

INFLUENCE DIAGNOSTICS FOR ROBUST P-SPLINES USING SCALE MIXTURE OF NORMAL DISTRIBUTIONS

FELIPE OSORIO

SUPPLEMENTARY MATERIAL

Appendix A presents details about the computational implementation. Proofs of Propositions 1 to 5 are established in Appendix B. Additional results obtained from a real data analysis and a Monte Carlo simulation study are presented in the Appendices C and D, respectively.

APPENDIX A. NOTES ABOUT THE COMPUTATIONAL IMPLEMENTATION

Below we describe the computational strategy adopted in the `heavy` R package that uses the estimation procedure proposed in this work. A set of routines has been written in C in order to accelerate some of the calculations. The implementation is fairly simple and most of the cases has been possible to draw upon the routines from the BLAS library (Lawson et al., 1979), Linpack (Dongarra et al., 1979) and Mathlib, included in the R software (R Core Team, 2014).

The M step of the algorithm, described in Equations (9) and (10) from Section 2.1 can be efficiently computed by considering $\mathbf{U} = \mathbf{W}^{(k)1/2} \mathbf{B}$ and decompose \mathbf{U} using a singular value decomposition like \mathbf{UDV}^\top , where \mathbf{U} is $n \times p$ matrix such that $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_p$, $\mathbf{D} = \text{diag}(\delta_1, \dots, \delta_p)$ and \mathbf{V} is an orthogonal $p \times p$ matrix. It is important to note that the form of decomposition highlights that it is possible to overwrite the matrix \mathbf{U} in order to keep storage to a minimum. A second singular value decomposition is then calculated:

$$\mathbf{KVD}^{-1} = \mathbf{QRS}^\top,$$

where only $\mathbf{R} = \text{diag}(r_1, \dots, r_p)$ and the orthogonal $p \times p$ matrix \mathbf{S} are stored. Let $\mathbf{Z} = \mathbf{U}^\top \mathbf{W}^{(k)1/2} \mathbf{Y}$ and $\mathbf{c} = \mathbf{S}^\top \mathbf{Z}$. Thus, we define

$$\boldsymbol{\alpha}_\lambda^{(k+1)} = \mathbf{S}(\mathbf{I} + \lambda \mathbf{R}^2)^{-1} \mathbf{c}, \quad \lambda > 0, \tag{A.1}$$

and proceed to calculate the fitted values and residuals for the current fit as

$$\mathbf{g}_\lambda^{(k+1)} = \mathbf{U} \boldsymbol{\alpha}_\lambda^{(k+1)} \quad \text{and} \quad \mathbf{e}_\lambda^{(k+1)} = \mathbf{Z} - \mathbf{g}_\lambda^{(k+1)},$$

respectively. Hence, Equation (10) adopts the form

$$\phi_\lambda^{(k+1)} = \frac{1}{n} \{ \|\mathbf{e}_\lambda^{(k+1)}\|^2 + \lambda \|\mathbf{R}(\mathbf{I} + \lambda \mathbf{R}^2)^{-1} \mathbf{c}\|^2 \}. \tag{A.2}$$

The author was supported by grants CONICYT 791100007 and FONDECYT 1140580.

Using the results above, it is possible to rewrite the WGCV criterion defined in (14) from Section 2.2 for the smoothing parameter selection as

$$V(\lambda) = \frac{1}{n} \frac{\|e_\lambda^{(k+1)}\|^2}{(1 - \text{edf}/n)^2}, \quad (\text{A.3})$$

where edf represents the effective number of parameters, defined as $\text{edf} = \text{tr } \mathbf{H}_W(\lambda)$, which can be calculated efficiently as

$$\text{edf} = \text{tr}(\mathbf{I} + \lambda \mathbf{R}^2)^{-1} = \sum_{j=1}^p \frac{1}{1 + \lambda r_j^2}.$$

We use the Brent's method (Brent, 1973) to carry out the minimization of $V(\lambda)$ given in (A.3). Finally, once convergences of the algorithm was reached, we take

$$\hat{\mathbf{a}}_\lambda = \mathbf{V} \mathbf{D}^{-1} \hat{\boldsymbol{\alpha}}_\lambda, \quad \hat{\mathbf{g}}_\lambda = \widehat{\mathbf{W}}^{-1/2} \mathbf{U} \hat{\boldsymbol{\alpha}}_\lambda.$$

The estimation procedure according to the algorithm described in this section is available in the `heavyPS` function of `heavy` package.

APPENDIX B. PROOFS OF MAIN RESULTS

In this appendix we derive the differentials $d_\theta^2 Q_\lambda(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}})$, $d_{\omega\theta}^2 Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\hat{\boldsymbol{\theta}})$ and $d_{\omega\theta}^2 V(\lambda, \boldsymbol{\omega})$ for scale and response perturbation schemes. The necessary matrices $\ddot{Q}(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}})$, $\boldsymbol{\Delta}(\boldsymbol{\omega})$ and the vector $\partial^2 V(\lambda, \boldsymbol{\omega})/\partial\lambda\partial\boldsymbol{\omega}^\top$ are obtained efficiently using the differentiation method and by applying some identification theorems discussed in Magnus and Neudecker (1999).

Proof of Proposition 1. The complete data penalized log-likelihood function given in Equation (6) from Section 2.1 is

$$Q_\lambda(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}) = -\frac{n}{2} \log \phi - \frac{1}{2\phi} [(\mathbf{Y} - \mathbf{B}\mathbf{a})^\top \widehat{\mathbf{W}}(\mathbf{Y} - \mathbf{B}\mathbf{a}) + \lambda \mathbf{a}^\top \mathbf{K}^\top \mathbf{K} \mathbf{a}],$$

Differentiating $Q_\lambda(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}})$ with respect to \mathbf{a} , we obtain

$$d_a Q_\lambda(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}) = \frac{1}{\phi} \{(\mathbf{Y} - \mathbf{B}\mathbf{a})^\top \widehat{\mathbf{W}} \mathbf{B} - \lambda \mathbf{a}^\top \mathbf{K}^\top \mathbf{K}\} d\mathbf{a},$$

$$d_a^2 Q_\lambda(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}) = -\frac{1}{\phi} (d\mathbf{a})^\top \{\mathbf{B}^\top \widehat{\mathbf{W}} \mathbf{B} + \lambda \mathbf{K}^\top \mathbf{K}\} d\mathbf{a},$$

and the differential of $d_a Q_\lambda(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}})$ with respect to ϕ leads to

$$d_{\phi a}^2 Q_\lambda(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}) = -\frac{1}{\phi^2} d\phi \{(\mathbf{Y} - \mathbf{B}\mathbf{a})^\top \widehat{\mathbf{W}} \mathbf{B} - \lambda \mathbf{a}^\top \mathbf{K}^\top \mathbf{K}\} d\mathbf{a}.$$

In addition, the differentials of $Q_\lambda(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}})$ with respect to ϕ are given by

$$d_\phi Q_\lambda(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}) = -\frac{n}{2\phi} d\phi + \frac{1}{2\phi^2} \{(\mathbf{Y} - \mathbf{B}\mathbf{a})^\top \widehat{\mathbf{W}}(\mathbf{Y} - \mathbf{B}\mathbf{a}) + \lambda \mathbf{a}^\top \mathbf{K}^\top \mathbf{K} \mathbf{a}\} d\phi,$$

$$d_\phi^2 Q_\lambda(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}) = \frac{n}{2\phi^2} d^2\phi - \frac{1}{\phi^3} \{(\mathbf{Y} - \mathbf{B}\mathbf{a})^\top \widehat{\mathbf{W}}(\mathbf{Y} - \mathbf{B}\mathbf{a}) + \lambda \mathbf{a}^\top \mathbf{K}^\top \mathbf{K} \mathbf{a}\} d^2\phi.$$

From the first-order conditions

$$\mathbf{B}^\top \widehat{\mathbf{W}} \mathbf{Y} - (\mathbf{B}^\top \widehat{\mathbf{W}} \mathbf{B} + \lambda \mathbf{K}^\top \mathbf{K}) \hat{\mathbf{a}} = \mathbf{0}, \quad (\text{B.1})$$

$$n\hat{\phi} - (RSS_{\widehat{\mathbf{W}}}(\hat{\mathbf{a}}) + \lambda \hat{\mathbf{a}}^\top \mathbf{K}^\top \mathbf{K} \hat{\mathbf{a}}) = 0, \quad (\text{B.2})$$

follow immediately that

$$\begin{aligned} d_a^2 Q_\lambda(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} &= -\frac{1}{\hat{\phi}}(d\mathbf{a})^\top \{\mathbf{B}^\top \widehat{\mathbf{W}} \mathbf{B} + \lambda \mathbf{K}^\top \mathbf{K}\} d\mathbf{a}, \\ d_{\phi a}^2 Q_\lambda(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} &= \mathbf{0}, \quad d_\phi^2 Q_\lambda(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} = -\frac{n}{2\hat{\phi}^2} d^2\phi. \end{aligned}$$

Applying the second identification theorem given in Magnus and Neudecker (1999) we obtain $\ddot{Q}(\hat{\boldsymbol{\theta}}|\hat{\boldsymbol{\theta}})$, and the proof is complete. \square

Proof of Proposition 2. The Q -function for the perturbed model introduced in Equation (18) from Subsection 3.2.1 assumes the form

$$Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\hat{\boldsymbol{\theta}}) = -\frac{n}{2} \log \phi - \frac{1}{2\hat{\phi}} [(\mathbf{Y} - \mathbf{B}\mathbf{a})^\top \widehat{\mathbf{W}}^{1/2} \text{diag}(\boldsymbol{\omega}) \widehat{\mathbf{W}}^{1/2} (\mathbf{Y} - \mathbf{B}\mathbf{a}) + \lambda \mathbf{a}^\top \mathbf{K}^\top \mathbf{K} \mathbf{a}],$$

where $\boldsymbol{\omega}$ is an n -dimensional perturbation vector, the perturbed model reduces to the postulated model when $\boldsymbol{\omega}_0 = \mathbf{1}$.

Taking differentials of the $Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\hat{\boldsymbol{\theta}})$ with respect to $\boldsymbol{\theta} = (\mathbf{a}^\top, \phi)^\top$, we get

$$\begin{aligned} d_a Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\hat{\boldsymbol{\theta}}) &= \frac{1}{\hat{\phi}} \{(\mathbf{Y} - \mathbf{B}\mathbf{a})^\top \widehat{\mathbf{W}} \text{diag}(\boldsymbol{\omega}) \mathbf{B} - \lambda \mathbf{a}^\top \mathbf{K}^\top \mathbf{K}\} d\mathbf{a}, \\ d_\phi Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\hat{\boldsymbol{\theta}}) &= -\frac{n}{2\hat{\phi}} d\phi + \frac{1}{2\hat{\phi}^2} \{(\mathbf{Y} - \mathbf{B}\mathbf{a})^\top \widehat{\mathbf{W}}^{1/2} \text{diag}(\boldsymbol{\omega}) \widehat{\mathbf{W}}^{1/2} (\mathbf{Y} - \mathbf{B}\mathbf{a}) \\ &\quad + \lambda \mathbf{a}^\top \mathbf{K}^\top \mathbf{K} \mathbf{a}\} d\phi, \end{aligned}$$

using that $\mathbf{z}^\top \text{diag}(\boldsymbol{\omega}) = \boldsymbol{\omega}^\top \text{diag}(\mathbf{z})$ and $\text{diag}(\boldsymbol{\omega}) \text{diag}(\mathbf{z}) = \text{diag}(\mathbf{z}) \text{diag}(\boldsymbol{\omega})$ for \mathbf{z} an n -dimensional vector, the differentials $d_a Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\hat{\boldsymbol{\theta}})$ and $d_\phi Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\hat{\boldsymbol{\theta}})$ can be written as

$$\begin{aligned} d_a Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\hat{\boldsymbol{\theta}}) &= \frac{1}{\hat{\phi}} \{\boldsymbol{\omega}^\top \text{diag}(\boldsymbol{\epsilon}) \widehat{\mathbf{W}} \mathbf{B} - \lambda \mathbf{a}^\top \mathbf{K}^\top \mathbf{K}\} d\mathbf{a}, \\ d_\phi Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\hat{\boldsymbol{\theta}}) &= -\frac{n}{2\hat{\phi}} d\phi + \frac{1}{2\hat{\phi}^2} \{\boldsymbol{\omega}^\top \text{diag}(\boldsymbol{\epsilon}) \widehat{\mathbf{W}} \boldsymbol{\epsilon} + \lambda \mathbf{a}^\top \mathbf{K}^\top \mathbf{K} \mathbf{a}\} d\phi, \end{aligned}$$

where $\boldsymbol{\epsilon} = \mathbf{Y} - \mathbf{B}\mathbf{a}$. To find $\boldsymbol{\Delta}(\boldsymbol{\omega})$, we differentiate $d_a Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\hat{\boldsymbol{\theta}})$ and $d_\phi Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\hat{\boldsymbol{\theta}})$ with respect to $\boldsymbol{\omega}$, thus

$$d_{\omega a}^2 Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\hat{\boldsymbol{\theta}}) = \frac{1}{\hat{\phi}} (d\boldsymbol{\omega})^\top \text{diag}(\boldsymbol{\epsilon}) \widehat{\mathbf{W}} \mathbf{B} d\mathbf{a}, \quad (\text{B.3})$$

$$d_{\omega \phi}^2 Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\hat{\boldsymbol{\theta}}) = \frac{1}{2\hat{\phi}^2} (d\boldsymbol{\omega})^\top \text{diag}(\boldsymbol{\epsilon}) \widehat{\mathbf{W}} \boldsymbol{\epsilon} d\phi. \quad (\text{B.4})$$

Applying the second identification theorem given in Magnus and Neudecker (1999) and evaluating (B.3) and (B.4) at $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}$ and $\boldsymbol{\omega} = \boldsymbol{\omega}_0$ we obtain $\boldsymbol{\Delta}(\boldsymbol{\omega}_0)$, and the proof is complete. \square

Proof of Proposition 3. Under the response perturbation $\mathbf{Y}(\boldsymbol{\omega}) = \mathbf{Y} + \boldsymbol{\omega}$, the Q -function for the perturbed model is given by

$$Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\hat{\boldsymbol{\theta}}) = -\frac{n}{2} \log \phi - \frac{1}{2\hat{\phi}} [(\mathbf{Y}(\boldsymbol{\omega}) - \mathbf{B}\mathbf{a})^\top \widehat{\mathbf{W}} (\mathbf{Y}(\boldsymbol{\omega}) - \mathbf{B}\mathbf{a}) + \lambda \mathbf{a}^\top \mathbf{K}^\top \mathbf{K} \mathbf{a}],$$

where $\boldsymbol{\omega}$ is an n -dimensional perturbation vector and the vector of null perturbation vector is given by $\boldsymbol{\omega}_0 = \mathbf{0}$.

The first differential of $Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\widehat{\boldsymbol{\theta}})$ with respect to $\boldsymbol{\theta} = (\mathbf{a}^\top, \phi)^\top$ is

$$d_a Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\widehat{\boldsymbol{\theta}}) = \frac{1}{\phi} \{(\mathbf{Y}(\boldsymbol{\omega}) - \mathbf{B}\mathbf{a})^\top \widehat{\mathbf{W}}\mathbf{B} - \lambda \mathbf{a}^\top \mathbf{K}^\top \mathbf{K}\} d\mathbf{a},$$

$$d_\phi Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\widehat{\boldsymbol{\theta}}) = -\frac{n}{2\phi} d\phi + \frac{1}{2\phi^2} \{(\mathbf{Y}(\boldsymbol{\omega}) - \mathbf{B}\mathbf{a})^\top \widehat{\mathbf{W}}(\mathbf{Y}(\boldsymbol{\omega}) - \mathbf{B}\mathbf{a}) + \lambda \mathbf{a}^\top \mathbf{K}^\top \mathbf{K}\mathbf{a}\} d\phi,$$

taking the differential of $d_a Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\widehat{\boldsymbol{\theta}})$ and $d_\phi Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\widehat{\boldsymbol{\theta}})$ with respect to $\boldsymbol{\omega}$ we have

$$d_{\boldsymbol{\omega}a}^2 Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\widehat{\boldsymbol{\theta}}) = \frac{1}{\phi} (d\boldsymbol{\omega})^\top \widehat{\mathbf{W}}\mathbf{B} d\mathbf{a}, \quad (\text{B.5})$$

$$d_{\boldsymbol{\omega}\phi}^2 Q_\lambda(\boldsymbol{\theta}, \boldsymbol{\omega}|\widehat{\boldsymbol{\theta}}) = \frac{1}{2\phi^2} (d\boldsymbol{\omega})^\top \widehat{\mathbf{W}}(\mathbf{Y}(\boldsymbol{\omega}) - \mathbf{B}\mathbf{a}) d\phi. \quad (\text{B.6})$$

The $\boldsymbol{\Delta}(\boldsymbol{\omega}_0)$ matrix can be obtained by applying the second identification theorem (Magnus and Neudecker, 1999) and evaluating (B.5) and (B.6) at $\boldsymbol{\theta} = \widehat{\boldsymbol{\theta}}$ and $\boldsymbol{\omega} = \boldsymbol{\omega}_0$. Thus, the proposition is verified. \square

Proof of Proposition 4. The weighted cross-validation criterion under the scale perturbation defined by Equation (19) from Subsection 3.2.3 is given by

$$V(\lambda, \boldsymbol{\omega}) = \frac{RSS_{\widehat{\mathbf{W}}}(\lambda, \boldsymbol{\omega})/n}{\{\text{tr}(\mathbf{I} - \mathbf{H}_{\widehat{\mathbf{W}}}(\lambda, \boldsymbol{\omega}))\}/n\}^2},$$

where

$$RSS_{\widehat{\mathbf{W}}}(\lambda, \boldsymbol{\omega}) = \mathbf{Y}^\top (\mathbf{I} - \mathbf{H}_{\widehat{\mathbf{W}}}(\lambda, \boldsymbol{\omega}))^\top \widehat{\mathbf{W}}^{1/2} \text{diag}(\boldsymbol{\omega}) \widehat{\mathbf{W}}^{1/2} (\mathbf{I} - \mathbf{H}_{\widehat{\mathbf{W}}}(\lambda, \boldsymbol{\omega})) \mathbf{Y},$$

and

$$\mathbf{H}_{\widehat{\mathbf{W}}}(\lambda, \boldsymbol{\omega}) = \mathbf{B}(\mathbf{B}^\top \widehat{\mathbf{W}} \text{diag}(\boldsymbol{\omega}) \mathbf{B} + \lambda \mathbf{K}^\top \mathbf{K})^{-1} \mathbf{B}^\top \widehat{\mathbf{W}} \text{diag}(\boldsymbol{\omega}).$$

Next, we shall write $RSS_{\widehat{\mathbf{W}}}(\lambda, \boldsymbol{\omega}) = RSS(\lambda, \boldsymbol{\omega})$ and $\mathbf{H}_{\widehat{\mathbf{W}}}(\lambda, \boldsymbol{\omega}) = \mathbf{H}(\lambda, \boldsymbol{\omega})$. Let $\text{edf}(\lambda, \boldsymbol{\omega}) = \text{tr} \mathbf{H}(\lambda, \boldsymbol{\omega})$ be the effective number of parameters for the perturbed model. Taking the differential of $V(\lambda, \boldsymbol{\omega})$ with respect to $\boldsymbol{\omega}$ leads to

$$d_\omega V(\lambda, \boldsymbol{\omega}) = \frac{1/n}{(1 - \text{edf}(\lambda, \boldsymbol{\omega})/n)^4} \left\{ (1 - \text{edf}(\lambda, \boldsymbol{\omega})/n)^2 d_\omega RSS(\lambda, \boldsymbol{\omega}) - RSS(\lambda, \boldsymbol{\omega}) d_\omega (1 - \text{edf}(\lambda, \boldsymbol{\omega})/n)^2 \right\}, \quad (\text{B.7})$$

is easy to see that,

$$\begin{aligned} d_\omega \mathbf{H}(\lambda, \boldsymbol{\omega}) &= -\mathbf{B}\mathbf{S}^{-1} \mathbf{B}^\top \widehat{\mathbf{W}} \text{diag}(d\boldsymbol{\omega}) \mathbf{B}\mathbf{S}^{-1} \mathbf{B}^\top \widehat{\mathbf{W}} \text{diag}(\boldsymbol{\omega}) + \mathbf{B}\mathbf{S}^{-1} \mathbf{B}^\top \widehat{\mathbf{W}} \text{diag}(d\boldsymbol{\omega}) \\ &= \mathbf{B}\mathbf{S}^{-1} \mathbf{B}^\top \widehat{\mathbf{W}} \text{diag}(d\boldsymbol{\omega}) (\mathbf{I} - \mathbf{H}(\lambda, \boldsymbol{\omega})) \end{aligned} \quad (\text{B.8})$$

where $\mathbf{S} = \mathbf{B}^\top \widehat{\mathbf{W}} \text{diag}(\boldsymbol{\omega}) \mathbf{B} + \lambda \mathbf{K}^\top \mathbf{K}$. Thus, using the properties of the trace operator, we obtain

$$d_\omega (1 - \text{edf}(\lambda, \boldsymbol{\omega})/n)^2 = -\frac{2}{n} (1 - \text{edf}(\lambda, \boldsymbol{\omega})/n) \times \mathbf{1}^\top \text{dg}((\mathbf{I} - \mathbf{H}(\lambda, \boldsymbol{\omega})) \mathbf{B}\mathbf{S}^{-1} \mathbf{B}^\top \widehat{\mathbf{W}}) d\boldsymbol{\omega}. \quad (\text{B.9})$$

Let $\mathbf{e}_\omega = (\mathbf{I} - \mathbf{H}(\lambda, \boldsymbol{\omega})) \mathbf{Y}$, direct calculations show that

$$\begin{aligned} d_\omega RSS(\lambda, \boldsymbol{\omega}) &= -\mathbf{Y}^\top (d_\omega \mathbf{H}(\lambda, \boldsymbol{\omega}))^\top \widehat{\mathbf{W}}^{1/2} \text{diag}(\boldsymbol{\omega}) \widehat{\mathbf{W}}^{1/2} \mathbf{e}_\omega \\ &\quad - \mathbf{e}_\omega^\top \widehat{\mathbf{W}}^{1/2} \text{diag}(\boldsymbol{\omega}) \widehat{\mathbf{W}}^{1/2} (d_\omega \mathbf{H}(\lambda, \boldsymbol{\omega})) \mathbf{Y} \\ &\quad + \mathbf{e}_\omega^\top \widehat{\mathbf{W}}^{1/2} \text{diag}(d\boldsymbol{\omega}) \widehat{\mathbf{W}}^{1/2} \mathbf{e}_\omega. \end{aligned} \quad (\text{B.10})$$

Substituting equation (B.8) into (B.10), and using some simple algebra we obtain

$$d_{\omega} RSS(\lambda, \boldsymbol{\omega}) = \mathbf{e}_{\omega}^{\top} (\mathbf{I} - 2\mathbf{H}(\lambda, \boldsymbol{\omega}))^{\top} \widehat{\mathbf{W}} \text{diag}(\mathbf{e}_{\omega}) d\boldsymbol{\omega}. \quad (\text{B.11})$$

Therefore, using Equations (B.9) and (B.10) the first differential of $V(\lambda, \boldsymbol{\omega})$ with respect to $\boldsymbol{\omega}$ given in Equation (B.7) can now be written as

$$d_{\omega} V(\lambda, \boldsymbol{\omega}) = \frac{1/n}{(1 - \text{edf}(\lambda, \boldsymbol{\omega})/n)^3} \left\{ (1 - \text{edf}(\lambda, \boldsymbol{\omega})/n) \mathbf{e}_{\omega}^{\top} (\mathbf{I} - 2\mathbf{H}(\lambda, \boldsymbol{\omega}))^{\top} \widehat{\mathbf{W}} \text{diag}(\mathbf{e}_{\omega}) + \frac{2}{n} RSS(\lambda, \boldsymbol{\omega}) \mathbf{1}^{\top} \text{dg}((\mathbf{I} - \mathbf{H}(\lambda, \boldsymbol{\omega})) \mathbf{B} \mathbf{S}^{-1} \mathbf{B}^{\top} \widehat{\mathbf{W}}) \right\} d\boldsymbol{\omega}. \quad (\text{B.12})$$

Applying the first identification theorem by Magnus and Neudecker (1999) and evaluating (B.12) at $\lambda = \widehat{\lambda}$ and $\boldsymbol{\omega} = \boldsymbol{\omega}_0$ leads to $\partial V(\widehat{\lambda}, \boldsymbol{\omega}_0)/\partial \boldsymbol{\omega}$ given in Equation (20) from Subsection 3.2.3

Using matrix differentiation, we have

$$d_{\lambda} \mathbf{H}(\lambda, \boldsymbol{\omega}) = -\mathbf{G} \widehat{\mathbf{W}} \text{diag}(\boldsymbol{\omega}) d\lambda, \quad (\text{B.13})$$

where $\mathbf{G} = \mathbf{B} \mathbf{S}^{-1} \mathbf{K}^{\top} \mathbf{K} \mathbf{S}^{-1} \mathbf{B}^{\top}$. Thus,

$$d_{\lambda} (1 - \text{edf}(\lambda, \boldsymbol{\omega})/n)^{-3} = -\frac{3}{n} (1 - \text{edf}(\lambda, \boldsymbol{\omega})/n)^{-4} \text{tr}(\mathbf{G} \widehat{\mathbf{W}} \text{diag}(\boldsymbol{\omega})) d\lambda, \quad (\text{B.14})$$

and

$$d_{\lambda} \mathbf{e}_{\omega} = -(\mathbf{e}_{\omega}^{\top} \mathbf{H}(\lambda, \boldsymbol{\omega})) \mathbf{Y} = \mathbf{G} \widehat{\mathbf{W}} \text{diag}(\boldsymbol{\omega}) \mathbf{Y} d\lambda. \quad (\text{B.15})$$

Therefore, taking the differential of $RSS(\lambda, \boldsymbol{\omega})$ with respect to λ we obtain

$$d_{\lambda} RSS(\lambda, \boldsymbol{\omega}) = 2 \mathbf{Y}^{\top} \text{diag}(\boldsymbol{\omega}) \widehat{\mathbf{W}} \mathbf{G} \widehat{\mathbf{W}} \text{diag}(\boldsymbol{\omega}) \mathbf{e}_{\omega} d\lambda \quad (\text{B.16})$$

Using Equations (B.12)-(B.16), and after some algebra yields to the differential of $d_{\omega} V(\lambda, \boldsymbol{\omega})$ with respect to λ which is given by

$$\begin{aligned} d_{\lambda}^2 V(\lambda, \boldsymbol{\omega}) = & -\frac{3}{1 - \text{edf}(\lambda, \boldsymbol{\omega})/n} \text{tr}(\mathbf{G} \widehat{\mathbf{W}} \text{diag}(\boldsymbol{\omega})) d\lambda d_{\omega} V(\lambda, \boldsymbol{\omega}) \\ & + \frac{1/n}{(1 - \text{edf}(\lambda, \boldsymbol{\omega})/n)^3} \left\{ \frac{1}{n} \text{tr}(\mathbf{G} \widehat{\mathbf{W}} \text{diag}(\boldsymbol{\omega})) \mathbf{e}_{\omega}^{\top} (\mathbf{I} - 2\mathbf{H}(\lambda, \boldsymbol{\omega}))^{\top} \widehat{\mathbf{W}} \text{diag}(\mathbf{e}_{\omega}) \right. \\ & + (1 - \text{edf}(\lambda, \boldsymbol{\omega})/n) [\mathbf{Y}^{\top} \text{diag}(\boldsymbol{\omega}) \widehat{\mathbf{W}} \mathbf{G} (\mathbf{I} - 2\mathbf{H}(\lambda, \boldsymbol{\omega}))^{\top} \widehat{\mathbf{W}} \text{diag}(\mathbf{e}_{\omega}) \\ & + 2 \mathbf{e}_{\omega}^{\top} \text{diag}(\boldsymbol{\omega}) \widehat{\mathbf{W}} \mathbf{G} \widehat{\mathbf{W}} \text{diag}(\mathbf{e}_{\omega}) + \mathbf{e}_{\omega}^{\top} (\mathbf{I} - 2\mathbf{H}(\lambda, \boldsymbol{\omega}))^{\top} \widehat{\mathbf{W}} \text{diag}(\mathbf{G} \widehat{\mathbf{W}} \text{diag}(\boldsymbol{\omega}) \mathbf{Y})] \\ & + \frac{4}{n} \mathbf{Y}^{\top} \text{diag}(\boldsymbol{\omega}) \widehat{\mathbf{W}} \mathbf{G} \widehat{\mathbf{W}} \text{diag}(\boldsymbol{\omega}) \mathbf{e}_{\omega} \mathbf{1}^{\top} \text{dg}((\mathbf{I} - \mathbf{H}(\lambda, \boldsymbol{\omega})) \mathbf{B} \mathbf{S}^{-1} \mathbf{B}^{\top} \widehat{\mathbf{W}}) \\ & + \frac{2}{n} RSS(\lambda, \boldsymbol{\omega}) \mathbf{1}_n^{\top} \text{dg}(\mathbf{G} \widehat{\mathbf{W}} \text{diag}(\boldsymbol{\omega}) \mathbf{B} \mathbf{S}^{-1} \mathbf{B}^{\top} \widehat{\mathbf{W}} \\ & \left. - (\mathbf{I} - \mathbf{H}(\lambda, \boldsymbol{\omega})) \mathbf{G} \widehat{\mathbf{W}} \right\} d\lambda d\boldsymbol{\omega}, \quad (\text{B.17}) \end{aligned}$$

where the differential $d_{\omega} V(\lambda, \boldsymbol{\omega})$ is given in Equation (B.12) and $\text{dg}(\mathbf{Z}) = \text{diag}(z_{11}, \dots, z_{nn}) = \mathbf{I} \odot \mathbf{Z}$ for $\mathbf{Z} = (z_{ij})$ a square matrix of order $n \times n$ with \odot being the hadamard product. Thus, evaluating (B.17) at $\lambda = \widehat{\lambda}$ and $\boldsymbol{\omega} = \boldsymbol{\omega}_0$ and applying the second identification theorem by Magnus and Neudecker (1999) we obtain $\partial^2 V(\lambda, \boldsymbol{\omega}_0)/\partial \lambda \partial \boldsymbol{\omega}^{\top}|_{\lambda=\widehat{\lambda}}$ and the proof is complete. \square

Proof of Proposition 5. Under the response perturbation scheme, $\mathbf{Y}(\boldsymbol{\omega}) = \mathbf{Y} + \boldsymbol{\omega}$, where $\boldsymbol{\omega} = (\omega_1, \dots, \omega_n)^\top$ with $\boldsymbol{\omega}_0 = \mathbf{0}$, the perturbed WGCV criterion assumes the form

$$V(\lambda, \boldsymbol{\omega}) = \frac{RSS_{\widehat{\mathbf{W}}}(\lambda, \boldsymbol{\omega})/n}{\{\text{tr}(\mathbf{I} - \mathbf{H}_{\widehat{\mathbf{W}}}(\lambda))/n\}^2},$$

where

$$RSS_{\widehat{\mathbf{W}}}(\lambda, \boldsymbol{\omega}) = \mathbf{Y}^\top(\boldsymbol{\omega})(\mathbf{I} - \mathbf{H}_{\widehat{\mathbf{W}}}(\lambda))^\top \widehat{\mathbf{W}}(\mathbf{I} - \mathbf{H}_{\widehat{\mathbf{W}}}(\lambda))\mathbf{Y}(\boldsymbol{\omega}),$$

and

$$\mathbf{H}_{\widehat{\mathbf{W}}}(\lambda) = \mathbf{B}(\mathbf{B}^\top \widehat{\mathbf{W}}\mathbf{B} + \lambda \mathbf{K}^\top \mathbf{K})^{-1} \mathbf{B}^\top \widehat{\mathbf{W}}.$$

In order to simplify the notation we will write $RSS_{\widehat{\mathbf{W}}}(\lambda, \boldsymbol{\omega}) = RSS(\lambda, \boldsymbol{\omega})$ and $\mathbf{H}_{\widehat{\mathbf{W}}}(\lambda) = \mathbf{H}(\lambda)$. The first differential of $V(\lambda, \boldsymbol{\omega})$ with respect to $\boldsymbol{\omega}$ assumes the form

$$d_{\boldsymbol{\omega}} V(\lambda, \boldsymbol{\omega}) = \frac{1}{n} \frac{d_{\boldsymbol{\omega}} RSS(\lambda, \boldsymbol{\omega})}{\{1 - \text{tr}(\mathbf{H}(\lambda))/n\}^2},$$

where

$$d_{\boldsymbol{\omega}} RSS(\lambda, \boldsymbol{\omega}) = 2 \mathbf{Y}^\top(\boldsymbol{\omega})(\mathbf{I} - \mathbf{H}(\lambda))^\top \widehat{\mathbf{W}}(\mathbf{I} - \mathbf{H}(\lambda)) d\boldsymbol{\omega}.$$

Let $\text{edf}(\lambda) = \text{tr} \mathbf{H}(\lambda)$ be the effective number of parameters, thus taking the differential of $d_{\boldsymbol{\omega}} V(\lambda, \boldsymbol{\omega})$ with respect to λ yields

$$\begin{aligned} d_{\lambda}^2 V(\lambda, \boldsymbol{\omega}) &= \frac{2}{n} d_{\lambda} (1 - \text{edf}(\lambda, \boldsymbol{\omega})/n)^{-2} \mathbf{Y}^\top(\boldsymbol{\omega})(\mathbf{I} - \mathbf{H}(\lambda))^\top \widehat{\mathbf{W}}(\mathbf{I} - \mathbf{H}(\lambda)) d\boldsymbol{\omega} \\ &\quad + \frac{2/n}{(1 - \text{edf}(\lambda, \boldsymbol{\omega})/n)^2} \mathbf{Y}^\top(\boldsymbol{\omega}) d_{\lambda} [(\mathbf{I} - \mathbf{H}(\lambda))^\top \widehat{\mathbf{W}}(\mathbf{I} - \mathbf{H}(\lambda))] d\boldsymbol{\omega}, \end{aligned} \tag{B.18}$$

Let $\mathbf{S} = \mathbf{B}^\top \widehat{\mathbf{W}}\mathbf{B} + \lambda \mathbf{K}^\top \mathbf{K}$, and using that

$$d_{\lambda} \mathbf{H}(\lambda) = -\mathbf{B}\mathbf{S}^{-1}(d_{\lambda} \mathbf{S})\mathbf{S}^{-1}\mathbf{B}^\top \widehat{\mathbf{W}} = -\mathbf{G}\widehat{\mathbf{W}} d\lambda,$$

where $\mathbf{G} = \mathbf{B}\mathbf{S}^{-1}\mathbf{K}^\top \mathbf{K}\mathbf{S}^{-1}\mathbf{B}^\top$, leads to

$$\begin{aligned} d_{\lambda}(\mathbf{I} - \mathbf{H}(\lambda))^\top \widehat{\mathbf{W}}(\mathbf{I} - \mathbf{H}(\lambda)) &= -(d_{\lambda} \mathbf{H}(\lambda))^\top \widehat{\mathbf{W}}(\mathbf{I} - \mathbf{H}(\lambda)) - (\mathbf{I} - \mathbf{H}(\lambda))^\top \widehat{\mathbf{W}} d_{\lambda} \mathbf{H}(\lambda) \\ &= (\widehat{\mathbf{W}}\mathbf{G}\widehat{\mathbf{W}}(\mathbf{I} - \mathbf{H}(\lambda)) + (\mathbf{I} - \mathbf{H}(\lambda))^\top \widehat{\mathbf{W}}\mathbf{G}\widehat{\mathbf{W}}) d\lambda, \end{aligned} \tag{B.19}$$

therefore

$$d_{\lambda} (1 - \text{edf}(\lambda)/n)^{-2} = \frac{2/n}{(1 - \text{edf}(\lambda)/n)^3} \text{tr}(d_{\lambda} \mathbf{H}(\lambda)) = -\frac{2/n}{(1 - \text{edf}(\lambda)/n)^3} \text{tr}(\mathbf{G}\widehat{\mathbf{W}}) d\lambda. \tag{B.20}$$

Substituting equations (B.19) and (B.20) into (B.18) yields

$$\begin{aligned} d_{\lambda}^2 V(\lambda, \boldsymbol{\omega}) &= \frac{2/n}{(1 - \text{edf}(\lambda)/n)^2} \left\{ \mathbf{Y}^\top(\boldsymbol{\omega})(\widehat{\mathbf{W}}\mathbf{G}\widehat{\mathbf{W}}(\mathbf{I} - \mathbf{H}(\lambda)) + (\mathbf{I} - \mathbf{H}(\lambda))^\top \widehat{\mathbf{W}}\mathbf{G}\widehat{\mathbf{W}}) \right. \\ &\quad \left. - \frac{2/n}{(1 - \text{edf}(\lambda)/n)} \text{tr}(\mathbf{G}\widehat{\mathbf{W}}) \mathbf{Y}^\top(\boldsymbol{\omega})(\mathbf{I} - \mathbf{H}(\lambda))^\top \widehat{\mathbf{W}}(\mathbf{I} - \mathbf{H}(\lambda)) \right\} d\lambda d\boldsymbol{\omega}. \end{aligned} \tag{B.21}$$

Evaluating (B.21) at $\lambda = \widehat{\lambda}$ and $\boldsymbol{\omega} = \boldsymbol{\omega}_0$ and applying the second identification theorem by Magnus and Neudecker (1999), we obtain $\partial^2 V(\lambda, \boldsymbol{\omega}_0)/\partial \lambda \partial \boldsymbol{\omega}^\top|_{\lambda=\widehat{\lambda}}$ and the Proposition 5 is verified. \square

APPENDIX C. REAL DATA EXAMPLE

We consider the balloon dataset reported by [Davies and Gather \(1993\)](#) and previously analyzed by [Kovac and Silverman \(2000\)](#), [Lee and Oh \(2007\)](#) and [Tharmaratnam et al. \(2010\)](#). The data are radiation measurements from the sun, taken during the flight of a weather balloon. Due to the rotation of the balloon, or for some other reasons, a large amount of outliers were introduced because the measuring instrument was occasionally blocked from the sun. The dataset is included in the R package `ftnonpar` available from [CRAN](#) repository. The sample size equals 4984. Following the model settings from [Tharmaratnam et al. \(2010\)](#), we applied P-splines using B-splines of third-degree and second-order of penalty, with 35 knots spread equally. The model was fitted using the routine `heavyPS` from the R package `heavy`, also available at [CRAN](#). All the computations were done on an IBM x3650 M4 Server with 2 Intel Xeon E5-2670 processors and 128 GB of RAM Quadro 6000.

TABLE 1. Estimation summary for the balloon dataset under three fitted models.

Model	$\hat{\nu}$	$\hat{\lambda}$	$\ell_\lambda(\hat{\theta}; \mathbf{Y}_{\text{obs}})$	edf	WGCV	iterations	time (sec.)
Normal	—	0.3284	2470.705	27.1735	0.0220	3	0.060
Slash	0.4306	0.0081	7503.993	30.2584	0.0001	101	3.529
Student- <i>t</i>	0.9712	0.0230	7513.882	30.3645	0.0003	84	3.799

As expected, the estimation procedure based on heavy-tailed distributions produces a curve estimate that is not affected by the outlying observations (see [Figure 1](#) and [Table 1](#)). On the other hand, the curve estimate obtained under the normality assumption suffers from the presence of the outliers. That is, the whole estimated curve was pulled downwardly away from the majority of the observations. Another interesting feature of the proposed estimation procedure is that the implementation is fairly simple and ensure a reasonable computational speed (Compare with the numerical experiments reported by [Staudenmayer et al., 2009](#)).

We conducted an additional experiment, considering a sequence of values for the number of knots. [Figures 2 to 7](#) show the fitted curves for the balloon dataset considering 15, 20, 25, 30, 35 and 40 knots. [Table 2](#) presents the estimation results for this dataset under normal, slash and Student-*t* errors. It is interesting to note that for any number of knots the fitted curve under gaussian errors is strongly affected by the outlying observations, however the fitted curve obtained using heavy-tailed distributions do not suffer from this phenomenon. An unexpected finding is that the number of knots have an extremely large influence on the fitted curve under the normality assumption, but this influence is reduced when we consider distributions with heavier tails than the normal ones

APPENDIX D. SIMULATION STUDY

A small simulation study was conducted to illustrate the practical performance of the proposed estimation procedure under a severe contamination scheme. The design of the simulation study is based on the used by [Tharmaratnam et al. \(2010\)](#) (see also [Cantoni and Ronchetti, 2001](#); [Lee and Oh, 2007](#)). We have used the following test function

$$Y_i = \sin(2\pi(1 - x_i)^2) + \epsilon_i, \quad i = 1, \dots, n.$$

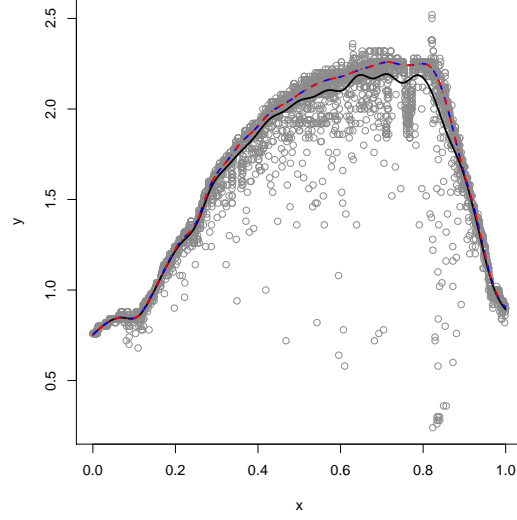


FIGURE 1. Fitted curve for the balloon dataset under three distributional assumptions: (—) normal, (---) Student- t and (-·-) slash models.

TABLE 2. Summary for the balloon dataset under three fitted models, considering different number of knots.

Number of knots	normal		slash			Student- t		
	$\hat{\lambda}$	$\ell_{\lambda}(\hat{\theta}; \mathbf{Y}_{\text{obs}})$	$\hat{\lambda}$	$\hat{\nu}$	$\ell_{\lambda}(\hat{\theta}; \mathbf{Y}_{\text{obs}})$	$\hat{\lambda}$	$\hat{\nu}$	$\ell_{\lambda}(\hat{\theta}; \mathbf{Y}_{\text{obs}})$
15	0.010	2408.03	0.000	0.499	6721.26	0.000	1.119	6742.48
20	0.011	2419.76	0.000	0.481	6945.77	0.001	1.076	6965.06
25	0.018	2439.09	0.000	0.432	7349.22	0.000	0.966	7375.85
30	0.657	2427.79	0.000	0.437	7382.80	0.001	0.976	7406.13
35	0.328	2470.71	0.008	0.431	7504.00	0.023	0.971	7513.88
40	0.089	2515.78	0.001	0.434	7554.08	0.001	0.967	7573.12

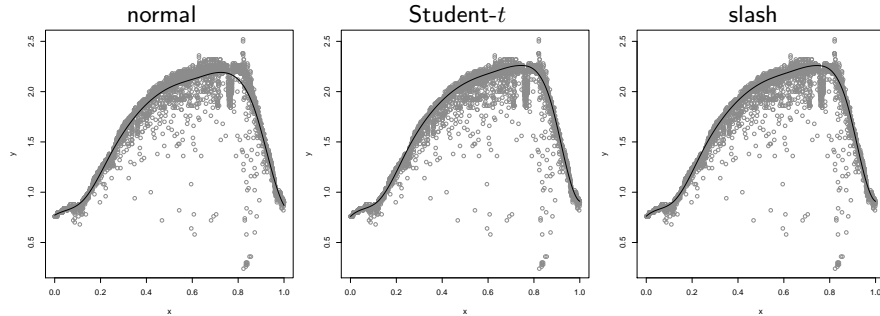


FIGURE 2. Fitted curve for the normal, Student- t and slash models using 15 knots.

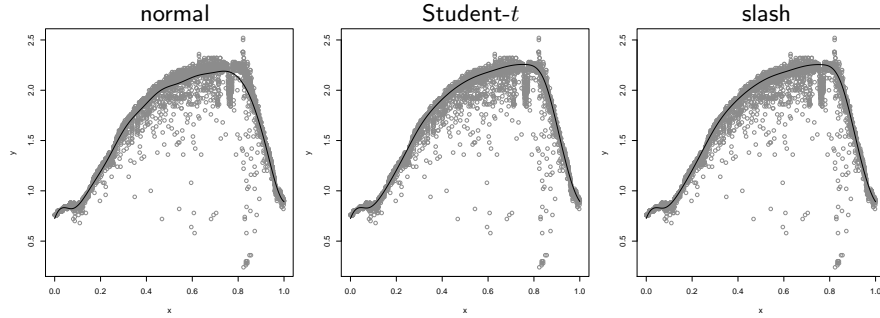


FIGURE 3. Fitted curve for the normal, Student- t and slash models using 20 knots.

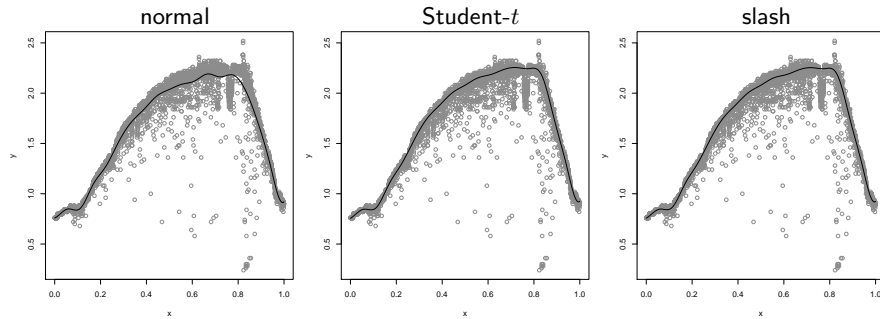


FIGURE 4. Fitted curve for the normal, Student- t and slash models using 25 knots.

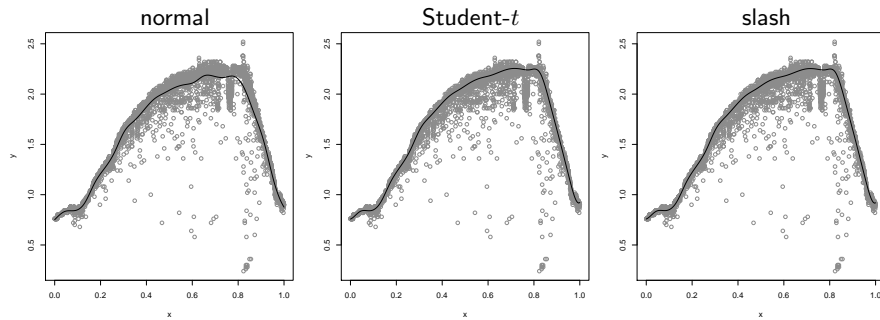


FIGURE 5. Fitted curve for the normal, Student- t and slash models using 30 knots.

We generated $M = 1000$ datasets of a sample size $n = 100$ from the model above with design points x_1, \dots, x_n generated independently from the uniform distribution $\mathcal{U}(-1, 1)$, the random disturbances $\{\epsilon_i\}$ were generated from the contaminated normal distribution

$$(1 - \delta)\mathcal{N}(0, 0.7^2) + \delta\mathcal{N}(20, 2^2),$$

for $\delta = 0\%, 5\%, 10\%, 25\%$ and 40% of outliers. We applied P-splines considering B-splines of third-degree and second-order of penalty, with 25 knots spread equally according to the quantiles of the data. The normal, slash and Student- t distributions were considered, fixing the degrees of freedom for the slash and Student- t

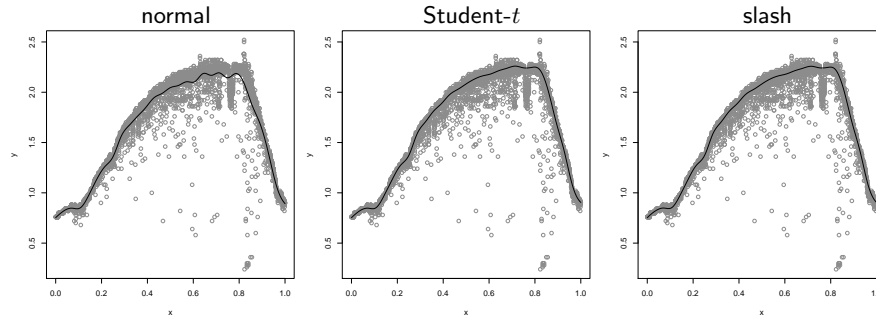


FIGURE 6. Fitted curve for the normal, Student- t and slash models using 35 knots.

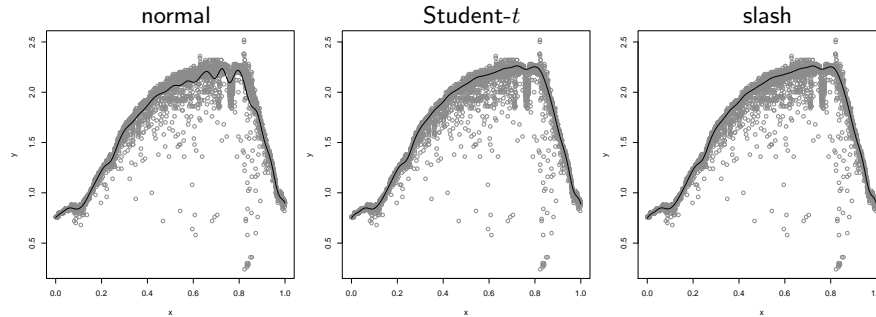


FIGURE 7. Fitted curve for the normal, Student- t and slash models using 40 knots.

at $1/2, 1$ and 2 , and $1, 2$, and 4 , respectively. The smoothing parameter is chosen by minimizing the weighted generalized cross-validation (WGCV) criterion defined in Section 2.2. Figure 8 displays two typical datasets together with the fitted curve from the gaussian and Student- t models (fitted curve for the slash model is similar to the obtained under Student- t errors and is not shown here), unlike the P-splines considering the assumption of normality the penalized spline estimator under Student- t errors with 1 degrees of freedom remains close to the true regression function (solid line) for both situations, without and in presence of outliers.

To assess the performance of the proposed procedure we compute the mean squared error (MSE) for each simulated dataset j ($j = 1, \dots, M$), as

$$MSE_j = \frac{1}{n} \sum_{i=1}^n (g(x_i) - \hat{g}_j(x_i))^2, \quad j = 1, \dots, M,$$

with $M = 1000$. We ran our simulation experiment on an IBM x3650 M4 Server with 2 Intel Xeon E5-2670 processors and 128 GB of RAM Quadro 6000. The total computation time was 16 hours, 56 minutes and 5.13 seconds.

Figures 9 to 11 present the boxplots of MSE in the logarithmic scale for several contamination percentages considering the P-splines estimation under normal, Student- t and slash distributions. From these plots it is clear that when we are in presence of outliers, the proposed procedure produces better estimates. In fact our findings are very similar to those reported by Tharmaratnam et al. (2010). However, it should be noted that the penalized spline estimator under heavy-tailed

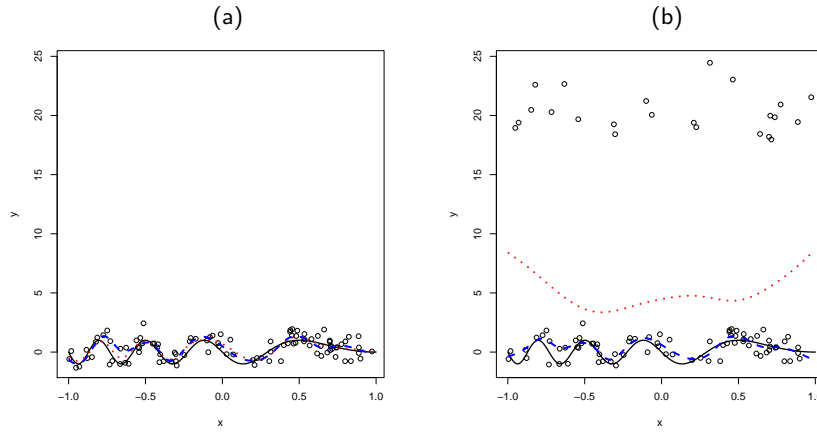


FIGURE 8. Fitted curve for the [Cantoni and Ronchetti \(2001\)](#) test function (a) without outliers and (b) with 25% of outliers from $\mathcal{N}(20, 2^2)$. True function $\sin(2\pi(1-x)^2)$ (solid line), fitted curve using normal errors (penalized LS-estimation) (dotted) and assuming Student- t model with $\nu = 1$ (dashed).

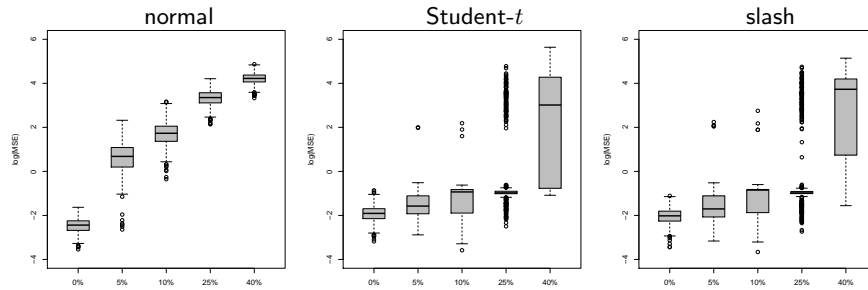


FIGURE 9. Boxplot of MSE using normal, Student- t ($\nu = 1$) and slash ($\nu = \frac{1}{2}$) models considering several contamination levels for the [Cantoni and Ronchetti \(2001\)](#) test function.

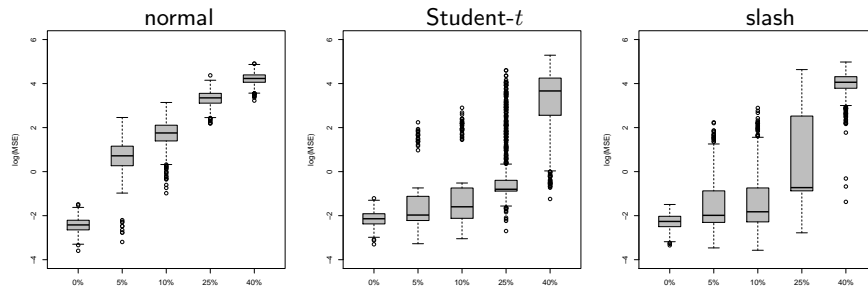


FIGURE 10. Boxplot of MSE using normal, Student- t ($\nu = 2$) and slash ($\nu = 1$) models considering several contamination levels for the [Cantoni and Ronchetti \(2001\)](#) test function.

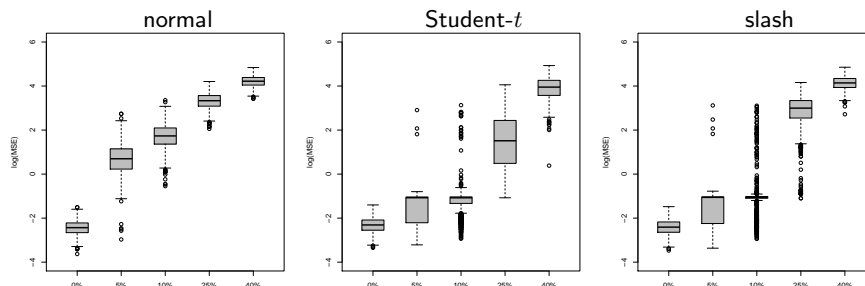


FIGURE 11. Boxplot of MSE using normal, Student- t ($\nu = 4$) and slash ($\nu = 2$) models considering several contamination levels for the [Cantoni and Ronchetti \(2001\)](#) test function.

distributions remains stable as the proportion of contamination increases only for small (and fixed) degrees of freedom. In our simulation experiment the penalized spline estimator's MSEs grows rapidly after 10% of outliers for degrees of freedom as small as 4 or 2, for the Student- t or slash distributions, respectively. In fact, as is discussed by [Lange et al. \(1989\)](#) modeling with distributions with heavier tails than the normal one is not the panacea for all robustness problems. In particular we can expect difficulties in P-splines under scale mixtures of normal distributions for data with extreme outliers.

REFERENCES

- BRENT, R.P. (1973). *Algorithms for Minimization Without Derivatives*. Dover, New York.
- CANTONI, E., RONCHETTI, E. (2001). Resistant selection of the smoothing parameter for smoothing splines. *Statistics and Computing* **11**, 141-146.
- DAVIES, L., GATHER, U. (1993). The identification of multiple outliers (with discussion). *Journal of the American Statistical Association* **88**, 782-801.
- DONGARRA, J.J., BUNCH, J.R., MOLER, C.B., STEWART, G.W. (1979). *Linpac Users Guide*. SIAM Publications, Philadelphia.
- KOVAC, A., SILVERMAN, B.W. (2000). Extending the scope of wavelet regression methods by coefficient-dependent thresholding. *Journal of the American Statistical Association* **95**, 172-183.
- LANGE, K., LITTLE, R.J.A., TAYLOR, J.M.G. (1989). Robust statistical modeling using the t distribution. *Journal of the American Statistical Association* **84**, 881-896.
- LAWSON, C.L., HANSON, R.J., KINCAID, D., KROGH, F.T. (1979). Basic linear algebra subprograms for FORTRAN usage. *ACM Transactions on Mathematical Software* **5**, 308-323.
- LEE, T.C.M., OH, H.S. (2007) Robust penalized regression spline fitting with application to additive mixed modeling. *Computational Statistics* **22**, 159-171.
- MAGNUS, J.R., NEUDECKER, H. (1999). *Matrix Differential Calculus with Applications in Statistics and Econometrics*. Wiley, Chichester.
- R CORE TEAM (2014). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. <http://www.R-project.org>
- RUPPERT, D. (2002) Selecting the number of knots for penalized splines. *Journal of Computational and Graphical Statistics* **11**, 735-757.
- STAUDENMAYER, J., LAKE, E.E., WAND, M.P. (2009). Robustness for general design mixed models using the t -distribution. *Statistical Modelling* **9**, 235-255.

THARMARATNAM, K., CLAESKENS, G., CROUX, C., SALIBIÁN-BARRERA, M. (2010). S-estimation for penalized regression splines. *Journal of Computational and Graphical Statistics* **19**, 609-625.

DEPARTAMENTO DE MATEMÁTICA, UNIVERSIDAD TÉCNICA FEDERICO SANTA MARÍA, CHILE
Current address: Departamento de Matemática, Universidad Técnica Federico Santa María,
Av. España 1680, Casilla 110-V, Valparaíso, Chile
E-mail address: felipe.osorios@usm.cl