

# 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

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

where `fl` 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.  $\Phi$  and  $\phi$  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  $\gamma_1$  and  $\gamma_2$  would have some components in common. Because the first  $Q$  components of  $\gamma_1$  are the same as the first  $Q$  components of  $\gamma_2$ , 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  $\frac{\partial\ell(\alpha)}{\partial\alpha_0} = \left[\frac{\partial\ell(\alpha)}{\partial\alpha_{011}} \dots \frac{\partial\ell(\alpha)}{\partial\alpha_{0QJ_Q}}\right]^\top$ ,  $\frac{\partial\ell(\alpha)}{\partial\alpha_1} = \left[\frac{\partial\ell(\alpha)}{\partial\alpha_{111}} \dots \frac{\partial\ell(\alpha)}{\partial\alpha_{1Q_1J_1Q_1}}\right]^\top$  and  $\frac{\partial\ell(\alpha)}{\partial\alpha_2} = \left[\frac{\partial\ell(\alpha)}{\partial\alpha_{21}} \dots \frac{\partial\ell(\alpha)}{\partial\alpha_{2Q_2J_2Q_2}}\right]^\top$ . 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

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

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

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

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

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

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

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

Since  $\alpha_1$  appears only in  $\xi_{1i}(\alpha_0, \alpha_1)$  and  $\alpha_2$  only in  $\xi_{2i}(\alpha_0, \alpha_2)$ , then  $\nabla_{\alpha_1 \alpha_2} \ell(\alpha) = \nabla_{\alpha_2 \alpha_1} \ell(\alpha) = 0$ . 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} \quad (10)$$

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

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



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

$$\begin{aligned} \frac{\partial^2 \ell(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}_0 \partial \boldsymbol{\alpha}_0^\top} &= \sum_{i=1}^n \left\{ \Phi_1 \frac{\partial \xi_{1i}}{\partial \boldsymbol{\alpha}_0} \left[ \frac{\partial \xi_{1i}}{\partial \boldsymbol{\alpha}_0} \right]^\top + \Phi_2 \frac{\partial \xi_{2i}}{\partial \boldsymbol{\alpha}_0} \left[ \frac{\partial \xi_{2i}}{\partial \boldsymbol{\alpha}_0} \right]^\top \right\}, \\ \frac{\partial^2 \ell(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}_0 \partial \boldsymbol{\alpha}_\nu^\top} &= \sum_{i=1}^n \left\{ \Phi_\nu \frac{\partial \xi_{\nu i}}{\partial \boldsymbol{\alpha}_0} \left[ \frac{\partial \xi_{\nu i}}{\partial \boldsymbol{\alpha}_\nu} \right]^\top \right\}, \\ \frac{\partial^2 \ell(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}_\nu \partial \boldsymbol{\alpha}_\nu^\top} &= \sum_{i=1}^n \left\{ \Phi_\nu \frac{\partial \xi_{\nu i}}{\partial \boldsymbol{\alpha}_\nu} \left[ \frac{\partial \xi_{\nu i}}{\partial \boldsymbol{\alpha}_\nu} \right]^\top + \Delta_\nu \frac{\partial^2 \xi_{\nu i}}{\partial \boldsymbol{\alpha}_\nu \partial \boldsymbol{\alpha}_\nu^\top} - \Psi_\nu \frac{\partial^2 \xi_{\nu i}}{\partial y_i \partial \boldsymbol{\alpha}_\nu} \left[ \frac{\partial^2 \xi_{\nu i}}{\partial y_i \partial \boldsymbol{\alpha}_\nu} \right]^\top + \Omega_\nu \frac{\partial^3 \xi_{\nu i}}{\partial y_i \partial \boldsymbol{\alpha}_\nu \partial \boldsymbol{\alpha}_\nu^\top} \right\}, \end{aligned}$$

where

$$\begin{aligned} \frac{\partial^2 \xi_{\nu i}}{\partial \boldsymbol{\alpha}_\nu \partial \boldsymbol{\alpha}_\nu^\top} &= \begin{bmatrix} \frac{\partial^2 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial \alpha_{\nu 11} \partial \alpha_{\nu 11}} & \cdots & \frac{\partial^2 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial \alpha_{\nu 11} \partial \alpha_{\nu Q_\nu J_\nu Q_\nu}} \\ \cdots & \ddots & \cdots \\ \frac{\partial^2 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial \alpha_{\nu Q_\nu J_\nu Q_\nu} \partial \alpha_{\nu 11}} & \cdots & \frac{\partial^2 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial \alpha_{\nu Q_\nu J_\nu Q_\nu} \partial \alpha_{\nu Q_\nu J_\nu Q_\nu}} \end{bmatrix}, \\ \frac{\partial^3 \xi_{\nu i}}{\partial y_i \partial \boldsymbol{\alpha}_\nu \partial \boldsymbol{\alpha}_\nu^\top} &= \begin{bmatrix} \frac{\partial^3 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial y_i \partial \alpha_{\nu 11} \partial \alpha_{\nu 11}} & \cdots & \frac{\partial^3 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial y_i \partial \alpha_{\nu 11} \partial \alpha_{\nu Q_\nu J_\nu Q_\nu}} \\ \cdots & \ddots & \cdots \\ \frac{\partial^3 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial y_i \partial \alpha_{\nu Q_\nu J_\nu Q_\nu} \partial \alpha_{\nu 11}} & \cdots & \frac{\partial^3 \xi_{\nu i}(\boldsymbol{\alpha}_\nu)}{\partial y_i \partial \alpha_{\nu Q_\nu J_\nu Q_\nu} \partial \alpha_{\nu Q_\nu J_\nu Q_\nu}} \end{bmatrix}. \end{aligned}$$

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

$$\begin{aligned} \frac{\partial^2 \xi_{\nu i}(\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_\nu)}{\partial \boldsymbol{\alpha}_{\nu q_\nu} \partial \boldsymbol{\alpha}_{\nu s_\nu}^\top} &= \begin{cases} \mathcal{Q}_{\nu 0}^{\iota \Delta \Delta}(y_i) & \text{if } \boldsymbol{\alpha}_{\nu q_\nu} = \boldsymbol{\alpha}_{\nu s_\nu} = \boldsymbol{\alpha}_{\nu 0} \\ 0 & \text{otherwise,} \end{cases} \\ \frac{\partial^3 \xi_{\nu i}(\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_\nu)}{\partial y_i \partial \boldsymbol{\alpha}_{\nu q_\nu} \boldsymbol{\alpha}_{\nu s_\nu}^\top} &= \begin{cases} \mathcal{Q}_{\nu 0}^{\iota \Delta \Delta'}(y_i) & \text{if } \boldsymbol{\alpha}_{\nu q_\nu} = \boldsymbol{\alpha}_{\nu s_\nu} = \boldsymbol{\alpha}_{\nu 0} \\ 0 & \text{otherwise,} \end{cases} \end{aligned}$$

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

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

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

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'''}{\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. \\
&\quad - \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}} \\
&\quad \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'''}{\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 + \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] \\
&\quad \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\},
\end{aligned} \tag{15}$$

$$\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'''}{\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. \\
&\quad - \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}} \\
&\quad \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'''}{\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. \\
&\quad + \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] \\
&\quad \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\}.
\end{aligned} \tag{16}$$

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

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

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

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  $\alpha$  and maintaining  $\lambda$  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

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

where  $\bar{\ell}_p(\alpha^{[a]}) = -\{\ell_p(\alpha^{[a]}) + \epsilon^\top \mathbf{g}_p(\alpha^{[a]}) + \frac{1}{2}\epsilon^\top \mathbf{H}_p(\alpha^{[a]})\epsilon\}$ ,  $\mathbf{g}_p(\alpha^{[a]}) = \mathbf{g}(\alpha^{[a]}) - \mathcal{S}\alpha^{[a]}$ ,  $\mathbf{H}_p(\alpha^{[a]}) = \mathbf{H}(\alpha^{[a]}) - \mathcal{S}$ . 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  $\mathbf{H}(\alpha^{[a]})_{l,j} = \nabla_{\alpha_l \alpha_j} \ell(\alpha)|_{\alpha_l=\alpha_l^{[a]}, \alpha_j=\alpha_j^{[a]}}$ , 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  $\lambda$  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  $\mathbf{H}(\alpha)$ . The main ideas and some useful results are given here.

To simplify the notation,  $\mathbf{g}_p(\alpha^{[a]})$ ,  $\mathbf{g}(\alpha^{[a]})$ ,  $\mathbf{H}_p(\alpha^{[a]})$  and  $\mathbf{H}(\alpha^{[a]})$  are denoted as  $\mathbf{g}_p^{[a]}$ ,  $\mathbf{g}^{[a]}$ ,  $\mathbf{H}_p^{[a]}$  and  $\mathbf{H}^{[a]}$ . First, it is necessary to express the parameter estimator in terms of  $\mathbf{g}_p^{[a]}$  and  $\mathbf{H}_p^{[a]}$ . To achieve this, a first order Taylor expansion of  $\mathbf{g}_p^{[a+1]}$  about  $\alpha^{[a]}$  is used, which yields the following expression:  $\mathbf{0} = \mathbf{g}_p^{[a+1]} \approx \mathbf{g}_p^{[a]}(\alpha^{[a+1]} - \alpha^{[a]})\mathbf{H}_p^{[a]}$ . After some manipulations,  $\alpha^{[a+1]} = (-\mathbf{H}^{[a]} + \mathcal{S})^{-1} \sqrt{-\mathbf{H}^{[a]}} [\sqrt{-\mathbf{H}^{[a]}} \alpha^{[a]} + \sqrt{-\mathbf{H}^{[a]}}^{-1} \mathbf{g}^{[a]}]$  is obtained, which then becomes  $\alpha^{[a+1]} = (-\mathbf{H}^{[a]} + \mathcal{S})^{-1} \sqrt{-\mathbf{H}^{[a]}} \mathbf{Z}^{[a]}$ , where  $\mathbf{Z}^{[a]} = \mathbf{v}_Z^{[a]} + \boldsymbol{\xi}_Z^{[a]}$ ,  $\mathbf{v}_Z^{[a]} = \sqrt{-\mathbf{H}^{[a]}} \alpha^{[a]}$  and  $\boldsymbol{\xi}_Z^{[a]} = \sqrt{-\mathbf{H}^{[a]}}^{-1} \mathbf{g}^{[a]}$ . Eigenvalue decomposition is used to obtain the square root of  $-\mathbf{H}^{[a]}$  and its inverse. Furthermore, from likelihood theory,  $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $\mathbf{Z} \sim \mathcal{N}(\mathbf{v}_Z, \mathbf{I})$ , where  $\mathbf{v}_Z = \sqrt{-\mathbf{H}} \alpha^0$ ,  $\alpha^0$  is the true parameter vector and  $\mathbf{I}$  is the identity matrix.  $\hat{\mathbf{v}}_Z = \sqrt{-\mathbf{H}} \hat{\alpha} = \mathbf{B} \mathbf{Z}$  is the predicted value vector for  $\mathbf{Z}$ , where  $\mathbf{B} = \sqrt{-\mathbf{H}}(-\mathbf{H} + \mathcal{S})^{-1} \sqrt{-\mathbf{H}}$ . Since our objective is to estimate  $\lambda$  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,  $\lambda$  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  $\lambda$  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  $\lambda$  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

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

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  $\log(\lambda)$ , because the smoothing parameters can only take positive values.

The methods for estimating  $\boldsymbol{\alpha}$  and  $\lambda$  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  $\theta$  to represents the generic vector of parameters. In particular,  $\theta = \alpha$  in Theorem 1 and  $\theta = \gamma$  in Theorem 2. Hence, the generic log-likelihood function can be written as

$$\ell(\theta) = \sum_{i=1}^n \log \left[ [f_1(y_i|\mathbf{z}_i; \theta_1) S_2(y_i|\mathbf{z}_i; \theta_2)]^{\delta_{1i}} [f_2(y_i|\mathbf{z}_i; \theta_2) S_1(y_i|\mathbf{z}_i; \theta_1)]^{\delta_{2i}} \right]. \quad (20)$$

In (20), it has been assumed that  $\mathbf{z}_{1i} = \mathbf{z}_{2i}$ . In what follows  $\ell(\theta) = \sum_{i=1}^n \log \omega(\mathbf{w}_i; \theta)$ , where  $\omega(\mathbf{w}_i; \theta) = \omega(y|\mathbf{z}; \theta) = \left[ [f_1(y_i|\mathbf{z}_i; \theta_1) S_2(y_i|\mathbf{z}_i; \theta_2)]^{\delta_{1i}} [f_2(y_i|\mathbf{z}_i; \theta_2) S_1(y_i|\mathbf{z}_i; \theta_1)]^{\delta_{2i}} \right]$  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,  $\ell(\mathbf{w}_i; \theta) = \log \omega(\mathbf{w}_i; \theta)$ ,  $\ell_n(\theta) = n^{-1} \sum_{i=1}^n \ell(\mathbf{w}_i; \theta)$ ,  $\nabla_\theta \ell(\mathbf{w}_i; \theta) = \frac{\partial \ell(\mathbf{w}_i; \theta)}{\partial \theta}$ ,  $\nabla_\theta \ell_n(\theta) = \frac{\partial \ell_n(\theta)}{\partial \theta}$ ,  $\nabla_{\theta\theta} \ell(\mathbf{w}_i; \theta) = \frac{\partial^2 \ell(\mathbf{w}_i; \theta)}{\partial \theta \partial \theta^\top}$  and  $\nabla_{\theta\theta} \ell_n(\theta) = \frac{\partial^2 \ell_n(\theta)}{\partial \theta \partial \theta^\top}$ . The penalised likelihood is  $\ell_p(\theta) = \ell_n(\theta) - \frac{1}{2} \theta^\top \mathcal{S} \theta$ .

**Assumption 1** (Regularity Conditions).

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

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

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

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

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

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

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

This holds for any  $\boldsymbol{\theta} \in \boldsymbol{\Theta}_0$ , in particular, for  $\boldsymbol{\theta}^*$ . Setting  $\boldsymbol{\theta} = \boldsymbol{\theta}^*$ , the following equation is obtained

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

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

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

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

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

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

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

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

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

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,  $\left| \frac{\partial^3 \ell_n(\bar{\theta})}{\partial \theta_e \partial \theta_f \partial \theta_h} \right| \leq \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{w}_i) \xrightarrow{p} \mathbb{E}[\phi(\mathbf{w}_i)] < \infty$ . 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

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

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

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

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  $\theta_n$  is a sequence of  $k$ -dimensional random vectors and  $\theta^*$  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  $\Omega$ . Let  $g : \mathbb{R}^k \rightarrow \mathbb{R}^l$  be continuously differentiable at  $\theta^*$ . 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  $\alpha^*$ . 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),  $\alpha^*$  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  $\alpha^*$  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  $\Delta_p$  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

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

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,  $n^{1/2}\nabla_{\alpha\alpha}\ell_n(\alpha^*) \xrightarrow{d} \mathcal{N}[0, \mathcal{I}(\alpha^*)]$ , 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

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

as required.

- (ii) Under Theorem 1,  $\sqrt{n}(\hat{\alpha} - \alpha^*) \xrightarrow{d} \mathcal{N}\{0, [\mathcal{I}(\alpha^*)]^{-1}\}$ . In particular, for  $\alpha_{\nu 0}^* \in \hat{\alpha}$  we have  $\sqrt{n}(\alpha_{\nu 0}^* - \alpha_{\nu 0}^*) \xrightarrow{d} \mathcal{N}\{0, [\mathcal{I}(\alpha_{\nu 0}^*)]^{-1}\}$ . 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

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

Furthermore, we know that  $\nabla_{\alpha_1\alpha_2}\ell(\alpha) = 0$ , therefore  $\mathbb{E}[-\nabla_{\alpha_1\alpha_2}\ell(\alpha_0)] = 0$ . This also implies that  $\mathbb{E}[-\nabla_{\alpha_{10}\alpha_{20}}\ell(\alpha_0)] = 0$ , which means that  $\alpha_{10}$  and  $\alpha_{20}$  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  $\gamma_\nu = (\gamma_\nu^\iota, \gamma_\nu^{n\iota})^\top$  so that  $\mathcal{Q}_i^\top \gamma_\nu = \mathcal{Q}_i^{0\top} \gamma_\nu^\iota + \mathcal{Q}_{\nu i}^{1\top} \gamma_\nu^{n\iota}$ . Where  $\gamma_\nu^\iota = (\gamma_{\nu 1}^{\iota\top}, \dots, \gamma_{\nu Q}^{\iota\top})^\top$  and  $\gamma_\nu^{n\iota} = (\gamma_{\nu(Q+1)}^{n\iota\top}, \dots, \gamma_{\nu Q_\nu}^{n\iota\top})^\top$  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

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

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

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

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

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  $\boldsymbol{\Sigma}_{\hat{\boldsymbol{\theta}}} = [\boldsymbol{\mathcal{H}}_p(\hat{\boldsymbol{\theta}})]^{-1}$ . 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  $\boldsymbol{Z}$  (described in Supplementary Material C), making the large sample assumption that  $\boldsymbol{\mathcal{H}}(\boldsymbol{\theta})$  can be treated as fixed, and making the usual prior Bayesian assumption for smooth models  $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\mathcal{S}}^{-1})$ , where  $\boldsymbol{\mathcal{S}}^{-1}$  is the Moore-Penrose pseudoinverse of  $\boldsymbol{\mathcal{S}}$  (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  $\boldsymbol{\Sigma}_{\tilde{\boldsymbol{\beta}}_{\nu 0}} = \text{diag}(\boldsymbol{\Gamma}_{\nu 0}) [\boldsymbol{\mathcal{H}}_p(\hat{\tilde{\boldsymbol{\beta}}}_{\nu 0})]^{-1} \text{diag}(\boldsymbol{\Gamma}_{\nu 0})$ . 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 ( $\Upsilon^{\text{KCV}}$ ). 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  $\Upsilon^{\text{KCV}}$  (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  $\hat{\boldsymbol{\alpha}}_\nu^{\setminus k}$ , 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  $\Upsilon^{\text{KCV}}$  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  $\Upsilon^{\text{KCV}}$  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  $\Upsilon^{\text{KCV}} = \sum_{k=1}^K \ell_k(\hat{\gamma}^{\setminus k})$ .

Model	Non-Inf. Covariates	Inf. Covariates	Link T <sub>1i</sub>	Link T <sub>2i</sub>	AIC	$\Upsilon^{\text{KCV}}$	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  $\Upsilon^{\text{KCV}}$ ) 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,  $s_{11}(z_{2i}) = s_{12}(z_{2i}) = \sin(2\pi z_i)$ ,  $\alpha_{01} = -0.10$ ,  $\alpha_{02} = -0.25$  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.

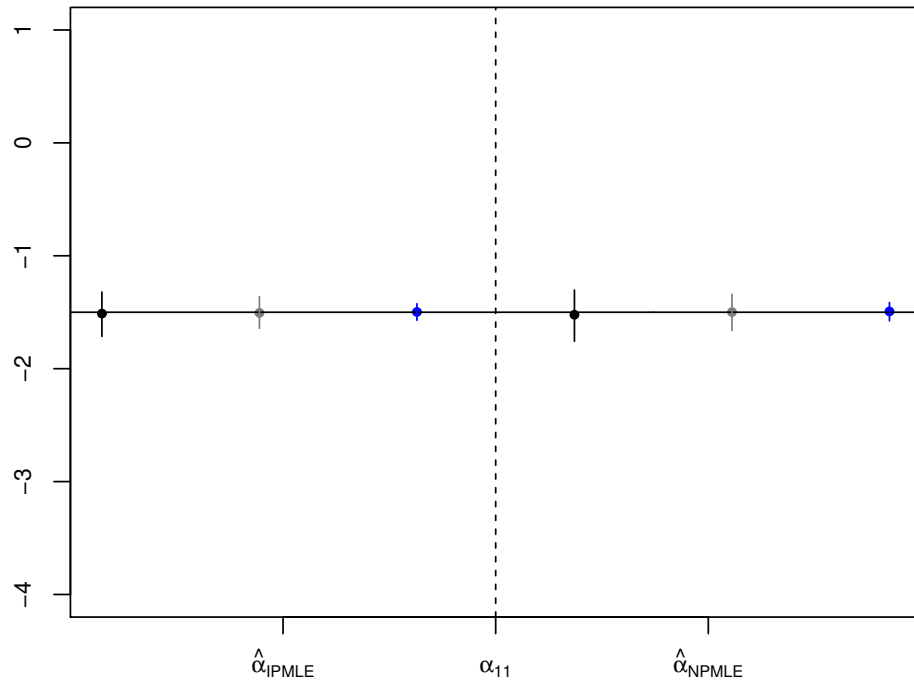


Figure 1: Linear coefficient estimates obtained by applying `gamlss()` to informative survival data simulated according to DGP3 characterised by a censoring rate of about 47%. Circles indicate mean estimates while bars represent the estimates' ranges resulting from 5% and 95% quantiles. True values are 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.

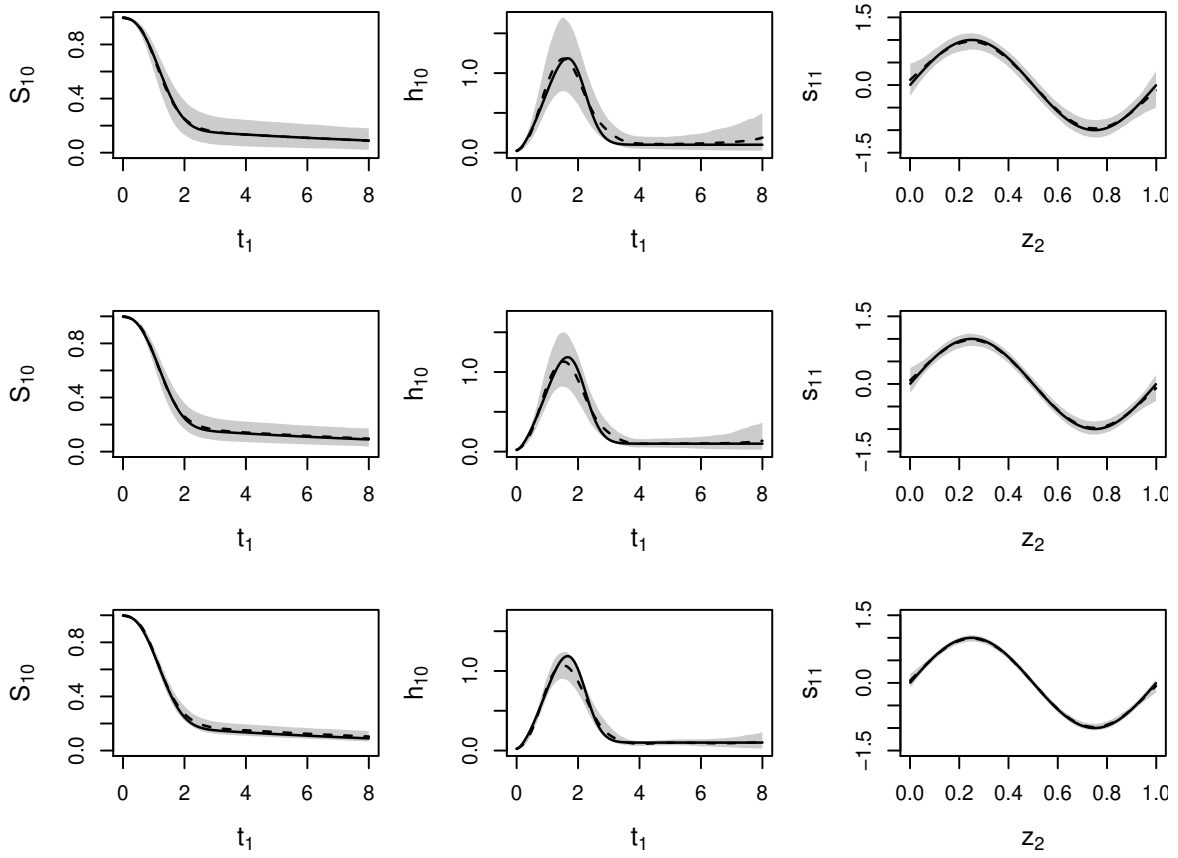


Figure 2: Smooth function estimates for the IPMLE obtained by applying `gamlss()` to informative survival data simulated according to DGP3 characterised by a censoring rate of about 47%. True functions are represented by black solid lines, mean estimates by dashed lines and pointwise ranges resulting 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
$\alpha_{11}$	-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
$\alpha_{11}$	-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.

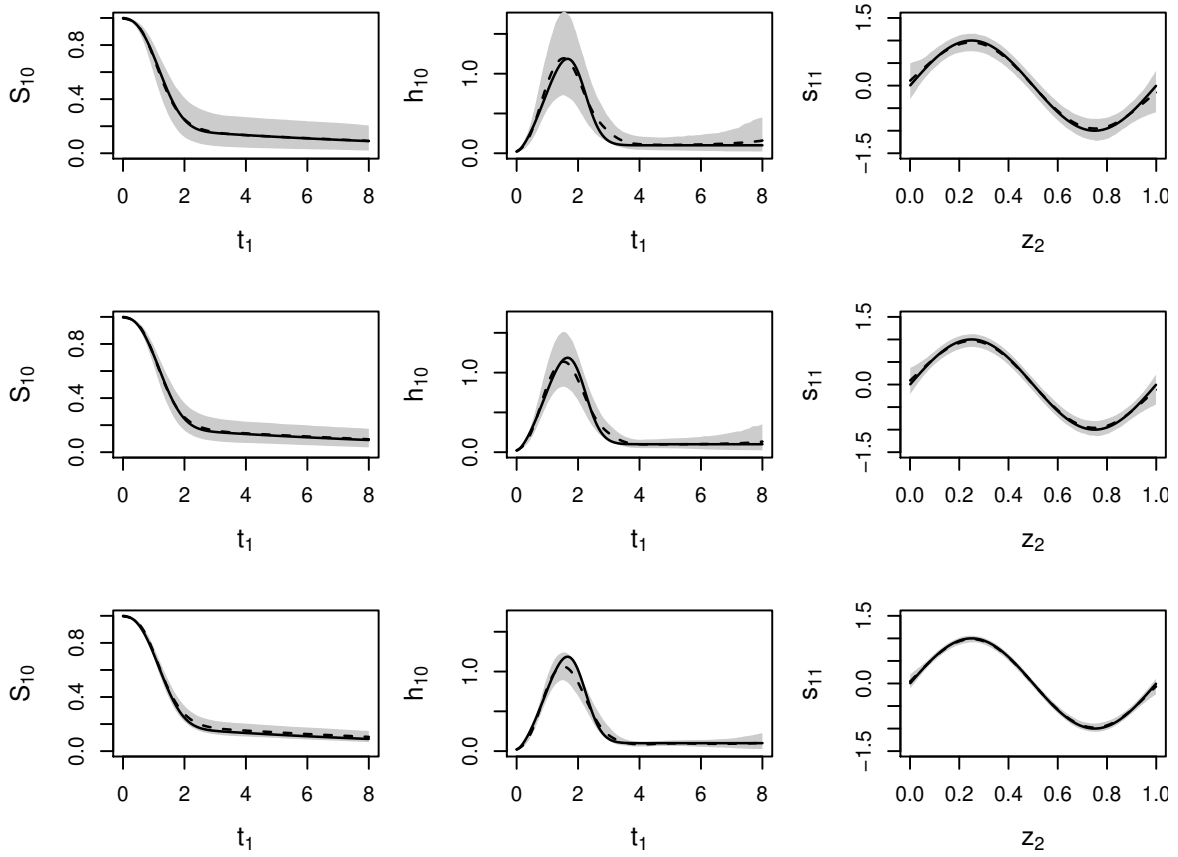


Figure 3: Smooth function estimates for the NPMLE obtained by applying `gamlss()` to informative survival data simulated according to DGP3 characterised by a censoring rate of about 47%. Further details are given in the caption of Figure 2.

(a) Informative Penalized Maximum Log-likelihood Estimator (IPMLE)						
	Bias			RMSE		
	n = 500	n = 1000	n = 4000	n = 500	n = 1000	n = 4000
$\alpha_{11}$	-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
$\alpha_{11}$	-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.

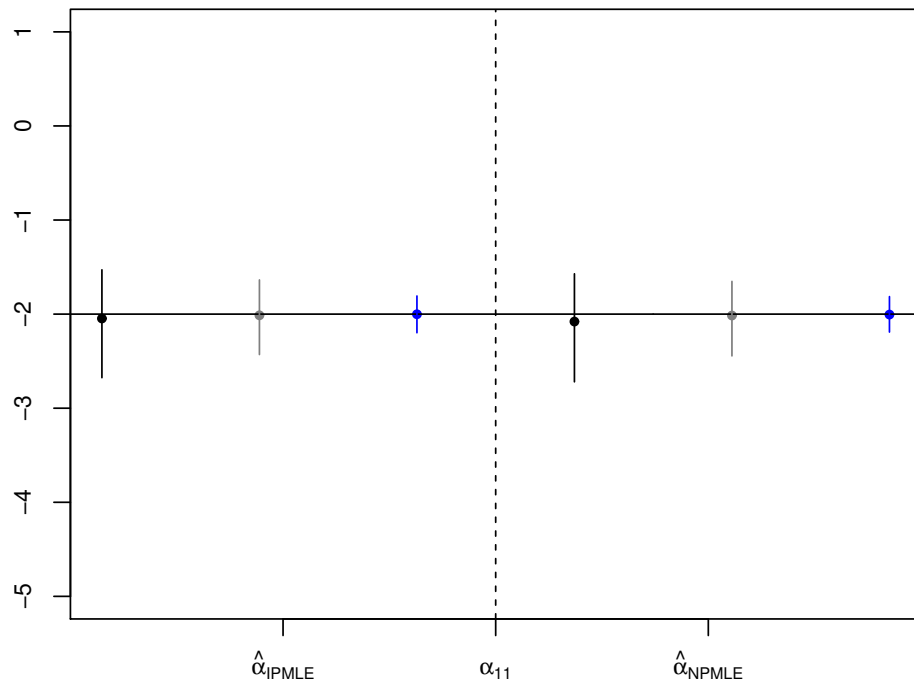


Figure 4: Linear coefficient estimates obtained by applying `gamlss()` to informative survival data simulated according to DGP1 which is characterised by a censoring rate of about 78%. Further details are given in the caption of Figure 1.

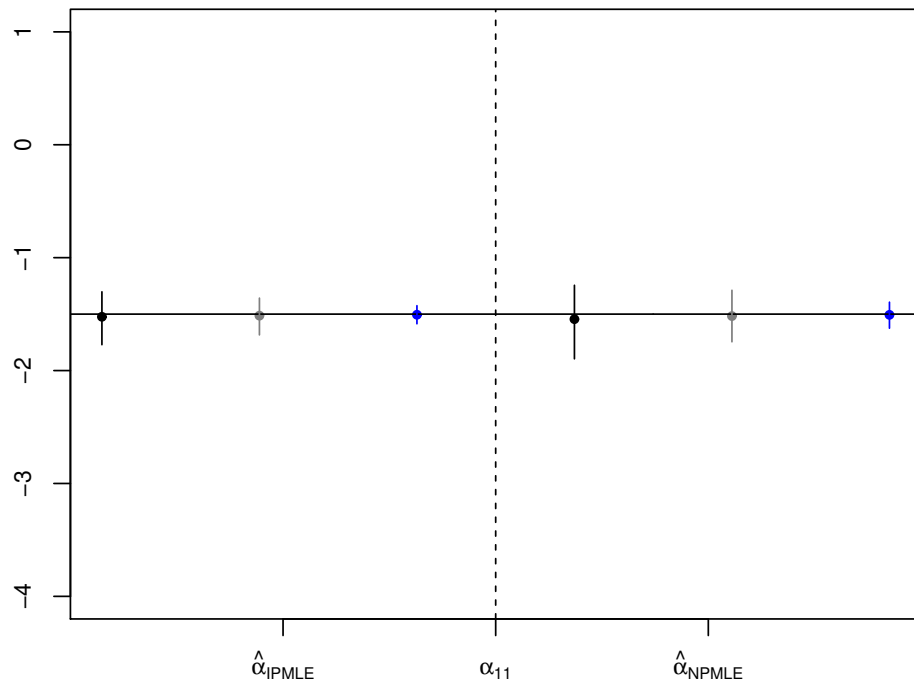


Figure 5: Linear coefficient estimates obtained by applying `gamlss()` to informative survival data simulated according to DGP2 which is characterised by a censoring rate of about 74%. Further details are given in the caption of Figure 1.

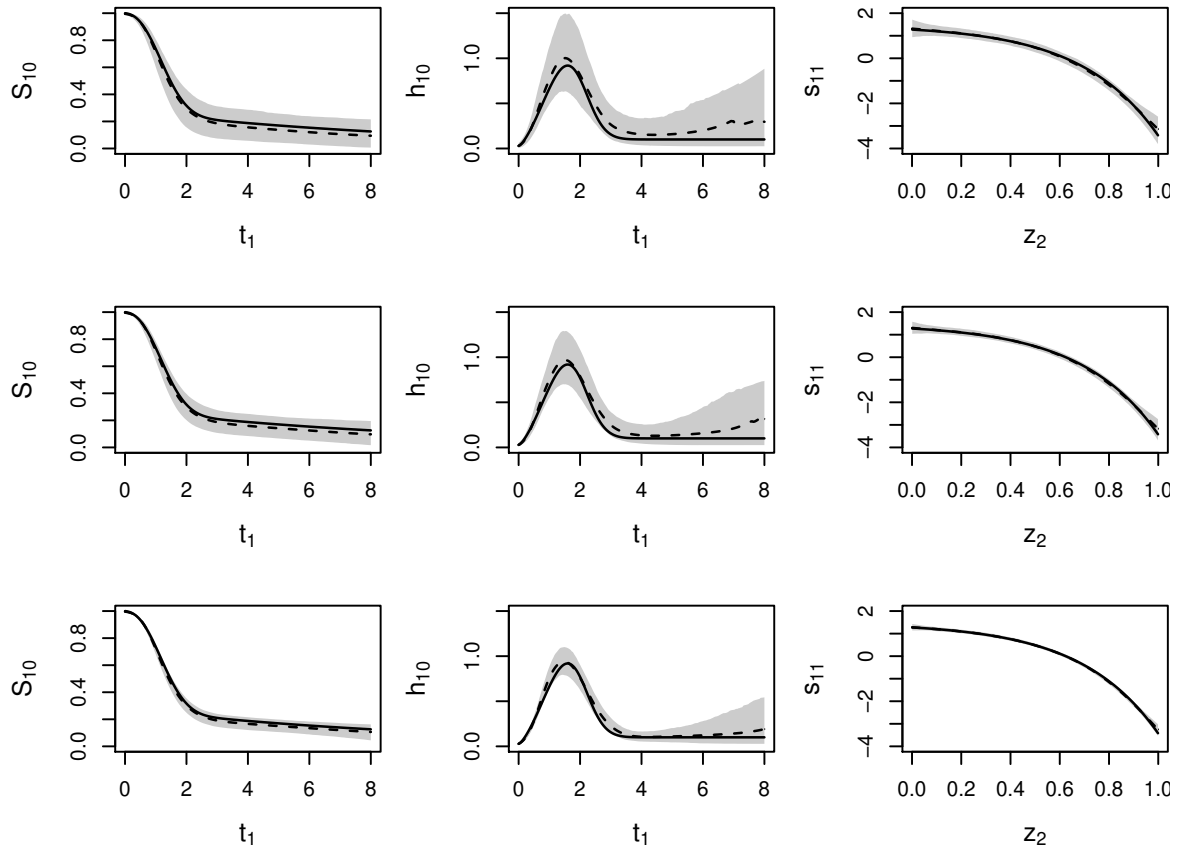


Figure 6: Smooth function estimates for the IPMLE obtained by applying `gam1ss()` to informative survival data simulated according to DGP1 characterised by a censoring rate of about 78%. Further details are given in the caption of Figure 2.

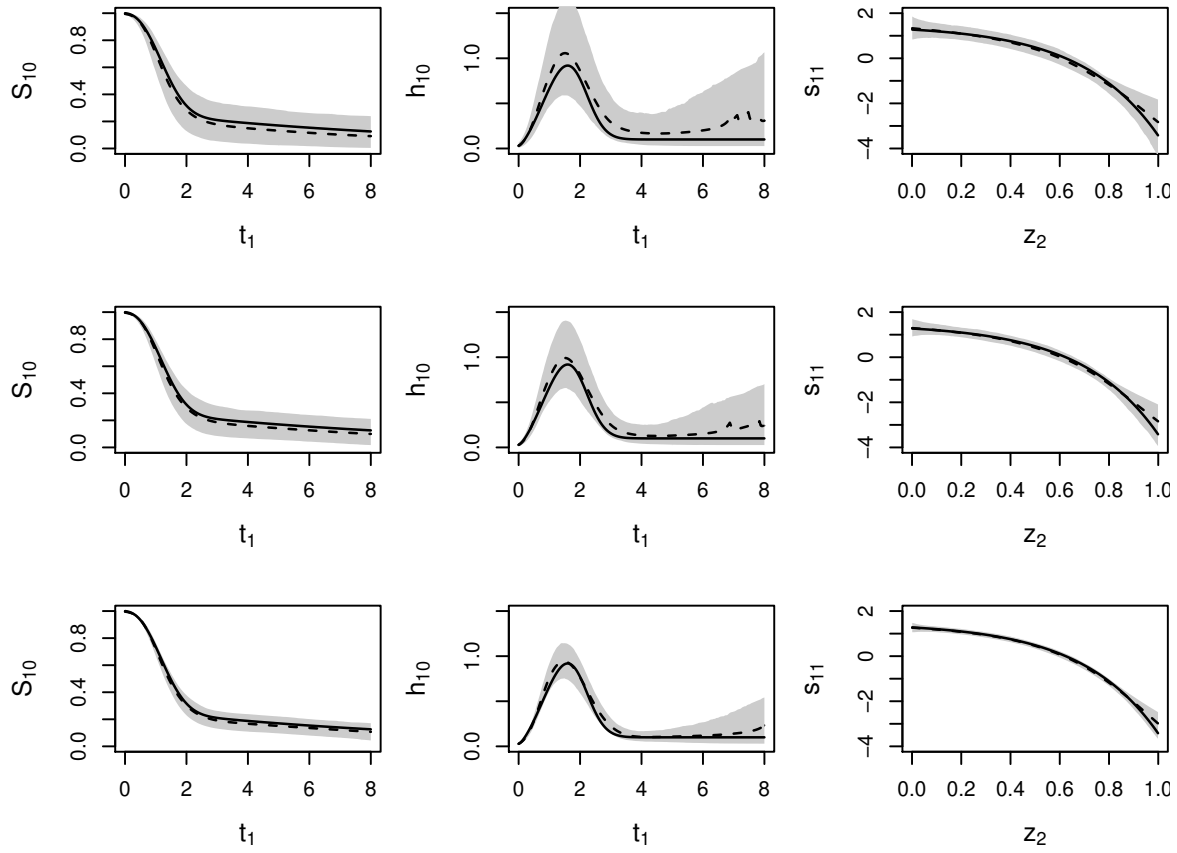


Figure 7: Smooth function estimates for the NPMLE obtained by applying `gamlss()` to informative survival data simulated according to DGP1 characterised by a censoring rate of about 78%. Further details are given in the caption of Figure 2.

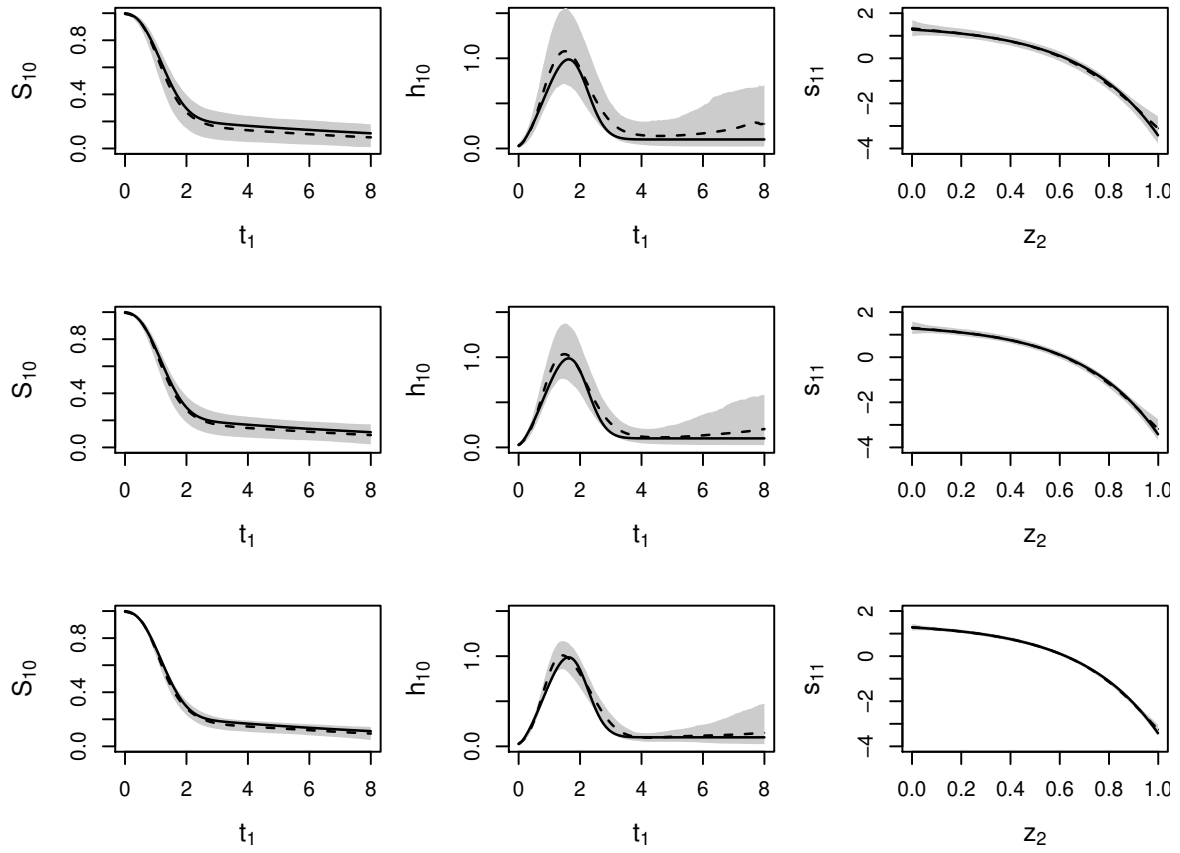


Figure 8: Smooth function estimates for the IPMLE obtained by applying `gamlss()` to informative survival data simulated according to DGP2 characterised by a censoring rate of about 74%. Further details are given in the caption of Figure 2

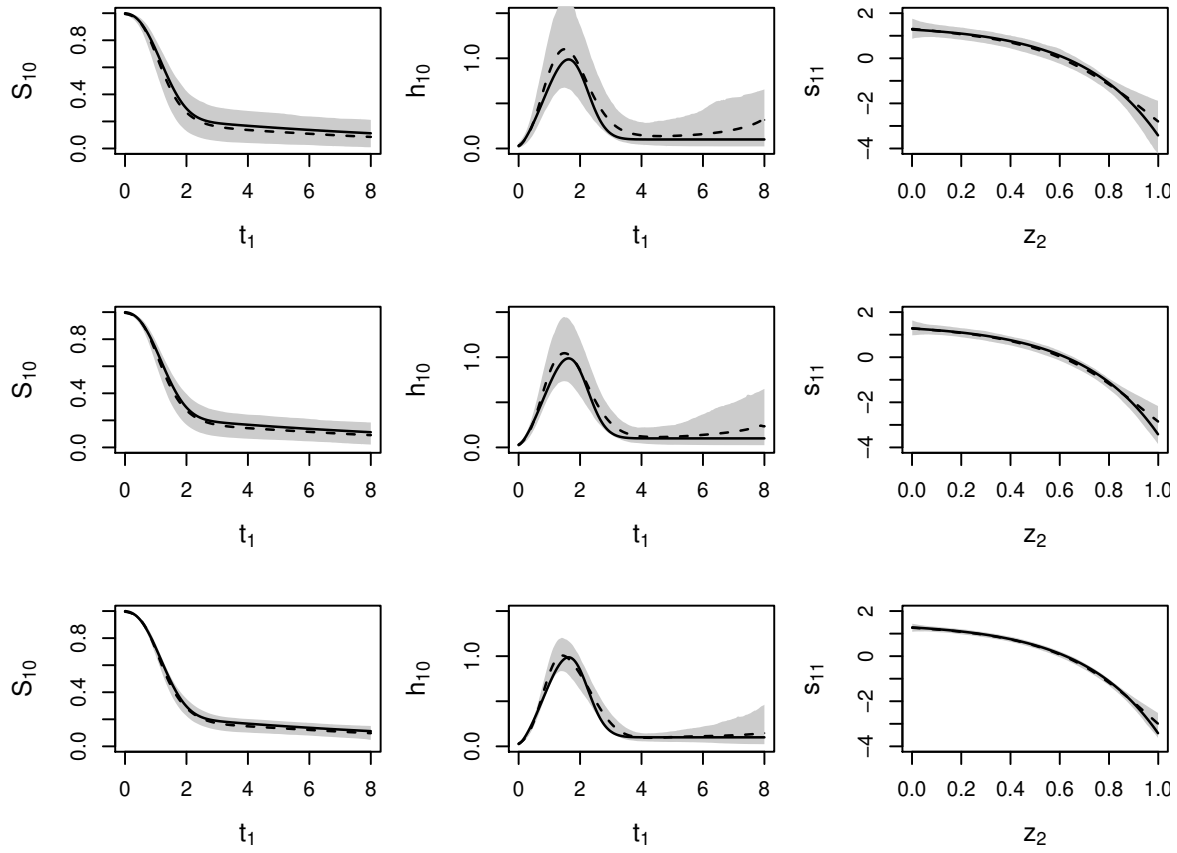


Figure 9: Smooth function estimates for the NPMLE obtained by applying `gamlss()` to informative survival data simulated according to DGP2 characterised by a censoring rate of about 74%. Further details are given in the caption of Figure 2

## References

- Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. *In: Petrov, B.N., Csaki, B.F. (eds.) Second International Symposium on Information Theory. Academiai Kiado, Budapest.*
- Gourieroux, C. & Monfort, A. (1995). *Statistics and econometric models*, volume 1. Cambridge University Press.
- Gu, C. (1992). Cross-validating non-gaussian data. *Journal of Computational and Graphical Statistics*, 1(2), 169–179.
- Hayashi, F. (2000). *Econometrics*. Princeton University Press Princeton, NJ.
- Marra, G. & Radice, R. (2019). *GJRM: Generalised Joint Regression Modelling*. R package version 0.2-1.
- Marra, G., Radice, R., Bärnighausen, T., Wood, S. N., & McGovern, M. E. (2017). A simultaneous equation approach to estimating hiv prevalence with nonignorable missing responses. *Journal of the American Statistical Association*, 112(518), 484–496.
- Marra, G. & Wood, S. N. (2012). Coverage properties of confidence intervals for generalized additive model components. *Scandinavian Journal of Statistics*, 39(1), 53–74.
- Newey, W. K. & McFadden, D. (1994). *Handbook of econometrics*, volume 4. Elsevier.
- Nocedal, J. & Wright, S. (2006). *Numerical optimization, series in operations research and financial engineering*. Springer, New York, USA, 2006.
- Pya, N. & Wood, S. N. (2015). Shape constrained additive models. *Statistics and Computing*, 25(3), 543–559.
- Silverman, B. W. (1985). Some aspects of the spline smoothing approach to non-parametric regression curve fitting. *Journal of the Royal Statistical Society: Series B (Methodological)*, 47(1), 1–21.

- Stone, M. (1974). Cross-validatory choice and assessment of statistical predictions. *Journal of the royal statistical society. Series B (Methodological)*, 36(2), 111–147.
- Vatter, T. & Chavez-Demoulin, V. (2015). Generalized additive models for conditional dependence structures. *Journal of Multivariate Analysis*, 141, 147–167.
- Wood, S. N. (2004). Stable and efficient multiple smoothing parameter estimation for generalized additive models. *Journal of the American Statistical Association*, 99(467), 673–686.
- Wood, S. N. (2013). On p-values for smooth components of an extended generalized additive model. *Biometrika*, 100(1), 221–228.
- Wood, S. N. (2017). *Generalized Additive Models: An Introduction With R*. Second Edition, Chapman & Hall/CRC, London.
- Xingwei, T., Tao, H., & Hengjian, C. (2010). Hazard regression with penalized spline: The smoothing parameter choice and asymptotics. *Acta Mathematica Scientia*, 30(5), 1759–1768.