldsymbol{\alpha}_\nu^\top}$ and $\frac{\partial^3 \xi_{\nu i}}{\partial y_i \partial \boldsymbol{\alpha}_\nu \partial \boldsymbol{\alpha}_\nu^\top}$ are calculated using

$$\frac{\partial^2 \xi_{\nu i}(\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_\nu)}{\partial \alpha_{\nu q_\nu} \partial \alpha_{\nu s_\nu}^\top} = \begin{cases} \mathcal{Q}_{\nu 0}^{\Delta \Delta}(y_i) & \text{if } \alpha_{\nu q_\nu} = \alpha_{\nu s_\nu} = \alpha_{\nu 0} \\ 0 & \text{otherwise,} \end{cases}$$

$$\frac{\partial^3 \xi_{\nu i}(\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_\nu)}{\partial y_i \partial \alpha_{\nu q_\nu} \alpha_{\nu s_\nu}^\top} = \begin{cases} \mathcal{Q}_{\nu 0}^{\Delta \Delta'}(y_i) & \text{if } \alpha_{\nu q_\nu} = \alpha_{\nu s_\nu} = \alpha_{\nu 0} \\ 0 & \text{otherwise,} \end{cases}$$

where $\mathcal{Q}'_{\nu 0}{}^{\Delta\Delta}(y_i)$ and $\mathcal{Q}'_{\nu 0}{}^{\Delta\Delta'}(y_i)$ are defined as

$$\mathcal{Q}'_{\nu 0}{}^{\Delta\Delta}(y_i) = \begin{cases} \frac{\partial^2 \xi_{\nu i}}{\partial \alpha_{\nu 0 j_{\nu 0}} \partial \alpha_{\nu 0 k_{\nu 0}}} = [\sum_{j_{\nu 0}}^{J_{\nu 0}} \mathcal{Q}_{\nu 0 j_{\nu 0}}(y_i)] \exp(\alpha_{\nu 0 j_{\nu 0}}) & \text{if } j = k \neq 1 \\ \frac{\partial^2 \xi_{\nu i}}{\partial \alpha_{\nu 0 j_{\nu 0}} \partial \alpha_{\nu 0 k_{\nu 0}}} = 0 & \text{otherwise,} \end{cases}$$

$$\mathcal{Q}'_{\nu 0}{}^{\Delta\Delta'}(y_i) = \begin{cases} \frac{\partial^3 \xi_{\nu i}}{\partial y_i \alpha_{\nu 0 j_{\nu 0}} \partial \alpha_{\nu 0 k_{\nu 0}}} = [\sum_{j_{\nu 0}}^{J_{\nu 0}} \mathcal{Q}'_{\nu 0 j_{\nu 0}}(y_i)] \exp(\alpha_{\nu 0 j_{\nu 0}}) & \text{if } j = k \neq 1 \\ \frac{\partial^3 \xi_{\nu i}}{\partial y_i \alpha_{\nu 0 j_{\nu 0}} \partial \alpha_{\nu 0 k_{\nu 0}}} = 0 & \text{otherwise.} \end{cases}$$

On the other hand, the non-informative penalized Hessian is

$$\nabla_{\gamma\gamma} \ell_p(\gamma) = \nabla_{\gamma\gamma} \ell(\gamma) - \mathcal{S}.$$

Since $\xi_{1i}(\gamma_1)$ and $\xi_{2i}(\gamma_2)$ do not have parameters in common, $\nabla_{\gamma\gamma} \ell(\gamma)$ can be written as

$$\nabla_{\gamma\gamma} \ell(\gamma) = \begin{bmatrix} \nabla_{\gamma_1 \gamma_1} \ell(\gamma) & \mathbf{0} \\ \mathbf{0} & \nabla_{\gamma_2 \gamma_2} \ell(\gamma) \end{bmatrix},$$

where $\nabla_{\gamma_\nu \gamma_\nu} \ell(\gamma) = \frac{\partial^2 \ell(\gamma)}{\partial \gamma_\nu \partial \gamma_\nu^\top}$. This expression can be obtained using

$$\nabla_{\gamma_\nu \gamma_\nu} \ell(\gamma) = \begin{bmatrix} \frac{\partial^2 \ell(\gamma)}{\partial \gamma_{\nu 11} \partial \gamma_{\nu 11}} & \cdots & \frac{\partial^2 \ell(\gamma)}{\partial \gamma_{\nu 11} \partial \gamma_{\nu K_\nu J_\nu K_\nu}} \\ \cdots & \ddots & \cdots \\ \frac{\partial^2 \ell(\gamma)}{\partial \gamma_{\nu K_\nu J_\nu K_\nu} \partial \gamma_{\nu 11}} & \cdots & \frac{\partial^2 \ell(\gamma)}{\partial \gamma_{\nu K_\nu J_\nu K_\nu} \partial \gamma_{\nu K_\nu J_\nu K_\nu}} \end{bmatrix}.$$

Furthermore, the scalar derivatives of $\nabla_{\gamma_1\gamma_1}\ell(\gamma)$ and $\nabla_{\gamma_2\gamma_2}\ell(\gamma)$ are

$$\begin{aligned}
\frac{\partial^2\ell(\gamma)}{\partial\gamma_{1j}\partial\gamma_{1k}} &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial\xi_{1i}}{\partial\gamma_{1j}} \frac{\partial\xi_{1i}}{\partial\gamma_{1k}} - \frac{\mathcal{G}'_1{}^2}{\mathcal{G}_1^2} \frac{\partial\xi_{1i}}{\partial\gamma_{1j}} \frac{\partial\xi_{1i}}{\partial\gamma_{1k}} + \frac{\mathcal{G}'_1{}^3}{\mathcal{G}_1} \delta_{1i} \frac{\partial\xi_{1i}}{\partial\gamma_{1j}} \frac{\partial\xi_{1i}}{\partial\gamma_{1k}} - \frac{\mathcal{G}'_1{}^2}{\mathcal{G}_1^2} \delta_{1i} \frac{\partial\xi_{1i}}{\partial\gamma_{1j}} \frac{\partial\xi_{1i}}{\partial\gamma_{1k}} \right. \\
&- \frac{\mathcal{G}'_1}{\mathcal{G}_1} \delta_{1i} \frac{\partial\xi_{1i}}{\partial\gamma_{1j}} \frac{\partial\xi_{1i}}{\partial\gamma_{1k}} + \frac{\mathcal{G}'_1{}^2}{\mathcal{G}_1^2} \delta_{1i} \frac{\partial\xi_{1i}}{\partial\gamma_{1j}} \frac{\partial\xi_{1i}}{\partial\gamma_{1k}} + \frac{\mathcal{G}'_1}{\mathcal{G}_1} \frac{\partial^2\xi_{1i}}{\partial\gamma_{1j}\partial\gamma_{1k}} + \frac{\mathcal{G}'_1}{\mathcal{G}_1} \delta_{1i} \frac{\partial^2\xi_{1i}}{\partial\gamma_{1j}\partial\gamma_{1k}} - \frac{\mathcal{G}'_1}{\mathcal{G}_1} \delta_{1i} \frac{\partial^2\xi_{1i}}{\partial\gamma_{1j}\partial\gamma_{1k}} \\
&+ \left. \frac{\partial^3\xi_{1i}}{\partial y_i \partial\gamma_{1j} \partial\gamma_{1k}} \delta_{1i} \left(\frac{\partial\xi_{1i}}{\partial y_i} \right)^{-1} - \frac{\partial^2\xi_{1i}}{\partial y_i \partial\gamma_{1k}} \frac{\partial^2\xi_{1i}}{\partial y_i \partial\gamma_{1j}} \delta_{1i} \left(\frac{\partial\xi_{1i}}{\partial y_i} \right)^{-2} \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\partial\xi_{1i}}{\partial\gamma_{1j}} \frac{\partial\xi_{1i}}{\partial\gamma_{1k}} \left[\left(\frac{\mathcal{G}'_1}{\mathcal{G}_1} - \frac{\mathcal{G}'_1{}^2}{\mathcal{G}_1^2} \right) + \delta_{1i} \left(\frac{\mathcal{G}'_1{}^3}{\mathcal{G}_1} - \frac{\mathcal{G}'_1{}^2}{\mathcal{G}_1^2} - \frac{\mathcal{G}'_1}{\mathcal{G}_1} + \frac{\mathcal{G}'_1{}^2}{\mathcal{G}_1^2} \right) \right] \right. \\
&+ \frac{\partial^2\xi_{1i}}{\partial\gamma_{1j}\partial\gamma_{1k}} \left[\frac{\mathcal{G}'_1}{\mathcal{G}_1} + \delta_{1i} \left(\frac{\mathcal{G}'_1}{\mathcal{G}_1} - \frac{\mathcal{G}'_1}{\mathcal{G}_1} \right) \right] - \frac{\partial^2\xi_{1i}}{\partial y_i \partial\gamma_{1k}} \frac{\partial^2\xi_{1i}}{\partial y_i \partial\gamma_{1j}} \left[\delta_{1i} \left(\frac{\partial\xi_{1i}}{\partial y_i} \right)^{-2} \right] \\
&+ \left. \frac{\partial^3\xi_{1i}}{\partial y_i \partial\gamma_{1j} \partial\gamma_{1k}} \left[\delta_{1i} \left(\frac{\partial\xi_{1i}}{\partial y_i} \right)^{-1} \right] \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\partial\xi_{1i}}{\partial\gamma_{1j}} \frac{\partial\xi_{1i}}{\partial\gamma_{1k}} \Phi_1 + \frac{\partial^2\xi_{1i}}{\partial\gamma_{1j}\partial\gamma_{1k}} \Delta_1 - \frac{\partial^2\xi_{1i}}{\partial y_i \partial\gamma_{1k}} \frac{\partial^2\xi_{1i}}{\partial y_i \partial\gamma_{1j}} \Psi_1 + \frac{\partial^3\xi_{1i}}{\partial y_i \partial\gamma_{1j} \partial\gamma_{1k}} \Omega_1 \right\}, \tag{15}
\end{aligned}$$

$$\begin{aligned}
\frac{\partial^2\ell(\gamma)}{\partial\gamma_{2j}\partial\gamma_{2k}} &= \sum_{i=1}^n \left\{ \frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial\xi_{2i}}{\partial\gamma_{2j}} \frac{\partial\xi_{2i}}{\partial\gamma_{2k}} - \frac{\mathcal{G}'_2{}^2}{\mathcal{G}_2^2} \frac{\partial\xi_{2i}}{\partial\gamma_{2j}} \frac{\partial\xi_{2i}}{\partial\gamma_{2k}} + \frac{\mathcal{G}'_2{}^3}{\mathcal{G}_2} \delta_{2i} \frac{\partial\xi_{2i}}{\partial\gamma_{2j}} \frac{\partial\xi_{2i}}{\partial\gamma_{2k}} - \frac{\mathcal{G}'_2{}^2}{\mathcal{G}_2^2} \delta_{2i} \frac{\partial\xi_{2i}}{\partial\gamma_{2j}} \frac{\partial\xi_{2i}}{\partial\gamma_{2k}} \right. \\
&- \frac{\mathcal{G}'_2}{\mathcal{G}_2} \delta_{2i} \frac{\partial\xi_{2i}}{\partial\gamma_{2j}} \frac{\partial\xi_{2i}}{\partial\gamma_{2k}} + \frac{\mathcal{G}'_2{}^2}{\mathcal{G}_2^2} \delta_{2i} \frac{\partial\xi_{2i}}{\partial\gamma_{2j}} \frac{\partial\xi_{2i}}{\partial\gamma_{2k}} + \frac{\mathcal{G}'_2}{\mathcal{G}_2} \frac{\partial^2\xi_{2i}}{\partial\gamma_{2j}\partial\gamma_{2k}} + \frac{\mathcal{G}'_2}{\mathcal{G}_2} \delta_{2i} \frac{\partial^2\xi_{2i}}{\partial\gamma_{2j}\partial\gamma_{2k}} - \frac{\mathcal{G}'_2}{\mathcal{G}_2} \delta_{2i} \frac{\partial^2\xi_{2i}}{\partial\gamma_{2j}\partial\gamma_{2k}} \\
&+ \left. \frac{\partial^3\xi_{2i}}{\partial y_i \partial\gamma_{2j} \partial\gamma_{2k}} \delta_{2i} \left(\frac{\partial\xi_{2i}}{\partial y_i} \right)^{-1} - \frac{\partial^2\xi_{2i}}{\partial y_i \partial\gamma_{2k}} \frac{\partial^2\xi_{2i}}{\partial y_i \partial\gamma_{2j}} \delta_{2i} \left(\frac{\partial\xi_{2i}}{\partial y_i} \right)^{-2} \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\partial\xi_{2i}}{\partial\gamma_{2j}} \frac{\partial\xi_{2i}}{\partial\gamma_{2k}} \left[\left(\frac{\mathcal{G}'_2}{\mathcal{G}_2} - \frac{\mathcal{G}'_2{}^2}{\mathcal{G}_2^2} \right) + \delta_{2i} \left(\frac{\mathcal{G}'_2{}^3}{\mathcal{G}_2} - \frac{\mathcal{G}'_2{}^2}{\mathcal{G}_2^2} - \frac{\mathcal{G}'_2}{\mathcal{G}_2} + \frac{\mathcal{G}'_2{}^2}{\mathcal{G}_2^2} \right) \right] \right. \\
&+ \frac{\partial^2\xi_{2i}}{\partial\gamma_{2j}\partial\gamma_{2k}} \left[\frac{\mathcal{G}'_2}{\mathcal{G}_2} + \delta_{2i} \left(\frac{\mathcal{G}'_2}{\mathcal{G}_2} - \frac{\mathcal{G}'_2}{\mathcal{G}_2} \right) \right] - \frac{\partial^2\xi_{2i}}{\partial y_i \partial\gamma_{2k}} \frac{\partial^2\xi_{2i}}{\partial y_i \partial\gamma_{2j}} \left[\delta_{2i} \left(\frac{\partial\xi_{2i}}{\partial y_i} \right)^{-2} \right] \\
&+ \left. \frac{\partial^3\xi_{2i}}{\partial y_i \partial\gamma_{2j} \partial\gamma_{2k}} \left[\delta_{2i} \left(\frac{\partial\xi_{2i}}{\partial y_i} \right)^{-1} \right] \right\} \\
&= \sum_{i=1}^n \left\{ \frac{\partial\xi_{2i}}{\partial\gamma_{2j}} \frac{\partial\xi_{2i}}{\partial\gamma_{2k}} \Phi_2 + \frac{\partial^2\xi_{2i}}{\partial\gamma_{2j}\partial\gamma_{2k}} \Delta_2 - \frac{\partial^2\xi_{2i}}{\partial y_i \partial\gamma_{2k}} \frac{\partial^2\xi_{2i}}{\partial y_i \partial\gamma_{2j}} \Psi_2 + \frac{\partial^3\xi_{2i}}{\partial y_i \partial\gamma_{2j} \partial\gamma_{2k}} \Omega_2 \right\}. \tag{16}
\end{aligned}$$

The last terms of equations (15) and (16) allow to express $\nabla_{\gamma_1 \gamma_1} \ell(\gamma)$ and $\nabla_{\gamma_2 \gamma_2} \ell(\gamma)$ as

$$\nabla_{\gamma_\nu \gamma_\nu} \ell(\gamma) = \sum_{i=1}^n \left\{ \Phi_{\nu i} \frac{\partial \xi_{\nu i}}{\partial \gamma_\nu} \left[\frac{\partial \xi_{\nu i}}{\partial \gamma_\nu} \right]^\top + \Delta_{\nu i} \frac{\partial^2 \xi_{\nu i}}{\partial \gamma_\nu \partial \gamma_\nu^\top} - \Psi_{\nu i} \frac{\partial^2 \xi_{\nu i}}{\partial y_i \partial \gamma_\nu} \left[\frac{\partial^2 \xi_{\nu i}}{\partial y_i \partial \gamma_\nu} \right]^\top + \Omega_{\nu i} \frac{\partial^3 \xi_{\nu i}}{\partial y_i \partial \gamma_\nu \partial \gamma_\nu^\top} \right\},$$

where

$$\frac{\partial^2 \xi_{\nu i}}{\partial \gamma_\nu \partial \gamma_\nu^\top} = \begin{bmatrix} \frac{\partial^2 \xi_{\nu i}(\gamma_\nu)}{\partial \gamma_{\nu 11} \partial \gamma_{\nu 11}} & \cdots & \frac{\partial^2 \xi_{\nu i}(\gamma_\nu)}{\partial \gamma_{\nu 11} \partial \gamma_{\nu K_\nu J_\nu K_\nu}} \\ \cdots & \ddots & \cdots \\ \frac{\partial^2 \xi_{\nu i}(\gamma_\nu)}{\partial \gamma_{\nu K_\nu J_\nu K_\nu} \partial \gamma_{\nu 11}} & \cdots & \frac{\partial^2 \xi_{\nu i}(\gamma_\nu)}{\partial \gamma_{\nu K_\nu J_\nu K_\nu} \partial \gamma_{\nu K_\nu J_\nu K_\nu}} \end{bmatrix},$$

$$\frac{\partial^3 \xi_{\nu i}}{\partial y_i \partial \gamma_\nu \partial \gamma_\nu^\top} = \begin{bmatrix} \frac{\partial^3 \xi_{\nu i}(\gamma_\nu)}{\partial y_i \partial \gamma_{\nu 11} \partial \gamma_{\nu 11}} & \cdots & \frac{\partial^3 \xi_{\nu i}(\gamma_\nu)}{\partial y_i \partial \gamma_{\nu 11} \partial \gamma_{\nu K_\nu J_\nu K_\nu}} \\ \cdots & \ddots & \cdots \\ \frac{\partial^3 \xi_{\nu i}(\gamma_\nu)}{\partial y_i \partial \gamma_{\nu K_\nu J_\nu K_\nu} \partial \gamma_{\nu 11}} & \cdots & \frac{\partial^3 \xi_{\nu i}(\gamma_\nu)}{\partial y_i \partial \gamma_{\nu K_\nu J_\nu K_\nu} \partial \gamma_{\nu K_\nu J_\nu K_\nu}} \end{bmatrix}.$$

In addition, the design sub-matrices of $\frac{\partial^2 \xi_{\nu i}(\gamma_\nu)}{\partial \gamma_\nu \partial \gamma_\nu^\top}$ and $\frac{\partial^3 \xi_{\nu i}(\gamma_\nu)}{\partial y_i \partial \gamma_\nu \partial \gamma_\nu^\top}$ can be obtained using the following equations

$$\frac{\partial^2 \xi_{\nu i}(\gamma_\nu)}{\partial \gamma_{\nu k_\nu} \partial \gamma_{\nu s_\nu}^\top} = \begin{cases} \mathcal{Q}_{\nu 0}^{\Delta \Delta}(y_i) & \text{if } \gamma_{\nu k_\nu} = \gamma_{\nu s_\nu} = \gamma_{\nu 0} \\ \mathbf{0} & \text{otherwise,} \end{cases}$$

$$\frac{\partial^3 \xi_{\nu i}(\gamma_\nu)}{\partial y_i \partial \gamma_{\nu k_\nu} \partial \gamma_{\nu s_\nu}^\top} = \begin{cases} \mathcal{Q}_{\nu 0}^{\Delta \Delta'}(y_i) & \text{if } \gamma_{\nu k_\nu} = \gamma_{\nu s_\nu} = \gamma_{\nu 0} \\ \mathbf{0} & \text{otherwise,} \end{cases}$$

where $\mathcal{Q}_{\nu 0}^{\Delta\Delta}(y_i)$ and $\mathcal{Q}_{\nu 0}^{\Delta\Delta'}(y_i)$ can be calculated as

$$\mathcal{Q}_{\nu 0}^{\Delta\Delta}(y_i) = \begin{cases} \frac{\partial^2 \xi_{\nu i}}{\partial \gamma_{\nu 0 j_{\nu 0}} \partial \gamma_{\nu 0 k_{\nu 0}}} = [\sum_{j_{\nu 0}}^{J_{\nu 0}} \mathcal{Q}_{\nu 0 j_{\nu 0}}(y_i)] \exp(\gamma_{\nu 0 j_{\nu 0}}) & \text{if } j = k \neq 1 \\ \frac{\partial^2 \xi_{\nu i}}{\partial \gamma_{\nu 0 j_{\nu 0}} \partial \gamma_{\nu 0 k_{\nu 0}}} = 0 & \text{otherwise,} \end{cases}$$

$$\mathcal{Q}_{\nu 0}^{\Delta\Delta'}(y_i) = \begin{cases} \frac{\partial^3 \xi_{\nu i}}{\partial y_i \gamma_{\nu 0 j_{\nu 0}} \partial \gamma_{\nu 0 k_{\nu 0}}} = [\sum_{j_{\nu 0}}^{J_{\nu 0}} \mathcal{Q}'_{\nu 0 j_{\nu 0}}(y_i)] \exp(\gamma_{\nu 0 j_{\nu 0}}) & \text{if } j = k \neq 1 \\ \frac{\partial^3 \xi_{\nu i}}{\partial y_i \gamma_{\nu 0 j_{\nu 0}} \partial \gamma_{\nu 0 k_{\nu 0}}} = 0 & \text{otherwise.} \end{cases}$$

Supplementary Material C: Estimation Algorithm

The optimization method used is the trust region algorithm. At iteration a , for a given vector α and maintaining λ fixed at a vector of values, equation (13) in the main paper (or generally, any of the models' likelihoods considered in the paper) is maximized using

$$\alpha^{[a+1]} = \arg \min_{\varepsilon: \|\varepsilon\| \leq \Xi^{[a]}} \bar{\ell}_p(\alpha^{[a]}),$$

where $\bar{\ell}_p(\alpha^{[a]}) = -\{\ell_p(\alpha^{[a]}) + \varepsilon^\top \mathbf{g}_p(\alpha^{[a]}) + \frac{1}{2} \varepsilon^\top \mathcal{H}_p(\alpha^{[a]}) \varepsilon\}$, $\mathbf{g}_p(\alpha^{[a]}) = \mathbf{g}(\alpha^{[a]}) - \mathcal{S}\alpha^{[a]}$, $\mathcal{H}_p(\alpha^{[a]}) = \mathcal{H}(\alpha^{[a]}) - \mathcal{S}$. Vector $\mathbf{g}(\alpha^{[a]})$ consists of $\mathbf{g}_0(\alpha^{[a]}) = \nabla_{\alpha_0} \ell(\alpha)|_{\alpha_0=\alpha_0^{[a]}}$ and $\mathbf{g}_\nu(\alpha^{[a]}) = \nabla_{\alpha_\nu} \ell(\alpha)|_{\alpha_\nu=\alpha_\nu^{[a]}}$, and $\mathcal{H}(\alpha^{[a]})_{l,j} = \nabla_{\alpha_l \alpha_j} \ell(\alpha)|_{\alpha_l=\alpha_l^{[a]}, \alpha_j=\alpha_j^{[a]}}$, where $l, j = 0, 1, 2$ and $\nu = 1, 2$. The euclidean norm is denoted by $\|\cdot\|$, and the radius of the trust region is represented by $\Xi^{[a]}$ which is adjusted through the iterations. Close to the solution, the trust region algorithms behaves as a classic Newton-Raphson unconstrained method (Nocedal & Wright, 2006).

Estimation of λ is achieved by adapting the general and automatic multiple smoothing parameter estimation method of (Marra et al., 2017) to the context of the proposed survival models. The smoothing criterion is based on the knowledge of $\mathbf{g}(\alpha)$ and $\mathcal{H}(\alpha)$. The main ideas and some useful results are given here.

To simplify the notation, $\mathbf{g}_p(\alpha^{[a]})$, $\mathbf{g}(\alpha^{[a]})$, $\mathcal{H}_p(\alpha^{[a]})$ and $\mathcal{H}(\alpha^{[a]})$ are denoted as $\mathbf{g}_p^{[a]}$, $\mathbf{g}^{[a]}$, $\mathcal{H}_p^{[a]}$ and $\mathcal{H}^{[a]}$. First, it is necessary to express the parameter estimator in terms of $\mathbf{g}_p^{[a]}$ and $\mathcal{H}_p^{[a]}$. To achieve this, a first order Taylor expansion of $\mathbf{g}_p^{[a+1]}$ about $\alpha^{[a]}$ is used, which yields the following expression: $\mathbf{0} = \mathbf{g}_p^{[a+1]} \approx \mathbf{g}_p^{[a]}(\alpha^{[a+1]} - \alpha^{[a]}) \mathcal{H}_p^{[a]}$. After some manipulations, $\alpha^{[a+1]} = (-\mathcal{H}^{[a]} + \mathcal{S})^{-1} \sqrt{-\mathcal{H}^{[a]}} [\sqrt{-\mathcal{H}^{[a]}} \alpha^{[a]} + \sqrt{-\mathcal{H}^{[a]}}^{-1} \mathbf{g}^{[a]}]$ is obtained, which then becomes $\alpha^{[a+1]} = (-\mathcal{H}^{[a]} + \mathcal{S})^{-1} \sqrt{-\mathcal{H}^{[a]}} \mathcal{Z}^{[a]}$, where $\mathcal{Z}^{[a]} = \mathbf{v}_{\mathcal{Z}}^{[a]} + \boldsymbol{\xi}_{\mathcal{Z}}^{[a]}$, $\mathbf{v}_{\mathcal{Z}}^{[a]} = \sqrt{-\mathcal{H}^{[a]}} \alpha^{[a]}$ and $\boldsymbol{\xi}_{\mathcal{Z}}^{[a]} = \sqrt{-\mathcal{H}^{[a]}}^{-1} \mathbf{g}^{[a]}$. Eigenvalue decomposition is used to obtain the square root of $-\mathcal{H}^{[a]}$ and its inverse. Furthermore, from likelihood theory, $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\mathcal{Z} \sim \mathcal{N}(\mathbf{v}_{\mathcal{Z}}, \mathbf{I})$, where $\mathbf{v}_{\mathcal{Z}} = \sqrt{-\mathcal{H}} \alpha^0$, α^0 is the true parameter vector and \mathbf{I} is the identity matrix. $\hat{\mathbf{v}}_{\mathcal{Z}} = \sqrt{-\mathcal{H}} \hat{\alpha} = \mathcal{B} \mathcal{Z}$ is the predicted value vector for \mathcal{Z} , where $\mathcal{B} = \sqrt{-\mathcal{H}} (-\mathcal{H} + \mathcal{S})^{-1} \sqrt{-\mathcal{H}}$. Since our objective is to estimate λ so that the smooth terms' complexity which is not supported by the data is removed, the following criterion

is used

$$\mathbb{E}(\|\mathbf{v}_{\mathbf{Z}} - \hat{\mathbf{v}}_{\mathbf{Z}}\|^2) = \mathbb{E}(\|\mathbf{Z} - \mathbf{B}\mathbf{Z}\|^2) - \bar{n} + 2\text{tr}(\mathbf{B}), \quad (17)$$

where $\bar{n} = 2n$ and $\text{tr}(\mathbf{B})$ represent the number of effective degrees of freedom of the penalized model. In applications, λ is estimated by minimizing an estimate of equation (17), in other words

$$\|\widehat{\mathbf{v}_{\mathbf{Z}}} - \hat{\mathbf{v}}_{\mathbf{Z}}\|^2 = \|\mathbf{Z} - \mathbf{B}\mathbf{Z}\|^2 - \bar{n} + 2\text{tr}(\mathbf{B}). \quad (18)$$

The RHS of equation (18) depends on λ through \mathbf{B} while \mathbf{Z} is associated with the un-penalized part of the model. Equation (17) is approximately equivalent to the AIC (Akaike, 1973). This implies that λ is estimated by minimizing what is effectively the AIC with number of parameters given by $\text{tr}(\mathbf{B})$. Holding the model's parameter vector value fixed at $\boldsymbol{\alpha}^{[a+1]}$, the following problem

$$\boldsymbol{\lambda}^{[a+1]} = \arg \min_{\boldsymbol{\lambda}} \|\mathbf{Z}^{[a+1]} - \mathbf{B}^{[a+1]}\mathbf{Z}^{[a+1]}\|^2 - \bar{n} + 2\text{tr}(\mathbf{B}^{[a+1]}) \quad (19)$$

is solved using the automatic efficient and stable computational method proposed by Wood (2004). This approach uses the performance iteration idea of Gu (1992), which is based on Newton's method and can evaluate in an efficient and stable way the components in (19) along with their first and second derivatives with respect to $\log(\boldsymbol{\lambda})$, because the smoothing parameters can only take positive values.

The methods for estimating $\boldsymbol{\alpha}$ and $\boldsymbol{\lambda}$ are iterated until the algorithm satisfies the criterion $|\ell(\boldsymbol{\alpha}^{[a+1]}) - \ell(\boldsymbol{\alpha}^{[a]})| / (0.1 + |\ell(\boldsymbol{\alpha}^{[a+1]})|) \leq (1e - 0.7)$. Starting values are obtained by fitting two non-informative models for the survival and censoring times.

Supplementary Material D: Proofs of the Theorems

This section provides the proofs of Theorems 1, 2 and 3 stated in Section 2.4. First, we establish the main set of assumptions (regularity conditions and vanishing penalties), then the main results are presented.

D.1. Assumptions

Since the same set of assumptions are used to proof Theorems 1 and 2, we use $\boldsymbol{\theta}$ to represents the generic vector of parameters. In particular, $\boldsymbol{\theta} = \boldsymbol{\alpha}$ in Theorem 1 and $\boldsymbol{\theta} = \boldsymbol{\gamma}$ in Theorem 2. Hence, the generic log-likelihood function can be written as

$$\ell(\boldsymbol{\theta}) = \sum_{i=1}^n \log \left[[f_1(y_i|\mathbf{z}_i; \boldsymbol{\theta}_1)S_2(y_i|\mathbf{z}_i; \boldsymbol{\theta}_2)]^{\delta_{1i}} [f_2(y_i|\mathbf{z}_i; \boldsymbol{\theta}_2)S_1(y_i|\mathbf{z}_i; \boldsymbol{\theta}_1)]^{\delta_{2i}} \right]. \quad (20)$$

In (20), it has been assumed that $\mathbf{z}_{1i} = \mathbf{z}_{2i}$. In what follows $\ell(\boldsymbol{\theta}) = \sum_{i=1}^n \log \omega(\mathbf{w}_i; \boldsymbol{\theta})$, where $\omega(\mathbf{w}_i; \boldsymbol{\theta}) = \omega(y|\mathbf{z}; \boldsymbol{\theta}) = \left[[f_1(y_i|\mathbf{z}_i; \boldsymbol{\theta}_1)S_2(y_i|\mathbf{z}_i; \boldsymbol{\theta}_2)]^{\delta_{1i}} [f_2(y_i|\mathbf{z}_i; \boldsymbol{\theta}_2)S_1(y_i|\mathbf{z}_i; \boldsymbol{\theta}_1)]^{\delta_{2i}} \right]$ and $\mathbf{w}_i = (y_i, \mathbf{z}_i^\top)^\top \in \mathbb{R}_+ \times \mathbb{R}^p$, and $\mathbb{R}_+ = (0, \infty)$. In addition, $\ell(\mathbf{w}_i; \boldsymbol{\theta}) = \log \omega(\mathbf{w}_i; \boldsymbol{\theta})$, $\ell_n(\boldsymbol{\theta}) = n^{-1} \sum_{i=1}^n \ell(\mathbf{w}_i; \boldsymbol{\theta})$, $\nabla_{\boldsymbol{\theta}} \ell(\mathbf{w}_i; \boldsymbol{\theta}) = \frac{\partial \ell(\mathbf{w}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$, $\nabla_{\boldsymbol{\theta}} \ell_n(\boldsymbol{\theta}) = \frac{\partial \ell_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$, $\nabla_{\boldsymbol{\theta}\boldsymbol{\theta}} \ell(\mathbf{w}_i; \boldsymbol{\theta}) = \frac{\partial^2 \ell(\mathbf{w}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}$ and $\nabla_{\boldsymbol{\theta}\boldsymbol{\theta}} \ell_n(\boldsymbol{\theta}) = \frac{\partial^2 \ell_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}$. The penalised likelihood is $\ell_p(\boldsymbol{\theta}) = \ell_n(\boldsymbol{\theta}) - \frac{1}{2} \boldsymbol{\theta}^\top \boldsymbol{S} \boldsymbol{\theta}$.

Assumption 1 (Regularity Conditions).

- (i) The parameter space $\Theta_{\boldsymbol{\theta}}$ is a compact subset of \mathbb{R}^p .
- (ii) For all \mathbf{w}_i , $\omega(\mathbf{w}_i; \boldsymbol{\theta})$ is continuous in $\boldsymbol{\theta}$. Furthermore, $\omega(\mathbf{w}_i; \boldsymbol{\theta})$ is measurable in \mathbf{w}_i for all $\boldsymbol{\theta} \in \Theta_{\boldsymbol{\theta}}$.
- (iii) Identification condition. $\mathbb{P}[\omega(\mathbf{w}_i; \boldsymbol{\theta}) \neq \omega(\mathbf{w}_i; \boldsymbol{\theta}^*)] > 0$ for all $\boldsymbol{\theta} \neq \boldsymbol{\theta}^* \in \Theta_{\boldsymbol{\theta}}$.
- (iv) Dominance. $\mathbb{E}\{\sup_{\boldsymbol{\theta} \in \Theta_{\boldsymbol{\theta}}} |\ell(\mathbf{w}_i; \boldsymbol{\theta})|\} < \infty$
- (v) The true vector of parameters $\boldsymbol{\theta}^*$ is in the interior of $\Theta_{\boldsymbol{\theta}}$, and Θ_0 is an open neighbourhood around $\boldsymbol{\theta}^*$.
- (vi) For all \mathbf{w}_i , $\omega(\mathbf{w}_i; \boldsymbol{\theta})$ is three times continuously differentiable in $\boldsymbol{\theta}$ in an open neighbourhood around $\boldsymbol{\theta}^*$. That is $\omega(\mathbf{w}_i; \boldsymbol{\theta}) \in \mathcal{C}^3(\Theta_0)$

$$(vii) \int \sup_{\theta \in \Theta_0} \|\nabla_{\theta} \ell(\mathbf{w}_i; \theta)\| d\mathbf{w}_i < \infty \text{ and } \int \sup_{\theta \in \Theta_0} \|\nabla_{\theta\theta} \ell(\mathbf{w}_i; \theta)\| d\mathbf{w}_i < \infty.$$

$$(viii) \text{ For } \theta \in \Theta_0, \mathcal{I}(\theta^*) = \text{Cov}\{\nabla_{\theta} \ell(\mathbf{w}_i; \theta^*)\} = \mathbb{E}\{\{\nabla_{\theta} \ell(\mathbf{w}_i; \theta^*) - \mathbb{E}[\nabla_{\theta} \ell(\mathbf{w}_i; \theta^*)]\}\{\nabla_{\theta} \ell(\mathbf{w}_i; \theta^*) - \mathbb{E}[\nabla_{\theta} \ell(\mathbf{w}_i; \theta^*)]\}^{\top}\} \text{ exists and is positive-definite.}$$

$$(ix) \text{ For all } 1 \leq e, f, h \leq p+1, \text{ there exist a function } \phi : \mathbb{R}_+ \times \mathbb{R}^p \rightarrow \mathbb{R} \text{ such that, for } \theta \in \Theta_0 \text{ and } \mathbf{w}_i \in \mathbb{R}_+ \times \mathbb{R}^p, \left| \frac{\partial^3 \ell(\mathbf{w}_i; \theta)}{\partial \theta_e \partial \theta_f \partial \theta_h} \right| \leq \phi(\mathbf{w}_i), \text{ with } \mathbb{E}[\phi(\mathbf{w}_i)] < \infty.$$

Assumption 2. $\lambda = o(n^{-1/2})$.

In addition, the following lemmas are required to prove Theorems 1, 2 and 3.

Lemma 1. Let $s(\mathbf{w}, \theta)$ be a continuously differentiable function, a.s. $d\mathbf{w}$, on $\theta \in \Theta_0$.

If $\int \sup_{\theta \in \Theta_0} \left\| \frac{\partial s(\mathbf{w}, \theta)}{\partial \theta} \right\| d\mathbf{w} < \infty$, then for $\theta \in \Theta_0$,

$$(i) \int s(\mathbf{w}, \theta) d\mathbf{w} \text{ is continuously differentiable.}$$

$$(ii) \int [\partial s(\mathbf{w}, \theta) / \partial \theta] d\mathbf{w} = \partial [\int s(\mathbf{w}, \theta) d\mathbf{w}] / \partial \theta.$$

Proof. Newey & McFadden (1994, Lemma 3.6). □

Lemma 2. If Assumption 1 hold, then

$$(i) \mathbb{E}[\nabla_{\theta} \ell(\mathbf{w}; \theta^*)] = \mathbf{0}$$

$$(ii) \mathbb{E}[-\nabla_{\theta\theta} \ell(\mathbf{w}; \theta^*)] = \mathcal{I}(\theta^*)$$

Proof.

(i) Since $\omega(y|\mathbf{z}; \theta)$ is a hypothetical density, its integral is unity:

$$\int \omega(y|\mathbf{z}; \theta) dy = 1.$$

This is an identity, valid for any $\theta \in \Theta_{\theta}$. Differentiating both sides of this identity with respect to θ , we obtain

$$\frac{\partial}{\partial \theta} \int \omega(y|\mathbf{z}; \theta) dy = \mathbf{0}.$$

Then, by Assumptions **1**(vi) and **1**(vii), and Lemma 1 (the order of differentiation and integration can be interchanged), the following expression is obtained

$$\frac{\partial}{\partial \boldsymbol{\theta}} \int \omega(y|\mathbf{z}; \boldsymbol{\theta}) dy = \int \frac{\partial}{\partial \boldsymbol{\theta}} \omega(y|\mathbf{z}; \boldsymbol{\theta}) dy. \quad (21)$$

By the definition of the score, we have $\nabla_{\boldsymbol{\theta}} \ell(\mathbf{w}; \boldsymbol{\theta}) \omega(y|\mathbf{z}; \boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \omega(y|\mathbf{z}; \boldsymbol{\theta})$. Substituting into (21), we obtain

$$\int \nabla_{\boldsymbol{\theta}} \ell(\mathbf{w}; \boldsymbol{\theta}) \omega(y|\mathbf{z}; \boldsymbol{\theta}) dy = \mathbf{0}. \quad (22)$$

This holds for any $\boldsymbol{\theta} \in \Theta_0$, in particular, for $\boldsymbol{\theta}^*$. Setting $\boldsymbol{\theta} = \boldsymbol{\theta}^*$, the following equation is obtained

$$\int \nabla_{\boldsymbol{\theta}} \ell(\mathbf{w}; \boldsymbol{\theta}^*) \omega(y|\mathbf{z}; \boldsymbol{\theta}^*) dy = \mathbb{E}[\nabla_{\boldsymbol{\theta}} \ell(\mathbf{w}; \boldsymbol{\theta}^*) | \mathbf{z}] = \mathbf{0}.$$

Then, applying the Law of Total Expectations, we obtain the required result

$$\mathbb{E}[\nabla_{\boldsymbol{\theta}} \ell(\mathbf{w}; \boldsymbol{\theta}^*)] = \mathbb{E}\{\mathbb{E}[\nabla_{\boldsymbol{\theta}} \ell(\mathbf{w}; \boldsymbol{\theta}^*) | \mathbf{z}]\} = \mathbf{0}.$$

(ii) Differentiating both sides of identity (22) and by Assumptions **1**(vi) and **1**(vii), and Lemma 1, we obtain

$$\int \frac{\partial}{\partial \boldsymbol{\theta}^\top} [\nabla_{\boldsymbol{\theta}} \ell(\mathbf{w}; \boldsymbol{\theta}) \omega(y|\mathbf{z}; \boldsymbol{\theta})] dy = \mathbf{0}. \quad (23)$$

The integrand of (23) can be written as $\frac{\partial}{\partial \boldsymbol{\theta}^\top} [\nabla_{\boldsymbol{\theta}} \ell(\mathbf{w}; \boldsymbol{\theta}) \omega(y|\mathbf{z}; \boldsymbol{\theta})] = \nabla_{\boldsymbol{\theta}\boldsymbol{\theta}} \ell(\mathbf{w}; \boldsymbol{\theta}) \omega(y|\mathbf{z}; \boldsymbol{\theta}) + \nabla_{\boldsymbol{\theta}} \ell(\mathbf{w}; \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \ell(\mathbf{w}; \boldsymbol{\theta})^\top \omega(y|\mathbf{z}; \boldsymbol{\theta})$. Substituting into (23), we obtain

$$- \int \nabla_{\boldsymbol{\theta}\boldsymbol{\theta}} \ell(\mathbf{w}; \boldsymbol{\theta}) \omega(y|\mathbf{z}; \boldsymbol{\theta}) dy = \int \nabla_{\boldsymbol{\theta}} \ell(\mathbf{w}; \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \ell(\mathbf{w}; \boldsymbol{\theta})^\top \omega(y|\mathbf{z}; \boldsymbol{\theta}) dy \quad (24)$$

Setting $\boldsymbol{\theta} = \boldsymbol{\theta}^*$, the following equation is obtained

$$\mathbb{E}[-\nabla_{\boldsymbol{\theta}\boldsymbol{\theta}}\ell(\mathbf{w}; \boldsymbol{\theta}^*)|\mathbf{z}] = \mathbb{E}[\nabla_{\boldsymbol{\theta}}\ell(\mathbf{w}; \boldsymbol{\theta}^*)\nabla_{\boldsymbol{\theta}}\ell(\mathbf{w}; \boldsymbol{\theta}^*)^\top|\mathbf{z}].$$

Then, applying the Law of Total Expectations, we obtain the desired result

$$\mathbb{E}\{\mathbb{E}[-\nabla_{\boldsymbol{\theta}\boldsymbol{\theta}}\ell(\mathbf{w}; \boldsymbol{\theta}^*)|\mathbf{z}]\} = \mathbb{E}\{\mathbb{E}[\nabla_{\boldsymbol{\theta}}\ell(\mathbf{w}; \boldsymbol{\theta}^*)\nabla_{\boldsymbol{\theta}}\ell(\mathbf{w}; \boldsymbol{\theta}^*)^\top|\mathbf{z}]\}.$$

$$\mathbb{E}[-\nabla_{\boldsymbol{\theta}\boldsymbol{\theta}}\ell(\mathbf{w}; \boldsymbol{\theta}^*)] = \mathbb{E}[\nabla_{\boldsymbol{\theta}}\ell(\mathbf{w}; \boldsymbol{\theta}^*)\nabla_{\boldsymbol{\theta}}\ell(\mathbf{w}; \boldsymbol{\theta}^*)^\top].$$

$$\mathbb{E}[-\nabla_{\boldsymbol{\theta}\boldsymbol{\theta}}\ell(\mathbf{w}; \boldsymbol{\theta}^*)] = \mathcal{I}(\boldsymbol{\theta}^*)$$

□

Lemma 3. Let $r \in \mathbb{R}_+$, and Θ_r be the surface of a sphere with radius $rn^{-1/2}$ and center $\boldsymbol{\theta}^*$, that is $\Theta_r = \{\boldsymbol{\theta} \in \Theta_{\boldsymbol{\theta}} : \boldsymbol{\theta} = \boldsymbol{\theta}^* + n^{-1/2}\mathbf{r}, \|\mathbf{r}\| = r\}$. For any $\epsilon > 0$, there exist r such that $\mathbb{P}\left(\sup_{\boldsymbol{\theta} \in \Theta_r} \ell_p(\boldsymbol{\theta}) < \ell_p(\boldsymbol{\theta}^*)\right) \geq 1 - \epsilon$, when n is large enough.

Proof. We define $n\ell_p(\boldsymbol{\theta}) - n\ell_p(\boldsymbol{\theta}^*) = n\ell_n(\boldsymbol{\theta}) - n\ell_n(\boldsymbol{\theta}^*) - \frac{n}{2}[\boldsymbol{\theta}^\top \mathbf{S}\boldsymbol{\theta} - \boldsymbol{\theta}^{*\top} \mathbf{S}\boldsymbol{\theta}^*]$. A Third Order Taylor expansion around $\boldsymbol{\theta}^*$ yields

$$\begin{aligned} n\ell_p(\boldsymbol{\theta}) - n\ell_p(\boldsymbol{\theta}^*) &= n\nabla_{\boldsymbol{\theta}}\ell_n(\boldsymbol{\theta}^*)^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \frac{n}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla_{\boldsymbol{\theta}\boldsymbol{\theta}}\ell_n(\boldsymbol{\theta}^*)(\boldsymbol{\theta} - \boldsymbol{\theta}^*) - n\boldsymbol{\theta}^{*\top} \mathbf{S}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \\ &\quad + \frac{n}{6} \sum_e \sum_f \sum_h (\boldsymbol{\theta} - \boldsymbol{\theta}^*)_e (\boldsymbol{\theta} - \boldsymbol{\theta}^*)_f (\boldsymbol{\theta} - \boldsymbol{\theta}^*)_h \frac{\partial^3 \ell_n(\bar{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}_e \partial \boldsymbol{\theta}_f \partial \boldsymbol{\theta}_h} - \frac{n}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \mathbf{S}(\boldsymbol{\theta} - \boldsymbol{\theta}^*). \end{aligned} \tag{25}$$

Let $\boldsymbol{\theta} = \boldsymbol{\theta}^* + n^{-1/2}\mathbf{r} \in \Theta_r$. Then (25) becomes in

$$\begin{aligned} n\ell_p(\boldsymbol{\theta}) - n\ell_p(\boldsymbol{\theta}^*) &= n^{1/2}\nabla_{\boldsymbol{\theta}}\ell_n(\boldsymbol{\theta}^*)^\top \mathbf{r} + \frac{1}{2}\mathbf{r}^\top \nabla_{\boldsymbol{\theta}\boldsymbol{\theta}}\ell_n(\boldsymbol{\theta}^*)\mathbf{r} + \frac{n^{-1/2}}{6} \sum_e \sum_f \sum_h \mathbf{r}_e \mathbf{r}_f \mathbf{r}_h \frac{\partial^3 \ell_n(\bar{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}_e \partial \boldsymbol{\theta}_f \partial \boldsymbol{\theta}_h} \\ &\quad - n^{1/2}\boldsymbol{\theta}^{*\top} \mathbf{S}\mathbf{r} - \frac{1}{2}\mathbf{r}^\top \mathbf{S}\mathbf{r} \\ n\ell_p(\boldsymbol{\theta}) - n\ell_p(\boldsymbol{\theta}^*) &= \sum_{i=1}^5 \mathcal{C}_{in}(\mathbf{r}), \end{aligned}$$

where $\bar{\boldsymbol{\theta}}$ lies between $\boldsymbol{\theta}^*$ and $\boldsymbol{\theta}^* + n^{-1/2}\mathbf{r}$. For the first term, $|\mathcal{C}_{1n}(\mathbf{r})| = \mathcal{O}_p(1)\|\mathbf{r}\|$ since by Lemma 2(i), Assumption 1(vii) and the CLT, $n^{1/2}\nabla_{\boldsymbol{\theta}}\ell_n(\boldsymbol{\theta}^*) \xrightarrow{d} \mathcal{N}[\mathbf{0}, \mathcal{I}(\boldsymbol{\theta}^*)]$. By Lemma

2(ii) and the LLN, $n^{1/2}\nabla_{\theta\theta}\ell_n(\theta^*) \xrightarrow{p} -\mathcal{I}(\theta^*)$, which (by the continuous mapping theorem) yields $\mathcal{C}_{2n}(\mathbf{r}) \xrightarrow{p} -\frac{1}{2}\mathbf{r}^\top\mathcal{I}(\theta^*)\mathbf{r}$. Thus, by Assumption 1(viii), $\mathcal{C}_{2n}(\mathbf{r}) \leq -\frac{1}{2}\zeta_{\min}\|\mathbf{r}\|^2$, where $\zeta_{\min} > 0$ is the smallest eigenvalue of $\mathcal{I}(\theta^*)$. By Assumption 1(ix) and the LLN, $\left|\frac{\partial^3\ell_n(\bar{\theta})}{\partial\theta_e\partial\theta_f\partial\theta_h}\right| \leq \frac{1}{n}\sum_1^n\phi(\mathbf{w}_i) \xrightarrow{p} \mathbb{E}[\phi(\mathbf{w}_i)] < \infty$. This fact and the Cauchy-Schwarz inequality imply that $|\mathcal{C}_{3n}(\mathbf{r})| \xrightarrow{p} 0$. Finally, by Assumption 2 we have that $|\mathcal{C}_{4n}(\mathbf{r})| \xrightarrow{p} 0$ and $|\mathcal{C}_{5n}(\mathbf{r})| \xrightarrow{p} 0$. Therefore, combining all of these results, we have

$$nl_p(\theta) - nl_p(\theta^*) \leq \mathcal{O}_p(1)\|\mathbf{r}\| - \frac{1}{2}\zeta_{\min}\|\mathbf{r}\|^2 \quad (26)$$

for large enough n . Since the choice of θ was arbitrary, (26) becomes in

$$\sup_{\theta \in \Theta_r} nl_p(\theta) - nl_p(\theta^*) \leq \mathcal{C},$$

where $\mathcal{C} = \mathcal{O}_p(1)\|\mathbf{r}\| - \frac{1}{2}\zeta_{\min}\|\mathbf{r}\|^2$. This implies that $\mathbb{P}\left(\sup_{\theta \in \Theta_r} \ell_p(\theta) < \ell_p(\theta^*)\right) \geq \mathbb{P}(\mathcal{C} < 0)$.

Therefore, because for all $\epsilon > 0$, there exists a $\|\mathbf{r}\| \in \mathbb{R}_+$ such that $\mathbb{P}[\mathcal{C} < 0] \geq 1 - \epsilon$, we obtain $\mathbb{P}\left(\sup_{\theta \in \Theta_r} \ell_p(\theta) < \ell_p(\theta^*)\right) \geq 1 - \epsilon$, as required. \square

Lemma 4. (Delta Method). Suppose that θ_n is a sequence of k -dimensional random vectors and θ^* be a constant k -vector such that $\sqrt{n}(\theta_n - \theta^*) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \Omega)$ for some $k \times k$ matrix Ω . Let $g: \mathbb{R}^k \rightarrow \mathbb{R}^l$ be continuously differentiable at θ^* . Then

$$\sqrt{n}(g(\theta_n) - g(\theta^*)) \xrightarrow{d} \mathcal{N}(\mathbf{0}, G\Omega G^\top)$$

where $G = \left.\frac{\partial g(\theta)}{\partial\theta^\top}\right|_{\theta=\theta^*}$ is the $l \times k$ Jacobian matrix.

Proof. Hayashi (2000, Lemma 2.5). \square

D.2. Theorems

Theorem 1 (Asymptotic properties of the IPMLE estimator).

Proof. Under Assumptions 1(i), 1(ii) and Gourieroux & Monfort (1995, Property 24.1), there exists a well defined random variable (measurable function) $\hat{\alpha}$ that solves the optimization problem

in equation (13). Due to Lemma 3, the informative penalized log-likelihood function has a local maximum $\hat{\alpha}$ in the interior of a sphere centered on α^* . Then, $\|\hat{\alpha} - \alpha^*\| = \mathcal{O}_p(n^{-1/2})$, implying that $\hat{\alpha}$ is a \sqrt{n} -consistent estimator. Furthermore, by Assumption 1(iii) and Newey & McFadden (1994, Lemma 2.2), α^* is the unique maximizer of $Q^*(\alpha) = \mathbb{E}[\ell(\mathbf{w}_i; \alpha)]$.

- (i) To prove the asymptotic normality of the informative penalized likelihood estimator, we take the derivative of the log-likelihood function in equation (13) to obtain

$$\mathbf{0} = \nabla_{\alpha} \ell_n(\hat{\alpha}) - \mathcal{S} \hat{\alpha}. \quad (27)$$

Applying a second order Taylor expansion in equation (27) yields

$$\mathbf{0} = \nabla_{\alpha} \ell_n(\alpha^*) - \mathcal{S} \alpha^* + \nabla_{\alpha \alpha} \ell_n(\alpha^*)(\hat{\alpha} - \alpha^*) - \mathcal{S}(\hat{\alpha} - \alpha^*) + \Delta, \quad (28)$$

where the last term is defined as

$$\Delta = \begin{bmatrix} (\hat{\alpha} - \alpha^*)^{\top} [\nabla^2 \nabla_{\alpha} \ell_n(\bar{\alpha})]_1 (\hat{\alpha} - \alpha^*) \\ \vdots \\ (\hat{\alpha} - \alpha^*)^{\top} [\nabla^2 \nabla_{\alpha} \ell_n(\bar{\alpha})]_p (\hat{\alpha} - \alpha^*) \end{bmatrix}, \quad (29)$$

and $\bar{\alpha}$ lies between α^* and $\hat{\alpha}$, therefore $\|\bar{\alpha} - \alpha^*\| \leq \|\hat{\alpha} - \alpha^*\|$. We can rewrite equation (28) to obtain

$$\mathbf{0} = \nabla_{\alpha} \ell_n(\alpha^*) - \mathcal{S} \alpha^* + \nabla_{\alpha \alpha} \ell_n(\alpha^*)(\hat{\alpha} - \alpha^*) - \mathcal{S}(\hat{\alpha} - \alpha^*) + \Delta_p(\hat{\alpha} - \alpha^*), \quad (30)$$

where Δ_p is defined as

$$\Delta_p = \begin{bmatrix} (\hat{\alpha} - \alpha^*)^{\top} [\nabla \nabla_{\alpha \alpha} \ell_n(\bar{\alpha})]_1 \\ \vdots \\ (\hat{\alpha} - \alpha^*)^{\top} [\nabla \nabla_{\alpha \alpha} \ell_n(\bar{\alpha})]_p \end{bmatrix}.$$

Multiplying the right hand side of equation (30) by \sqrt{n} , leads

$$[\nabla_{\alpha\alpha}\ell_n(\alpha^*) - \mathcal{S} + \Delta_p]\sqrt{n}(\hat{\alpha} - \alpha^*) = \sqrt{n}[\mathcal{S}\alpha^* - \nabla_{\alpha}\ell_n(\alpha^*)] \quad (31)$$

By assumption **2**, $\mathcal{S} \xrightarrow{p} 0$ and $\mathcal{S}\alpha^* \xrightarrow{p} 0$. Furthermore, by assumption **1**(ix), $\Delta_p \xrightarrow{p} 0$. As earlier mentioned, by Lemma 2(i), Assumption **1**(vii) and the CLT, $n^{1/2}\nabla_{\alpha}\ell_n(\alpha^*) \xrightarrow{d} \mathcal{N}[\mathbf{0}, \mathcal{I}(\alpha^*)]$, and by Lemma 2(ii) and the LLN, $n^{1/2}\nabla_{\alpha\alpha}\ell_n(\alpha^*) \xrightarrow{p} -\mathcal{I}(\alpha^*)$. Finally, by Slutsky's theorem, we obtain

$$\sqrt{n}(\hat{\alpha} - \alpha^*) \xrightarrow{d} \mathcal{N}\{\mathbf{0}, [\mathcal{I}(\alpha^*)]^{-1}\},$$

as required.

- (ii) Under Theorem 1, $\sqrt{n}(\hat{\alpha} - \alpha^*) \xrightarrow{d} \mathcal{N}\{\mathbf{0}, [\mathcal{I}(\alpha^*)]^{-1}\}$. In particular, for $\hat{\alpha}_{\nu 0} \in \hat{\alpha}$ we have $\sqrt{n}(\hat{\alpha}_{\nu 0} - \alpha_{\nu 0}^*) \xrightarrow{d} \mathcal{N}\{\mathbf{0}, [\mathcal{I}(\alpha_{\nu 0}^*)]^{-1}\}$. In addition, $S : \mathbb{R}^k \rightarrow \mathbb{R}$ is continuously differentiable at $\alpha_{\nu 0}^*$, with gradient defined as $\nabla_{\alpha_{\nu 0}}S(\alpha_{\nu 0}^*) = \mathcal{G}'_{\nu 0}[s(\alpha_{\nu 0}^*)]\nabla_{\alpha_{\nu 0}}s(\alpha_{\nu 0}^*)$. Then, we can applied Lemma 4 to obtain

$$\sqrt{n}[\hat{S}_{\nu 0}(\hat{\alpha}_{\nu 0}) - S_{\nu 0}(\alpha_{\nu 0}^*)] \xrightarrow{d} \mathcal{N}\{\mathbf{0}, \mathcal{G}'_{\nu 0}[s(\alpha_{\nu 0}^*)]\nabla_{\alpha_{\nu 0}}s(\alpha_{\nu 0}^*)[\mathcal{I}(\alpha_{\nu 0}^*)]^{-1}\nabla_{\alpha_{\nu 0}}s(\alpha_{\nu 0}^*)^{\top}\mathcal{G}'_{\nu 0}[s(\alpha_{\nu 0}^*)]\}.$$

Furthermore, we know that $\nabla_{\alpha_1\alpha_2}\ell(\alpha) = \mathbf{0}$, therefore $\mathbb{E}[-\nabla_{\alpha_1\alpha_2}\ell(\alpha_0)] = \mathbf{0}$. This also implies that $\mathbb{E}[-\nabla_{\alpha_{10}\alpha_{20}}\ell(\alpha_0)] = \mathbf{0}$, which means that α_{10} and α_{20} are independent. Then, $S(\alpha_{10})$ and $S(\alpha_{20})$ are also independent, as required.

□

Theorem 2 (Asymptotic properties of the NPMLE estimator).

Proof. This proof follows similar arguments of Theorem 1. □

Theorem 3 (Efficiency of the IPMLE estimator).

Proof. For $\nu = 1, 2$, we define $\gamma_{\nu} = (\gamma_{\nu}^t, \gamma_{\nu}^{nu})^{\top}$ so that $\mathcal{Q}_i^{\top}\gamma_{\nu} = \mathcal{Q}_i^{0\top}\gamma_{\nu}^t + \mathcal{Q}_{\nu i}^{1\top}\gamma_{\nu}^{nu}$. Where $\gamma_{\nu}^t = (\gamma_{\nu 1}^{t\top}, \dots, \gamma_{\nu Q}^{t\top})^{\top}$ and $\gamma_{\nu}^{nu} = (\gamma_{\nu(Q+1)}^{nu\top}, \dots, \gamma_{\nu Q_{\nu}}^{nu\top})^{\top}$ are the informative and non-informative

parameters of the non-informative model respectively. Thus, under Assumption 1(viii) and Lemma 2(ii), $\mathcal{I}(\gamma^*)$ can be written as

$$\mathcal{I}(\gamma^*) = \begin{bmatrix} \mathcal{I}_{\gamma_1^t} & \mathcal{I}_{\gamma_1^t \gamma_1^{nu}} & \mathbf{0} & \mathbf{0} \\ \mathcal{I}_{\gamma_1^{nu} \gamma_1^t} & \mathcal{I}_{\gamma_1^{nu}} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathcal{I}_{\gamma_2^t} & \mathcal{I}_{\gamma_2^t \gamma_2^{nu}} \\ \mathbf{0} & \mathbf{0} & \mathcal{I}_{\gamma_2^{nu} \gamma_2^t} & \mathcal{I}_{\gamma_2^{nu}} \end{bmatrix}, \quad (32)$$

where $\mathcal{I}_{\gamma_\nu^t} = \mathcal{I}(\gamma_\nu^{*t})$, $\mathcal{I}_{\gamma_\nu^{nu}} = \mathcal{I}(\gamma_\nu^{*nu})$ and $\mathcal{I}_{\gamma_\nu^t \gamma_\nu^{nu}} = \mathcal{I}(\gamma_\nu^{*t}, \gamma_\nu^{*nu})$. Taking the inverse of (32), we obtain

$$[\mathcal{I}(\gamma^*)]^{-1} = \begin{bmatrix} \Sigma_{\gamma_1^{*t}} & \Sigma_{\gamma_1^{*t} \gamma_1^{*nu}} & \mathbf{0} & \mathbf{0} \\ \Sigma_{\gamma_1^{*nu} \gamma_1^{*t}} & \Sigma_{\gamma_1^{*nu}} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \Sigma_{\gamma_2^{*t}} & \Sigma_{\gamma_2^{*t} \gamma_2^{*nu}} \\ \mathbf{0} & \mathbf{0} & \Sigma_{\gamma_2^{*nu} \gamma_2^{*t}} & \Sigma_{\gamma_2^{*nu}} \end{bmatrix}, \quad (33)$$

where $\Sigma_{\gamma_\nu^{*t}} = [\mathcal{I}_{\gamma_\nu^t} - \mathcal{I}_{\gamma_\nu^t \gamma_\nu^{nu}} \mathcal{I}_{\gamma_\nu^{nu}}^{-1} \mathcal{I}_{\gamma_\nu^{nu} \gamma_\nu^t}]^{-1}$, $\Sigma_{\gamma_\nu^{*t} \gamma_\nu^{*nu}} = -\Sigma_{\gamma_\nu^{*t}} \mathcal{I}_{\gamma_\nu^t \gamma_\nu^{nu}} \mathcal{I}_{\gamma_\nu^{nu}}^{-1}$, $\Sigma_{\gamma_\nu^{*nu} \gamma_\nu^{*t}} = -\mathcal{I}_{\gamma_\nu^{nu}}^{-1} \mathcal{I}_{\gamma_\nu^{nu} \gamma_\nu^t} \Sigma_{\gamma_\nu^{*t}}$ and $\Sigma_{\gamma_\nu^{*nu}} = \mathcal{I}_{\gamma_\nu^{nu}}^{-1} + \mathcal{I}_{\gamma_\nu^{nu}}^{-1} \mathcal{I}_{\gamma_\nu^{nu} \gamma_\nu^t} \Sigma_{\gamma_\nu^{*t}} \mathcal{I}_{\gamma_\nu^t \gamma_\nu^{nu}} \mathcal{I}_{\gamma_\nu^{nu}}^{-1}$.

On the other hand, also by Assumption 1(viii) and Lemma 2(ii), $\mathcal{I}(\alpha^*)$ can be written as

$$\mathcal{I}(\alpha^*) = \begin{bmatrix} \mathcal{I}_{\alpha_0} & \mathcal{I}_{\alpha_0 \alpha_1} & \mathcal{I}_{\alpha_0 \alpha_2} \\ \mathcal{I}_{\alpha_1 \alpha_0} & \mathcal{I}_{\alpha_1} & \mathbf{0} \\ \mathcal{I}_{\alpha_2 \alpha_0} & \mathbf{0} & \mathcal{I}_{\alpha_2} \end{bmatrix}, \quad (34)$$

where $\mathcal{I}_{\alpha_0} = \mathcal{I}(\alpha_0^*)$, $\mathcal{I}_{\alpha_\nu} = \mathcal{I}(\alpha_\nu^*)$, $\mathcal{I}_{\alpha_0 \alpha_\nu} = \mathcal{I}(\alpha_0^*, \alpha_\nu^*)$ and $\mathcal{I}_{\alpha_\nu \alpha_0} = \mathcal{I}(\alpha_\nu^*, \alpha_0^*)$. Taking the inverse of (34), yields

$$[\mathcal{I}(\alpha^*)]^{-1} = \begin{bmatrix} \Sigma_{\alpha_0^*} & \Sigma_{\alpha_0^* \alpha_1^*} & \Sigma_{\alpha_0^* \alpha_2^*} \\ \Sigma_{\alpha_1^* \alpha_0^*} & \Sigma_{\alpha_1^*} & \mathbf{0} \\ \Sigma_{\alpha_2^* \alpha_0^*} & \mathbf{0} & \Sigma_{\alpha_2^*} \end{bmatrix}, \quad (35)$$

where $\Sigma_{\alpha_0^*} = [\mathcal{I}_{\alpha_0} - \mathcal{I}_{\alpha_0 \alpha_1} \mathcal{I}_{\alpha_1}^{-1} \mathcal{I}_{\alpha_1 \alpha_0} - \mathcal{I}_{\alpha_0 \alpha_2} \mathcal{I}_{\alpha_2}^{-1} \mathcal{I}_{\alpha_2 \alpha_0}]^{-1}$, $\Sigma_{\alpha_0^* \alpha_\nu^*} = -\Sigma_{\alpha_0^*} \mathcal{I}_{\alpha_0 \alpha_\nu} \mathcal{I}_{\alpha_\nu}^{-1}$, $\Sigma_{\alpha_\nu^* \alpha_0^*} =$

$$-\mathcal{I}_{\alpha_\nu}^{-1} \mathcal{I}_{\alpha_\nu \alpha_0} \Sigma_{\alpha_0^*} \text{ and } \Sigma_{\alpha_\nu^*} = \mathcal{I}_{\alpha_\nu}^{-1} + \mathcal{I}_{\alpha_\nu}^{-1} \mathcal{I}_{\alpha_\nu \alpha_0} \Sigma_{\alpha_0^*} \mathcal{I}_{\alpha_0 \alpha_\nu} \mathcal{I}_{\alpha_\nu}^{-1}.$$

Thus, by (14), (15), (16), (17), (18) and using that $\gamma_{\nu 0}^{n\iota} = \alpha_{\nu 0}$, we obtain $\mathcal{I}_{\alpha_0} = \mathcal{I}_{\gamma_1^\iota} + \mathcal{I}_{\gamma_2^\iota}$, $\mathcal{I}_{\alpha_0 \alpha_\nu} = \mathcal{I}_{\gamma_\nu^\iota \gamma_\nu^{n\iota}}$, $\mathcal{I}_{\alpha_\nu \alpha_0} = \mathcal{I}_{\gamma_\nu^{n\iota} \gamma_\nu^\iota}$ and $\mathcal{I}_{\alpha_\nu} = \mathcal{I}_{\gamma_\nu^{n\iota}}$. This and the fact that $\Sigma_{\alpha_0^*}^{-1}$ and $\Sigma_{\gamma_\nu^{*\iota}}^{-1}$ are positive definite matrices, imply that $[\Sigma_{\gamma_\nu^{*\iota}} - \Sigma_{\alpha_0^*}]$ is positive definite. Therefore, $\Sigma_{\alpha_0^*} < \Sigma_{\gamma_\nu^{*\iota}}$. Using this reasoning, we conclude that $\Sigma_{\alpha_0^* \alpha_\nu^*} < \Sigma_{\gamma_\nu^{*\iota} \gamma_\nu^{*n\iota}}$, $\Sigma_{\alpha_\nu^* \alpha_0^*} < \Sigma_{\gamma_\nu^{*n\iota} \gamma_\nu^{*\iota}}$ and $\Sigma_{\alpha_\nu^*} < \Sigma_{\gamma_\nu^{*n\iota}}$, as required. \square

The proof of Lemma 3 in the context of informative and non-informative censoring models was adapted from Xingwei et al. (2010) and Vatter & Chavez-Demoulin (2015). The proofs of the asymptotic normality (part (i) of Theorems 1 and 2) are based on Vatter & Chavez-Demoulin (2015).

Supplementary Material E: Confidence Intervals

At convergence, point-wise intervals for linear and non-linear functions for both the non-informative and informative models' parameters can be obtained using the following Bayesian large sample approximation

$$\boldsymbol{\theta} \sim \mathcal{N}(\hat{\boldsymbol{\theta}}, \boldsymbol{\Sigma}_{\hat{\boldsymbol{\theta}}}), \quad (36)$$

where $\boldsymbol{\Sigma}_{\hat{\boldsymbol{\theta}}} = [\mathcal{H}_p(\hat{\boldsymbol{\theta}})]^{-1}$. For generalised additive models, intervals derived using equation (36) have good frequentist properties, since they account for both smoothing bias and sampling variability (Marra & Wood, 2012). For the non-informative and informative models, equation (36) can be verified using the distribution of \mathcal{Z} (described in Supplementary Material C), making the large sample assumption that $\mathcal{H}(\boldsymbol{\theta})$ can be treated as fixed, and making the usual prior Bayesian assumption for smooth models $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{S}^{-1})$, where \boldsymbol{S}^{-1} is the Moore-Penrose pseudoinverse of \boldsymbol{S} (Silverman, 1985; Wood, 2017). In equation (36), smoothing parameter uncertainty is neglected. Nevertheless, according to Marra & Wood (2012) this is not problematic if heavy over-smoothing is avoided so that the smoothing bias is not a large proportion of the sampling variability. See also Marra et al. (2017) for an application of this approach to a more general smoothing spline context.

Following Pya & Wood (2015), confidence interval estimates for the monotonic smooth terms in the models can be obtained using the distribution of $\tilde{\boldsymbol{\beta}}_{\nu 0}$ (defined in Section 2.3 of the main paper) since all smooth components would then depend linearly on $\tilde{\boldsymbol{\beta}}_{\nu 0}$. Such distribution is

$$\tilde{\boldsymbol{\beta}}_{\nu 0} \sim \mathcal{N}(\hat{\tilde{\boldsymbol{\beta}}}_{\nu 0}, \boldsymbol{\Sigma}_{\tilde{\boldsymbol{\beta}}_{\nu 0}}),$$

where $\boldsymbol{\Sigma}_{\tilde{\boldsymbol{\beta}}_{\nu 0}} = \text{diag}(\boldsymbol{\Gamma}_{\nu 0}) [\mathcal{H}_p(\hat{\tilde{\boldsymbol{\beta}}}_{\nu 0})]^{-1} \text{diag}(\boldsymbol{\Gamma}_{\nu 0})$. The derivation of this result can be found in Pya & Wood (2015).

P-values for the smooth components in the non-informative and informative models are obtained by adapting the results discussed in Wood (2013) to the present context, where $\boldsymbol{\Sigma}_{\tilde{\boldsymbol{\beta}}_{\nu 0}}$ is used for the calculations. The reader is referred to the above citation for the definition of reference degrees of freedom.

Supplementary Material F: Model Selection

In practical situations, it is important to detect if $\sum_{k_1=1}^{K_1} s_{1k_1}(\mathbf{x}_{1k_1i})$ and $\sum_{k_2=1}^{K_2} s_{2k_2}(\mathbf{x}_{2k_2i})$ have components in common. This is basically a model selection problem and, to this end, we propose using the AIC, BIC and K-Fold Cross validation criterion (Υ^{KCV}). The AIC and BIC can be defined as

$$\text{AIC} = -2\ell(\hat{\boldsymbol{\theta}}) + 2 \text{EDF},$$

$$\text{BIC} = -2\ell(\hat{\boldsymbol{\theta}}) + \log(n) \text{EDF},$$

where the log-likelihood is evaluated at the penalized parameter estimates and $\text{EDF} = \text{tr}(\hat{\mathbf{B}})$ with $\hat{\mathbf{B}}$ defined in Supplementary Material C.

As for Υ^{KCV} (Stone, 1974), we first randomly divide the set of observations in K groups (folds) of approximately equal size. Each fold is then in turn treated as a validation set, and the IPMLE for a given model is used to estimate the vector of parameters $\boldsymbol{\alpha}$ using the remaining $K - 1$ folds. The so obtained estimates are denoted as $\hat{\boldsymbol{\alpha}}_0^{\setminus k}$ and $\hat{\boldsymbol{\alpha}}_1^{\setminus k}$, and the log-likelihood function is calculated as

$$\begin{aligned} \ell_k(\hat{\boldsymbol{\alpha}}^{\setminus k}) = & \left\{ \log \mathcal{G}_1 \left[\xi_{1i}(\hat{\boldsymbol{\alpha}}_0^{\setminus k}, \hat{\boldsymbol{\alpha}}_1^{\setminus k}) \right] + \delta_{1i} \log \left\{ -\frac{\mathcal{G}'_1 \left[\xi_{1i}(\hat{\boldsymbol{\alpha}}_0^{\setminus k}, \hat{\boldsymbol{\alpha}}_1^{\setminus k}) \right]}{\mathcal{G}_1 \left[\xi_{1i}(\hat{\boldsymbol{\alpha}}_0^{\setminus k}, \hat{\boldsymbol{\alpha}}_1^{\setminus k}) \right]} \frac{\partial \xi_{1i}(\hat{\boldsymbol{\alpha}}_0^{\setminus k}, \hat{\boldsymbol{\alpha}}_1^{\setminus k})}{\partial y_i} \right\} \right\} \\ & + \left\{ \log \mathcal{G}_2 \left[\xi_{2i}(\hat{\boldsymbol{\alpha}}_0^{\setminus k}, \hat{\boldsymbol{\alpha}}_2^{\setminus k}) \right] + \delta_{2i} \log \left\{ -\frac{\mathcal{G}'_2 \left[\xi_{2i}(\hat{\boldsymbol{\alpha}}_0^{\setminus k}, \hat{\boldsymbol{\alpha}}_2^{\setminus k}) \right]}{\mathcal{G}_2 \left[\xi_{2i}(\hat{\boldsymbol{\alpha}}_0^{\setminus k}, \hat{\boldsymbol{\alpha}}_2^{\setminus k}) \right]} \frac{\partial \xi_{2i}(\hat{\boldsymbol{\alpha}}_0^{\setminus k}, \hat{\boldsymbol{\alpha}}_2^{\setminus k})}{\partial y_i} \right\} \right\}, \end{aligned}$$

and Υ^{KCV} given by

$$\Upsilon^{\text{KCV}} = \sum_{k=1}^K \ell_k(\hat{\boldsymbol{\alpha}}^{\setminus k}). \quad (37)$$

We choose the model which maximizes (37). The same procedure is used when Υ^{KCV} is calculated for the non-informative model. In such a case we have

$$\begin{aligned} \ell_k(\hat{\boldsymbol{\gamma}}^{\setminus k}) = & \left\{ \log \mathcal{G}_1 \left[\xi_{1i}(\hat{\boldsymbol{\gamma}}_1^{\setminus k}) \right] + \delta_{1i} \log \left\{ -\frac{\mathcal{G}'_1 \left[\xi_{1i}(\hat{\boldsymbol{\gamma}}_1^{\setminus k}) \right]}{\mathcal{G}_1 \left[\xi_{1i}(\hat{\boldsymbol{\gamma}}_1^{\setminus k}) \right]} \frac{\partial \xi_{1i}(\hat{\boldsymbol{\gamma}}_1^{\setminus k})}{\partial y_i} \right\} \right\} \\ & + \left\{ \log \mathcal{G}_2 \left[\xi_{2i}(\hat{\boldsymbol{\gamma}}_2^{\setminus k}) \right] + \delta_{2i} \log \left\{ -\frac{\mathcal{G}'_2 \left[\xi_{2i}(\hat{\boldsymbol{\gamma}}_2^{\setminus k}) \right]}{\mathcal{G}_2 \left[\xi_{2i}(\hat{\boldsymbol{\gamma}}_2^{\setminus k}) \right]} \frac{\partial \xi_{2i}(\hat{\boldsymbol{\gamma}}_2^{\setminus k})}{\partial y_i} \right\} \right\}, \end{aligned}$$

and therefore $\Upsilon^{\text{KCV}} = \sum_{k=1}^K \ell_k(\hat{\gamma}^k)$.

Model	Non-Inf. Covariates	Inf. Covariates	Link T _{1i}	Link T _{2i}	AIC	Υ^{KCV}	BIC
1	s(wmonth) s(mthage) region alcohol nsibs	...	PH	PH	13775.68	-6924.20	14015.53
2	s(wmonth) s(mthage) region alcohol nsibs	...	PO	PH	13776.87	-8396.57	14016.51
3	s(wmonth) s(mthage) nsibs	alcohol region	PH	PH	13772.60	-6922.63	13981.42
4	s(wmonth) s(mthage) nsibs	alcohol region	PO	PH	13773.80	-8392.31	13982.51

Table 2: Values of three model selection criteria (AIC, BIC and Υ^{KCV}) for the best informative and non-informative models fitted to the real data application of this paper. The models were fitted using `gam1ss()` in GJRM by employing different combinations of covariates and link functions.

Supplementary Material G: Additional simulation results for DGP1 and DGP2 and findings from a simulation study with mild censoring rate

In the DGP presented in this section (DGP3), z_{1i} is informative, z_{2i} is informative and a mild censoring rate (about 47%) is considered. T_{1i} and T_{2i} were generated using the model defined in equation (19) of the main paper. The baseline survival functions were defined as $S_{10}(t_{1i}) = 0.8 \exp(-0.4t_{1i}^{2.5}) + 0.2 \exp(-0.1t_{1i}^{1.0})$ and $S_{20}(t_{2i}) = 0.99 \exp(-0.05t_{2i}^{2.3}) + 0.01 \exp(-0.4t_{2i}^{1.1})$. The informative covariates, z_{1i} and z_{2i} , were generated using a binomial and a uniform distribution respectively. Also, $s_{11}(z_{2i}) = s_{12}(z_{2i}) = \sin(2\pi z_{2i})$, $\alpha_{01} = -0.10$, $\alpha_{02} = -0.25$ and $\alpha_{11} = \alpha_{12} = -1.5$.

The main findings are:

- Figure 1 and Table 4 show that overall the mean estimates for the two estimators are very close to the respective true values and improve as the sample size increases. However, even though the variability of the estimates (IPMLE and NPMLE) decreases as the sample size grows large, the IPMLE is slightly more efficient than the NPMLE in recovering the true linear effects for all sample sizes examined here. In particular, the RMSE of the IPMLE is slightly smaller than the RMSE of the NPMLE for all sample sizes considered.
- Figures 2 and 3, and Table 4 show that overall the true functions are recovered well by the IPMLE and NPMLE and that the results improve in terms of bias and efficiency as the sample size increases. Furthermore, the IPMLE is slightly more efficient than the NPMLE in recovering the non-linear covariate effects for all sample sizes examined in this section (Table 4). However, this gain in efficiency by the IPMLE is not too significant when a mild censoring rate (47%) is examined.

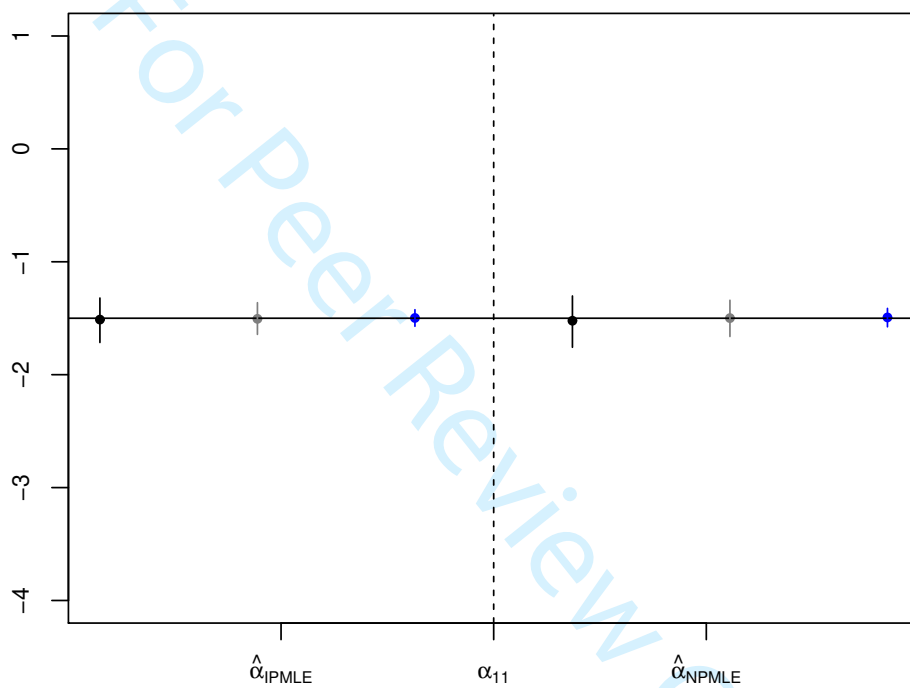


Figure 1: Linear coefficient estimates obtained by applying `gam1ss()` to informative survival data simulated according to DGP3 characterised by a censoring rate of about 47%. Circles indicate mean estimates while bars represent the estimates' ranges resulting from 5% and 95% quantiles. True values are indicated by black solid horizontal lines. Black circles and vertical bars refer to the results obtained for $n = 500$, whereas those for $n = 1000$ and $n = 4000$ are given in dark gray and blue, respectively.

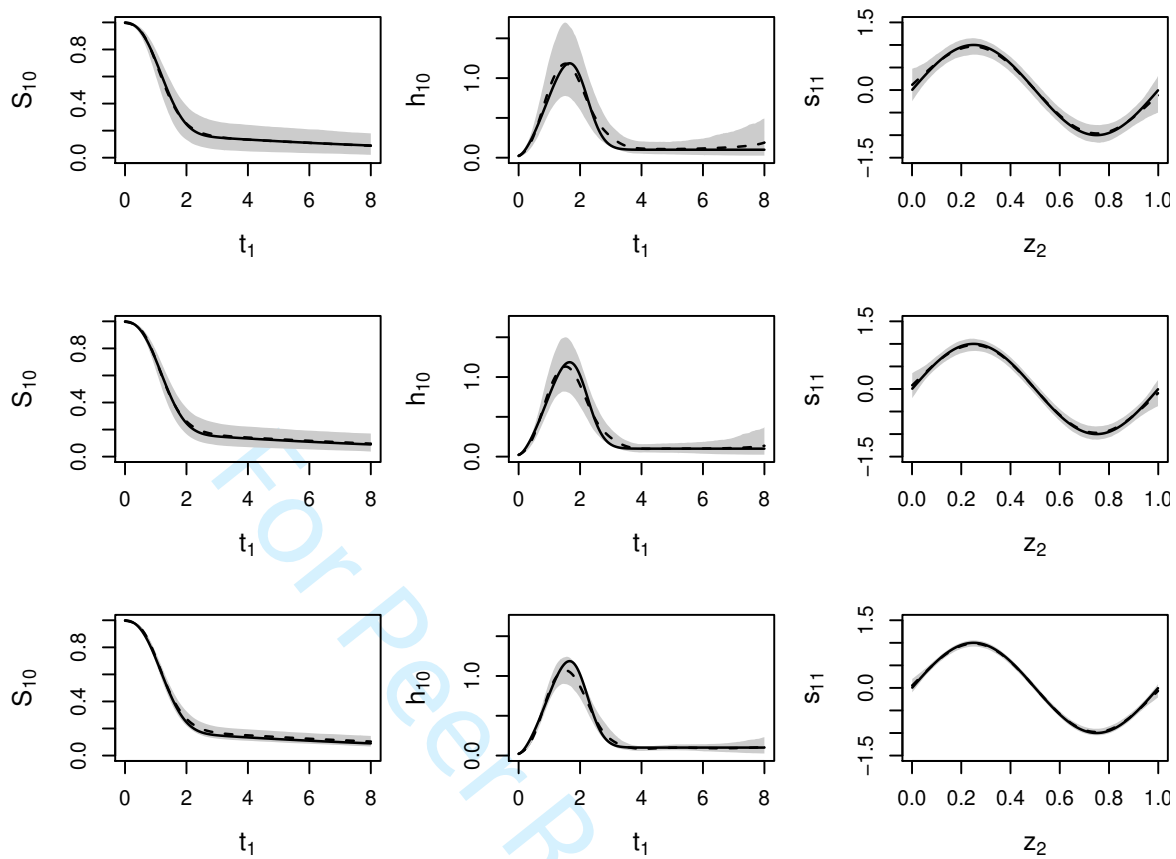


Figure 2: Smooth function estimates for the IPMLE obtained by applying `gamlss()` to informative survival data simulated according to DGP3 characterised by a censoring rate of about 47%. True functions are represented by black solid lines, mean estimates by dashed lines and pointwise ranges resulting from 5% and 95% quantiles by shaded areas. The results in the first row refer to $n = 500$, whereas those in the second and third rows to $n = 1000$ and $n = 4000$.

(a) Informative Penalized Maximum Log-likelihood Estimator (IPMLE)						
	Bias			RMSE		
	$n = 500$	$n = 1000$	$n = 4000$	$n = 500$	$n = 1000$	$n = 4000$
α_{11}	-0.024	-0.014	-0.006	0.138	0.100	0.049
s_1	0.039	0.025	0.012	0.154	0.114	0.059
h_{10}	0.084	0.048	0.035	0.262	0.144	0.083
S_{10}	0.028	0.020	0.017	0.063	0.050	0.031
(b) Non-informative Penalized Maximum Log-likelihood Estimator (NPMLE)						
	Bias			RMSE		
	$n = 500$	$n = 1000$	$n = 4000$	$n = 500$	$n = 1000$	$n = 4000$
α_{11}	-0.045	-0.017	-0.007	0.208	0.144	0.071
s_1	0.085	0.068	0.044	0.191	0.206	0.111
h_{10}	0.085	0.057	0.033	0.195	0.292	0.083
S_{10}	0.027	0.021	0.015	0.058	0.068	0.033

Table 3: Bias and root mean squared error (RMSE) for the IPMLE and NPMLE obtained by applying the `gamlss()` to informative survival data simulated according to DGP2 characterised by a censoring rate of about 74%. Further details are given in the caption of Table 1.

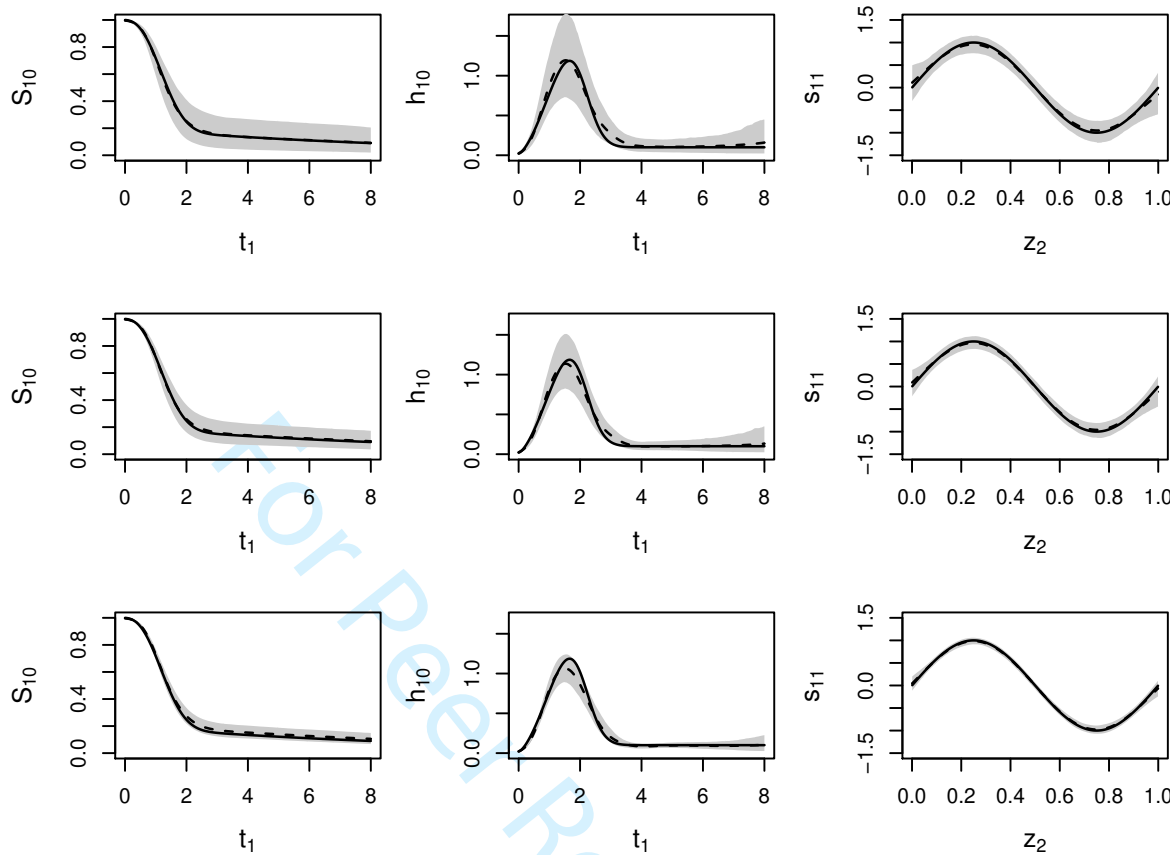


Figure 3: Smooth function estimates for the NPMLE obtained by applying `gamlss()` to informative survival data simulated according to DGP3 characterised by a censoring rate of about 47%. Further details are given in the caption of Figure 2.

(a) Informative Penalized Maximum Log-likelihood Estimator (IPMLE)						
	Bias			RMSE		
	n = 500	n = 1000	n = 4000	n = 500	n = 1000	n = 4000
α_{11}	-0.012	-0.006	0.003	0.121	0.058	0.045
s_1	0.031	0.021	0.015	0.124	0.091	0.051
h_{10}	0.040	0.027	0.026	0.135	0.088	0.058
S_{10}	0.003	0.008	0.015	0.057	0.047	0.030
(b) Non-informative Penalized Maximum Log-likelihood Estimator (NPMLE)						
	Bias			RMSE		
	n = 500	n = 1000	n = 4000	n = 500	n = 1000	n = 4000
α_{11}	-0.022	0.001	0.007	0.140	0.100	0.050
s_1	0.036	0.027	0.014	0.142	0.104	0.055
h_{10}	0.037	0.027	0.027	0.131	0.089	0.056
S_{10}	0.004	0.008	0.017	0.065	0.047	0.032

Table 4: Bias and root mean squared error (RMSE) for the IPMLE and NPMLE obtained by applying `gamlss()` to informative survival data simulated according to DGP3 characterised by a censoring rate of about 47%. Further details are given in the caption of Table 1.

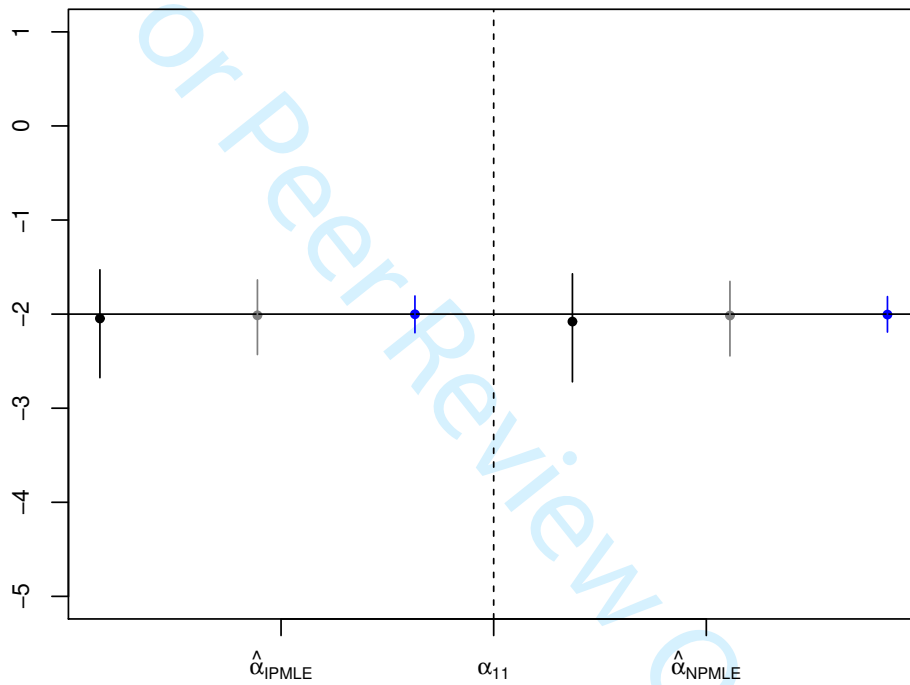


Figure 4: Linear coefficient estimates obtained by applying `gamlss()` to informative survival data simulated according to DGP1 which is characterised by a censoring rate of about 78%. Further details are given in the caption of Figure 1.

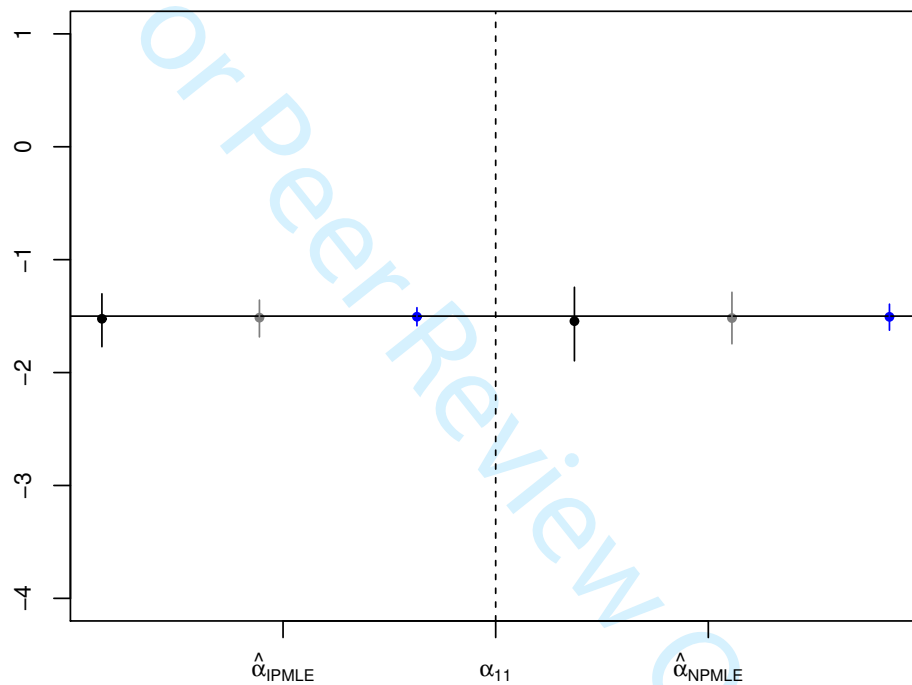


Figure 5: Linear coefficient estimates obtained by applying `gamlss()` to informative survival data simulated according to DGP2 which is characterised by a censoring rate of about 74%. Further details are given in the caption of Figure 1.

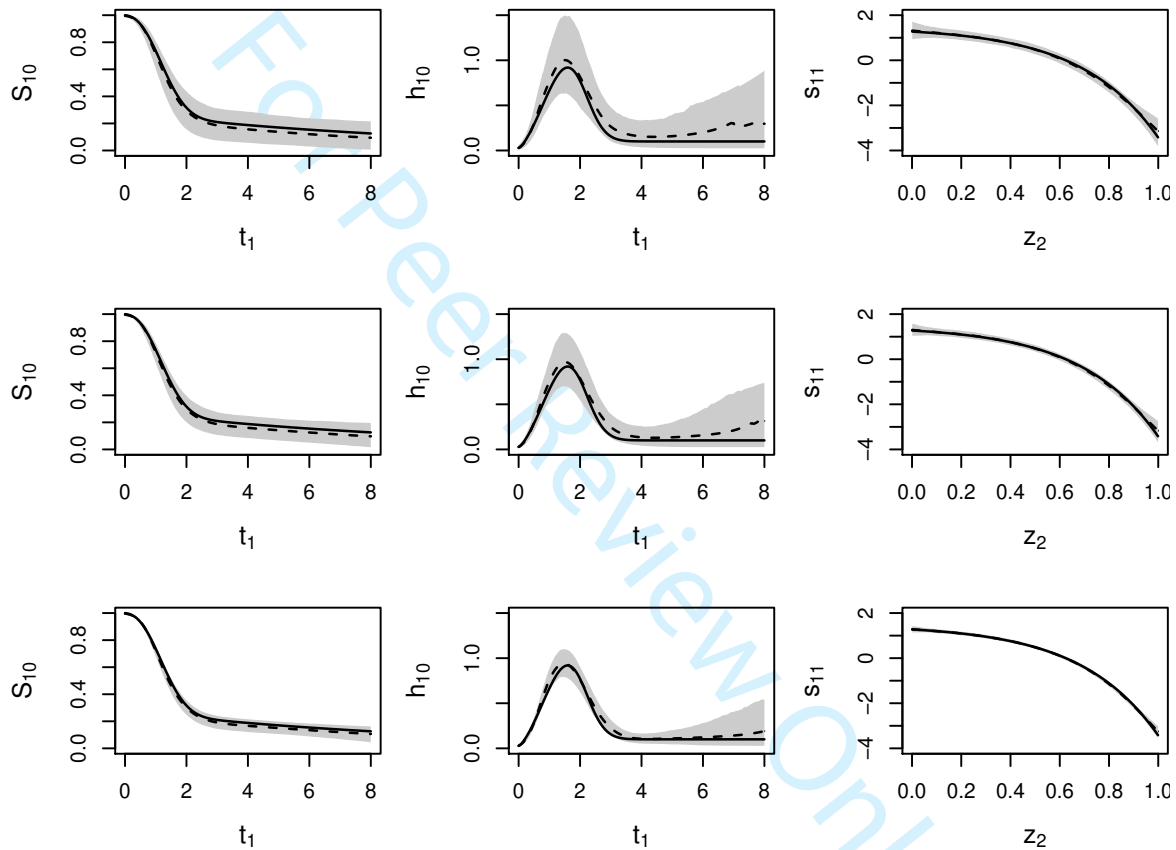


Figure 6: Smooth function estimates for the IPMLE obtained by applying `gamlss()` to informative survival data simulated according to DGP1 characterised by a censoring rate of about 78%. Further details are given in the caption of Figure 2.

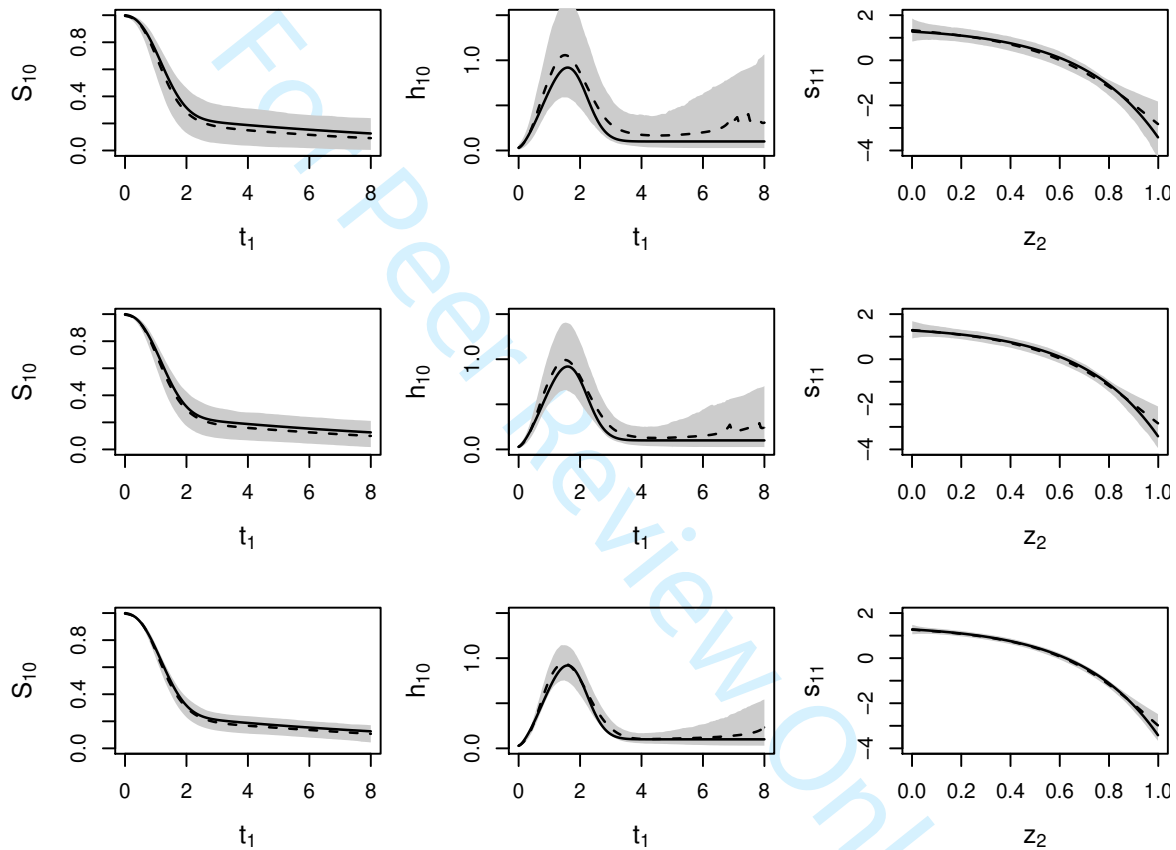


Figure 7: Smooth function estimates for the NPMLE obtained by applying `gamlss()` to informative survival data simulated according to DGP1 characterised by a censoring rate of about 78%. Further details are given in the caption of Figure 2.

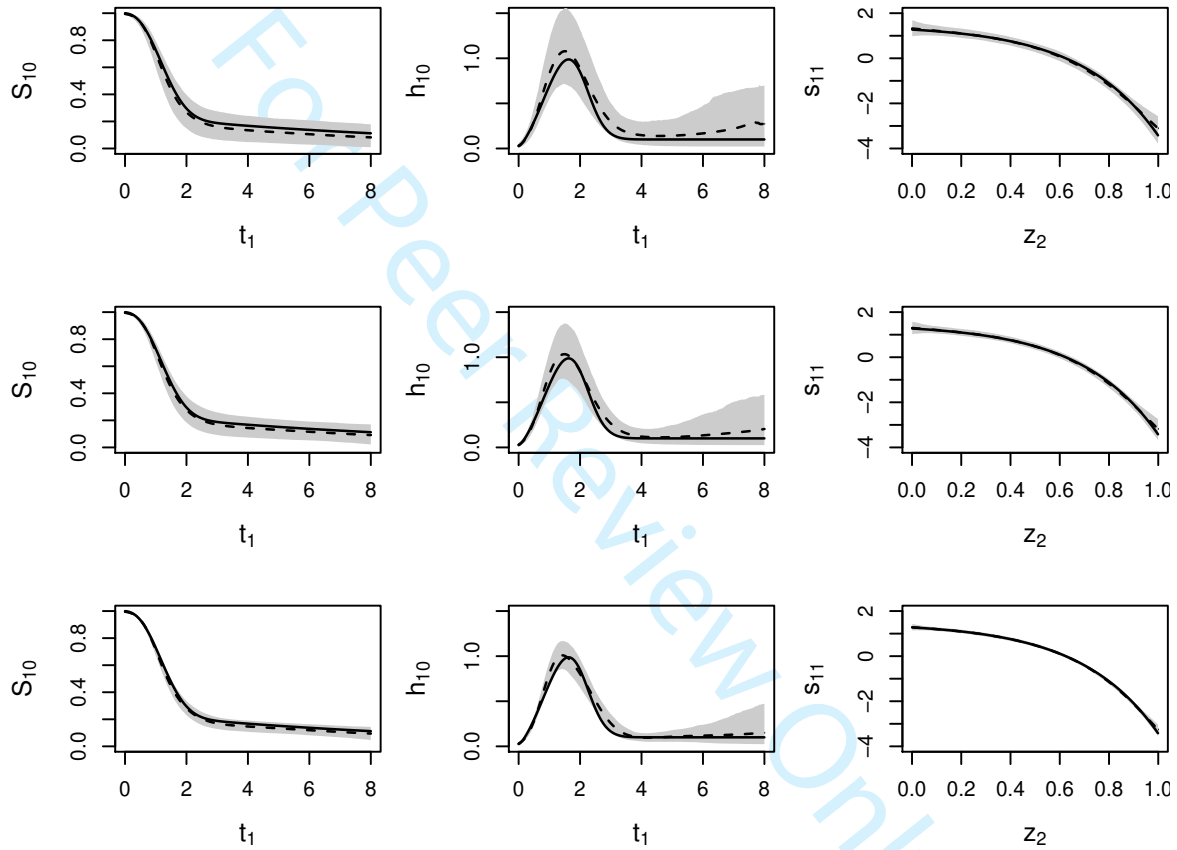


Figure 8: Smooth function estimates for the IPMLE obtained by applying `gamlss()` to informative survival data simulated according to DGP2 characterised by a censoring rate of about 74%. Further details are given in the caption of Figure 2

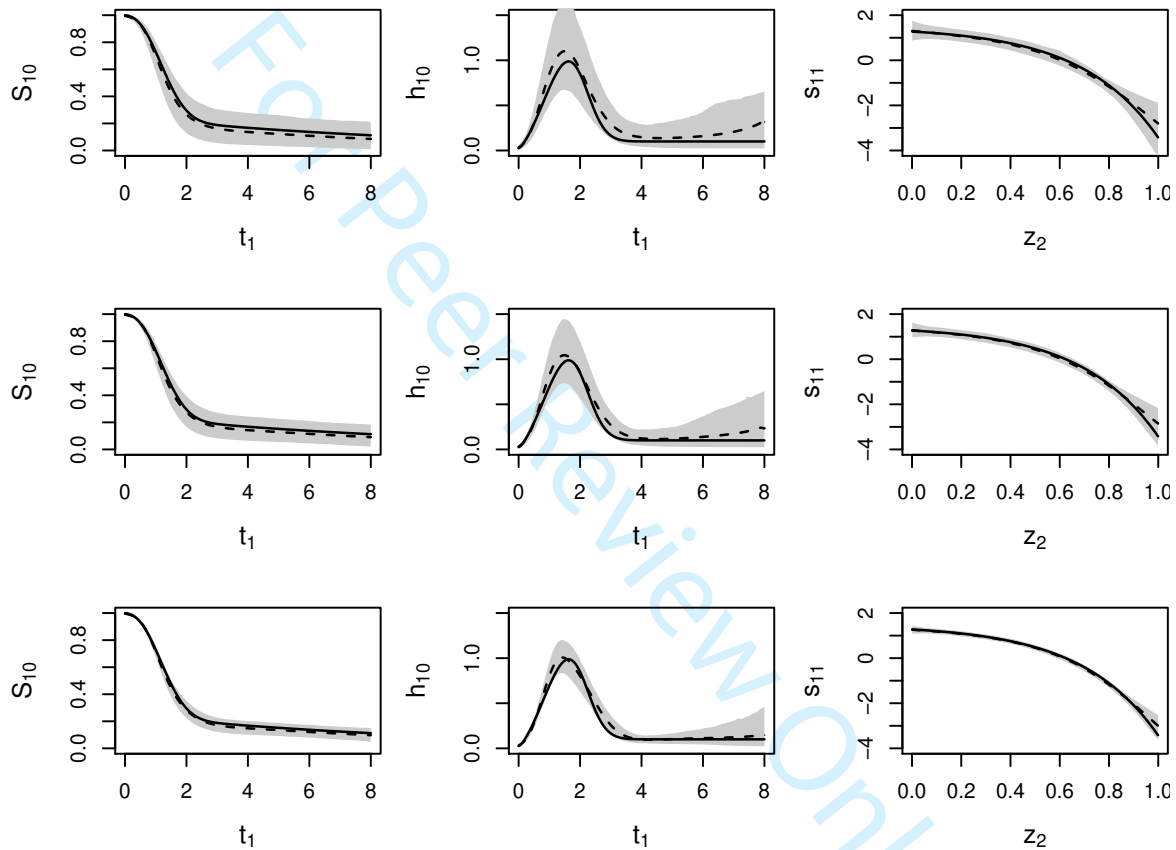


Figure 9: Smooth function estimates for the NPMLE obtained by applying `gamlss()` to informative survival data simulated according to DGP2 characterised by a censoring rate of about 74%. Further details are given in the caption of Figure 2

References

- 1
2
3 Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle.
4
5 *In: Petrov, B.N., Csaki, B.F. (eds.) Second International Symposium on Information Theory.*
6
7 *Academiai Kiado, Budapest.*
8
9
- 10 Gourieroux, C. & Monfort, A. (1995). *Statistics and econometric models*, volume 1. Cambridge
11
12 University Press.
13
- 14 Gu, C. (1992). Cross-validating non-gaussian data. *Journal of Computational and Graphical*
15
16 *Statistics*, 1(2), 169–179.
17
18
- 19 Hayashi, F. (2000). *Econometrics*. Princeton University Press Princeton, NJ.
20
21
- 22 Marra, G. & Radice, R. (2019). *GJRM: Generalised Joint Regression Modelling*. R package
23
24 version 0.2-1.
25
26
- 27 Marra, G., Radice, R., Bärnighausen, T., Wood, S. N., & McGovern, M. E. (2017). A simultaneous
28
29 equation approach to estimating hiv prevalence with nonignorable missing responses. *Journal*
30
31 *of the American Statistical Association*, 112(518), 484–496.
32
33
- 34 Marra, G. & Wood, S. N. (2012). Coverage properties of confidence intervals for generalized
35
36 additive model components. *Scandinavian Journal of Statistics*, 39(1), 53–74.
37
38
- 39 Newey, W. K. & McFadden, D. (1994). *Handbook of econometrics*, volume 4. Elsevier.
40
41
- 42 Nocedal, J. & Wright, S. (2006). *Numerical optimization, series in operations research and finan-*
43
44 *cial engineering*. Springer, New York, USA, 2006.
45
46
- 47 Pya, N. & Wood, S. N. (2015). Shape constrained additive models. *Statistics and Computing*,
48
49 25(3), 543–559.
50
- 51 Silverman, B. W. (1985). Some aspects of the spline smoothing approach to non-parametric re-
52
53 gression curve fitting. *Journal of the Royal Statistical Society: Series B (Methodological)*,
54
55 47(1), 1–21.
56
57
58
59
60

- 1 Stone, M. (1974). Cross-validatory choice and assessment of statistical predictions. *Journal of the*
2 *royal statistical society. Series B (Methodological)*, 36(2), 111–147.
3
4
5 Vatter, T. & Chavez-Demoulin, V. (2015). Generalized additive models for conditional dependence
6 structures. *Journal of Multivariate Analysis*, 141, 147–167.
7
8
9
10 Wood, S. N. (2004). Stable and efficient multiple smoothing parameter estimation for generalized
11 additive models. *Journal of the American Statistical Association*, 99(467), 673–686.
12
13
14
15 Wood, S. N. (2013). On p-values for smooth components of an extended generalized additive
16 model. *Biometrika*, 100(1), 221–228.
17
18
19
20 Wood, S. N. (2017). *Generalized Additive Models: An Introduction With R*. Second Edition,
21 Chapman & Hall/CRC, London.
22
23
24 Xingwei, T., Tao, H., & Hengjian, C. (2010). Hazard regression with penalized spline: The
25 smoothing parameter choice and asymptotics. *Acta Mathematica Scientia*, 30(5), 1759–1768.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60