



On the estimation of variance parameters in non-standard generalised linear mixed models: application to penalised smoothing

María Xosé Rodríguez-Álvarez^{1,2} · Maria Durban³ · Dae-Jin Lee¹ · Paul H. C. Eilers⁴

Received: 16 March 2018 / Accepted: 5 June 2018
© Springer Science+Business Media, LLC, part of Springer Nature 2018

Abstract

We present a novel method for the estimation of variance parameters in generalised linear mixed models. The method has its roots in Harville (J Am Stat Assoc 72(358):320–338, 1977)'s work, but it is able to deal with models that have a precision matrix for the random effect vector that is linear in the inverse of the variance parameters (i.e., the precision parameters). We call the method SOP (separation of overlapping precision matrices). SOP is based on applying the method of successive approximations to easy-to-compute estimate updates of the variance parameters. These estimate updates have an appealing form: they are the ratio of a (weighted) sum of squares to a quantity related to effective degrees of freedom. We provide the sufficient and necessary conditions for these estimates to be strictly positive. An important application field of SOP is penalised regression estimation of models where multiple quadratic penalties act on the same regression coefficients. We discuss in detail two of those models: penalised splines for locally adaptive smoothness and for hierarchical curve data. Several data examples in these settings are presented.

Keywords Generalised linear mixed models · Generalised additive models · Variance parameters · Smoothing parameters · REML · Effective degrees of freedom

1 Introduction

The estimation of variance parameters is a statistical problem that has received extensive attention for more than 50 years. It originated with the ANOVA methodology proposed by Fisher in the 1920's, where estimates were obtained equating mean squared error to its expected value. However, the results yielded by this method were not optimal in some situations, for example, in the case of unbalanced data. Later on, Crump (1951) applied maximum likelihood (ML) under the assumption of normally distributed errors and random effects. But it was not until the 1970's when

the estimation of variance parameters based on ML methods gained interest. The method of *Restricted Maximum Likelihood* (REML) (Patterson and Thompson 1971) gave a solution to the problem of biased estimators of the variance parameters. However, one of the main obstacles to the use of this technique, at the time, was the fact that the calculation of ML/REML estimates requires the numerical solution of a nonlinear problem. Patterson and Thompson (1971) proposed an iterative solution using the Fisher Scoring algorithm, but it was Harville (1977) who proposed the first numerical algorithm to compute REML estimates of the variance parameters. His proposal is the inspiration of our work.

Along the years, several computational approaches have appeared with the aim of improving the computational burden of solving the score equations for the variance parameters: Smith (1990) proposed the use of the EM algorithm, Graser et al. (1987) suggested the use of the simplex algorithm to obtain the estimates directly from the likelihood, and Gilmour et al. (1995) developed a method based on the use of an average information matrix.

In the context of generalised linear mixed models (GLMMs), estimation based on iterative re-weighted REML has been proposed independently by a number of authors

✉ María Xosé Rodríguez-Álvarez
mxrodriguez@bcamath.org

¹ BCAM - Basque Center for Applied Mathematics, Alameda de Mazarredo, 14, 48009 Bilbao, Basque Country, Spain

² IKERBASQUE, Basque Foundation for Science, Bilbao, Spain

³ Department of Statistics and Econometrics, Universidad Carlos III de Madrid, Leganés, Spain

⁴ Erasmus University Medical Centre, Rotterdam, The Netherlands

(e.g., Schall 1991; Engel and Keen 1994), as an extension of the iterative re-weighted least squares algorithm for Generalised Linear Models (GLM, McCullagh and Nelder 1989). Breslow and Clayton (1993) proposed a general method based on penalised quasi-likelihood (PQL) for the estimation of the fixed and random effects, and pseudo-likelihood for the variance parameters. As noted by Engel and Buist (1996), the estimation procedures discussed in all these papers are equivalent, although motivated from quite different starting points.

The majority of the methods mentioned above impose a strong restriction on the vector of random effects: its variance-covariance matrix has to be linear in the variance parameters. The results we present in this paper relax that assumption to the case in which the linearity in the parameters is necessary on the precision matrix and not on the variance-covariance matrix. Our contribution is motivated by the need to estimate smoothing parameters in the context of penalised regression models with non-standard quadratic penalties.

Penalised spline regression (P-splines, Eilers and Marx 1996) has become a popular method for estimating models in which the mean response (or linear predictor in the non-Gaussian case) is a smooth unknown function of one or more covariates. The method is based on the representation of the smooth component in terms of basis functions and the estimation of the parameters by modifying the likelihood with a quadratic penalty on the coefficients. The size of the penalty is controlled by the so-called *smoothing parameter*. The connection between penalised smoothing and linear mixed models was first established a long time ago (Green 1987), and it has become of common use in the last 15 years (Currie and Durban 2002; Currie et al. 2006; Lee 2010; Wand 2003). The key point of the equivalence is that the smoothing parameter becomes the ratio between two variance parameters. Therefore, the methods mentioned above can be used to estimate directly the amount of smoothing needed in the model, instead of using methods based on the minimisation of some prediction error such as the Akaike information criterion (AIC), generalised cross-validation (GCV) or Mallows' C_p (see, e.g., Eilers and Marx 1996; Wood 2008). In Krivobokova (2009), the asymptotic properties of REML and Mallows' C_p smoothing parameter estimators are considered. The paper shows that, in the frequentist framework, the REML oracle smoothing parameter is asymptotically sub-optimal. However, the variance of the REML-based estimator is found to be smaller than that of the Mallows' C_p estimator. Also, Reiss and Ogden (2009) show that at finite sample sizes AIC and GCV are prone to undersmoothing and are more likely to develop multiple minima than REML (see also Wood 2017, Chapter 6). An additional argument in favour of REML is that the inclusion of random effects into the penalised regression model is straightforward.

In the smoothing context, standard methods based on REML can be applied when simple penalties are used, i.e.,

each regression coefficient is affected by a single penalty (by a single smoothing parameter). However, in some circumstances, the penalties present an overlapping structure, with the same coefficients being penalised simultaneously by several smoothing parameters. This includes important cases such as multidimensional penalised splines with anisotropic penalties or adaptive penalised splines. Estimation methods that can deal with this situation have been proposed in the smoothing literature (e.g., Wood 2011; Wood et al. 2016), but they have the drawback of being very computationally demanding, especially when the number of smoothing parameters is large.

This work addresses this problem and presents a fast method for estimating the variance parameters/smoothing parameters in generalised linear mixed models/generalised additive models. The method can be used whenever the precision matrix of the random component (or the penalty matrix of the P-spline model) is a linear combination defined over the inverse of the variance parameters (smoothing parameters). We obtain simple expressions for the estimates of the variance parameters that are ratios between a sum of squares and a quantity related to the notion of effective degrees of freedom in the smoothing context (Hastie and Tibshirani 1990). We show the sufficient and necessary conditions that guarantee the positiveness of these estimates and discuss several situations where these conditions can be easily verified. Particular cases of the method presented here have been introduced in Rodríguez-Álvarez et al. (2015b), which solved the problem in the case of anisotropic multidimensional smoothing, and in Rodríguez-Álvarez et al. (2015a), where results for adaptive P-splines were first discussed. More recently, Wood and Fasiolo (2017) extended the above-mentioned works to more general penalised spline models. The proposal discussed here presents two main advantages with respect to Wood and Fasiolo (2017)'s approach. First, the smoothing/variance parameter estimates described in Wood and Fasiolo (2017) rely on Moore–Penrose pseudoinverses of the penalty matrices, which, in our experience, may present numerical instabilities. Second, our proposal establishes an explicit connection between variance component estimates and effective degrees of freedom, which lacks in Wood and Fasiolo (2017). Effective degrees of freedom are key components in smoothing models. They help in summarising a model, as partial effective degrees of freedom are measures of model components' complexity with strong intuitive appeal (see, e.g., Rodríguez-Álvarez et al. 2018 for an example in the agricultural field).

The rest of this paper is organised as follows: Sect. 2 introduces the work by Harville (1977), which constitutes the foundation of the work presented here. Section 3 is the core of the paper: the new method, called SOP (separation of overlapping precision matrices), is presented; and the connection between SOP and the notion of effective degrees of

freedom is discussed. Section 4 describes several P-splines models whose estimation can be approached using SOP. We focus in this paper on adaptive P-splines and P-splines for hierarchical curve data. Illustrations with data examples are provided in Sect. 5. A discussion closes the paper. Some technical details have been added as appendices. The estimating algorithm is detailed there.

2 Estimation of variance parameters in generalised linear mixed models: Harville (1977)'s work and extensions

This section is our little tribute to Harville (1977)'s paper, which was the inspiration for this work. Harville (1977)'s paper deals with ML/REML approaches to variance parameters estimation in linear mixed models (LMM) for Gaussian data. Nonetheless, estimation of GLMM can be done by repeated use of LMM methodology on a *working*-dependent variable (see, e.g., Schall 1991; Engel and Keen 1994 where use is made of the results by Harville 1977). This is the approach we follow in this paper.

Let $\mathbf{y} = (y_1, \dots, y_n)^\top$ be a vector of n observations. A GLMM can be written as

$$g(\boldsymbol{\mu}) = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\alpha}, \quad \text{with } \boldsymbol{\alpha} \sim N(\mathbf{0}, \mathbf{G}), \quad (1)$$

where $\mu_i = \mathbb{E}(y_i | \boldsymbol{\alpha})$ and $g(\cdot)$ is the link function. The model assumes that, conditional on the random effects, the observations y_i are independently distributed with mean μ_i and variance $\text{Var}(y_i | \boldsymbol{\alpha}) = \phi v(\mu_i)$. Here, $v(\cdot)$ is a specified variance function, and ϕ is the dispersion parameter that may be known or unknown. In model (1), \mathbf{X} and \mathbf{Z} represent column-partitioned matrices, associated, respectively, with the fixed and random effects. We assume that \mathbf{X} has full rank, $\mathbf{Z} = [\mathbf{Z}_1, \dots, \mathbf{Z}_c]$, and $\boldsymbol{\alpha} = (\boldsymbol{\alpha}_1^\top, \dots, \boldsymbol{\alpha}_c^\top)^\top$. Each \mathbf{Z}_k corresponds to the design matrix of the k -th random component $\boldsymbol{\alpha}_k$, with $\boldsymbol{\alpha}_k$ being a $(q_k \times 1)$ vector ($k = 1, \dots, c$). We assume further that $\boldsymbol{\alpha}_k \sim N(\mathbf{0}, \mathbf{G}_k)$ and that

$$\mathbf{G} = \bigoplus_{k=1}^c \mathbf{G}_k = \bigoplus_{k=1}^c \sigma_k^2 \mathbf{I}_{q_k} = \text{diag}(\sigma_1^2 \mathbf{I}_{q_1}, \dots, \sigma_c^2 \mathbf{I}_{q_c}),$$

where \mathbf{I}_m is an identity matrix of order $m \times m$ and \bigoplus denotes the direct sum of matrices. Note that the variance-covariance matrix \mathbf{G} is linear in the variance parameters σ_k^2 .

As noted before, estimation of model (1) can be approached by iterative fitting of a LMM that involves a *working* dependent variable \mathbf{z} and a weight matrix \mathbf{W} (updated at each iteration). The specific form of \mathbf{z} and \mathbf{W} is given in "Appendix C". If ϕ and σ_k^2 ($k = 1, \dots, c$) are known, at each iteration, the updates for $\boldsymbol{\beta}$ and $\boldsymbol{\alpha}$ follow from the so-called Henderson equations (Henderson 1963)

$$\underbrace{\begin{bmatrix} \mathbf{X}^\top \mathbf{R}^{-1} \mathbf{X} & \mathbf{X}^\top \mathbf{R}^{-1} \mathbf{Z} \\ \mathbf{Z}^\top \mathbf{R}^{-1} \mathbf{X} & \mathbf{Z}^\top \mathbf{R}^{-1} \mathbf{Z} + \mathbf{G}^{-1} \end{bmatrix}}_{\mathbf{C}} \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\alpha}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}^\top \mathbf{R}^{-1} \mathbf{z} \\ \mathbf{Z}^\top \mathbf{R}^{-1} \mathbf{z} \end{bmatrix}, \quad (2)$$

which give rise to closed-form expressions

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{V}^{-1} \mathbf{z}, \\ \hat{\boldsymbol{\alpha}}_k &= \mathbf{G}_k \mathbf{Z}_k^\top \mathbf{P} \mathbf{z} \quad (k = 1, \dots, c), \end{aligned} \quad (3)$$

where $\mathbf{P} = \mathbf{V}^{-1} - \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{V}^{-1}$ with $\mathbf{V} = \mathbf{R} + \mathbf{Z} \mathbf{G} \mathbf{Z}^\top$ and $\mathbf{R} = \phi \mathbf{W}^{-1}$. The Henderson equations are of little use if the variances parameters ϕ and σ_k^2 ($k = 1, \dots, c$) are unknown. In his 1977 paper, Harville shows how to estimate them by REML by an elegant iterative algorithm. Let's first define

$$\mathbf{T} = (\mathbf{I} + \mathbf{Z}^\top \mathbf{S} \mathbf{Z} \mathbf{G})^{-1},$$

where $\mathbf{S} = \mathbf{R}^{-1} - \mathbf{R}^{-1} \mathbf{X} (\mathbf{X}^\top \mathbf{R}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{R}^{-1}$. We note that \mathbf{T} can be partitioned as follows

$$\mathbf{T} = \begin{bmatrix} \mathbf{T}_{11} & \mathbf{T}_{12} & \cdots & \mathbf{T}_{1c} \\ \mathbf{T}_{21} & \mathbf{T}_{22} & \cdots & \mathbf{T}_{2c} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{T}_{c1} & \mathbf{T}_{c2} & \cdots & \mathbf{T}_{cc} \end{bmatrix},$$

where \mathbf{T}_{ij} are matrices of order $q_i \times q_j$. In Harville (1977), the updated estimate of σ_k^2 ($k = 1, \dots, c$) is

$$\hat{\sigma}_k^2 = \frac{\hat{\boldsymbol{\alpha}}_k^{[t] \top} \hat{\boldsymbol{\alpha}}_k^{[t]}}{\text{ED}_k^{[t]}}, \quad (4)$$

where

$$\text{ED}_k^{[t]} = q_k - \text{trace}(\mathbf{T}_{kk}^{[t]}), \quad (5)$$

and the superscript $[t]$ denotes quantities evaluated at current estimates of the variance parameters. From the estimates of $\boldsymbol{\beta}$ and $\boldsymbol{\alpha}$ follow an estimate for \mathbf{z} : $\hat{\mathbf{z}} = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{Z}\hat{\boldsymbol{\alpha}}$. The residuals are $\mathbf{z} - \hat{\mathbf{z}}$. Harville uses

$$\hat{\phi} = \frac{\mathbf{z}^\top \mathbf{W} (\mathbf{z} - \hat{\mathbf{z}}^{[t]})}{n - \text{rank}(\mathbf{X})}, \quad (6)$$

to estimate the dispersion parameter (not always needed in GLMM). An alternative expression is (see, e.g., Engel 1990; Rodríguez-Álvarez et al. 2015b)

$$\hat{\phi} = \frac{(\mathbf{z} - \hat{\mathbf{z}}^{[t]})^\top \mathbf{W} (\mathbf{z} - \hat{\mathbf{z}}^{[t]})}{n - \text{rank}(\mathbf{X}) - \sum_{k=1}^c \text{ED}_k^{[t]}}. \quad (7)$$

Here $\text{rank}(\mathbf{X}) + \sum_{k=1}^c \text{ED}_k^{[t]}$ can be interpreted as the effective model dimension. At convergence, Eqs. (6) and (7) give identical numerical values.

2.1 Effective degrees of freedom in Harville's method

As noted by Harville (1977), the iterates derived from expression (4) have an intuitively appealing form. On each iteration, σ_k^2 is estimated by the ratio between the sum of squares of the estimates for α_k and a number between zero and q_k . We now show that the denominator in expression (4) can, in fact, be interpreted as effective degrees of freedom in smoothing sensu, i.e., as the trace of a “hat” matrix (Hastie and Tibshirani 1990).

First note that expression (3) reveals that $\mathbf{Z}_k \hat{\alpha}_k = \mathbf{Z}_k \mathbf{G}_k \mathbf{Z}_k^\top \mathbf{P} \mathbf{z}$. Thus, the “hat” matrix corresponding to the k -th random component α_k is

$$\mathbf{H}_k = \mathbf{Z}_k \mathbf{G}_k \mathbf{Z}_k^\top \mathbf{P},$$

i.e., $\mathbf{H}_k \mathbf{z} = \mathbf{Z}_k \hat{\alpha}_k$. We now show that $\text{trace}(\mathbf{H}_k) = \text{ED}_k$. It is easy to verify that

$$\mathbf{T} = (\mathbf{I} + \mathbf{Z}^\top \mathbf{S} \mathbf{Z} \mathbf{G})^{-1} = \mathbf{G}^{-1} (\mathbf{G}^{-1} + \mathbf{Z}^\top \mathbf{S} \mathbf{Z})^{-1}, \quad (8)$$

where $(\mathbf{G}^{-1} + \mathbf{Z} \mathbf{S} \mathbf{Z}^\top)^{-1}$ is that partition of the inverse of \mathbf{C} in (2) corresponding to the random vector α (Harville 1977; Johnson and Thompson 1995). Exploiting the block structure of \mathbf{Z} and \mathbf{G} , and making use of result (8) and (A4) in Johnson and Thompson (1995), we have that

$$\mathbf{Z}_k^\top \mathbf{P} \mathbf{Z}_k \mathbf{G}_k = \mathbf{I}_{q_k} - \mathbf{G}_k^{-1} \mathbf{C}_{kk}^* = \mathbf{I}_{q_k} - \mathbf{T}_{kk}, \quad (9)$$

where, to ease the notation, \mathbf{C}^* denotes the inverse of \mathbf{C} and \mathbf{C}_{kk}^* denotes that partition of \mathbf{C}^* corresponding to the k -th random component α_k . Thus,

$$\begin{aligned} \text{trace}(\mathbf{H}_k) &= \text{trace}(\mathbf{Z}_k \mathbf{G}_k \mathbf{Z}_k^\top \mathbf{P}) = \text{trace}(\mathbf{Z}_k^\top \mathbf{P} \mathbf{Z}_k \mathbf{G}_k) \\ &= \text{trace}(\mathbf{I}_{q_k} - \mathbf{T}_{kk}) = q_k - \text{trace}(\mathbf{T}_{kk}) \\ &= \text{ED}_k. \end{aligned}$$

3 Separation of overlapping precision matrices: the SOP method

In the previous section, we have discussed an estimating method for generalised linear mixed models where the variance–covariance matrix of the random component is linear in the variance parameters. However, more complex

structures of the variance–covariance matrix appear in practice. The present research was motivated by our work on penalised spline regression. In spite of that, the method to be discussed in this section is not confined to this area: it can be seen as a general estimating method for generalised linear mixed models with a precision matrix of a specific structure. As in Sect. 2, we consider the generalised linear mixed model

$$g(\mu) = \mathbf{X} \beta + \mathbf{Z} \alpha = \mathbf{X} \beta + \sum_{k=1}^c \mathbf{Z}_k \alpha_k, \quad (10)$$

with $\alpha_k \sim N(\mathbf{0}, \mathbf{G}_k)$, $\alpha \sim N(\mathbf{0}, \mathbf{G})$, and $\mathbf{G} = \bigoplus_{k=1}^c \mathbf{G}_k$. The main difference with respect to Sect. 2 is that we do not assume that $\mathbf{G}_k = \sigma_k^2 \mathbf{I}_{q_k}$, but we consider precision matrices of the form

$$\mathbf{G}_k^{-1} = \sum_{l=1}^{p_k} \sigma_{kl}^{-2} \mathbf{\Lambda}_{kl}, \quad (11)$$

where σ_{kl}^2 ($l = 1, \dots, p_k$ and $k = 1, \dots, c$) are the variance parameters and $\mathbf{\Lambda}_{kl}$ are known symmetric positive semi-definite matrices of dimension $q_k \times q_k$. Note that we do not require $\mathbf{\Lambda}_{kl}$ to be positive definite. The only requirement we need is that \mathbf{G}_k^{-1} ($k = 1, \dots, c$) are positive definite and so are \mathbf{G}^{-1} and its inverse, the variance–covariance matrix \mathbf{G} .

Expression (11) deserves some detailed discussion. Firstly, it is worth noting that we do not work with variance–covariance matrices, but with their inverses, the precision matrices. As said, the developments in this work have their origin on penalised spline methods. In Sect. 4, the need to work with precision matrices will become clear, or, in the terminology of penalised splines, with penalty matrices. Secondly, what constitutes the main contribution of this paper is that we assume that each random component α_k ($k = 1, \dots, c$) in model (10) may be “affected” (shrunk) by several variance parameters. A particular case would be when $p_k = 1 \forall k$, in which case we are in the situation discussed in Sect. 2.

For the sake of simplicity, in some cases we will rewrite the precision matrix \mathbf{G}^{-1} as follows

$$\mathbf{G}^{-1} = \sum_{l=1}^p \sigma_l^{-2} \tilde{\mathbf{\Lambda}}_l, \quad (12)$$

where $p = \sum_{k=1}^c p_k$. By a slight abuse of notation, let $\mathbf{\Lambda}_l$ denote the matrices involved in expression (11). The matrix $\tilde{\mathbf{\Lambda}}_l$ is $\mathbf{\Lambda}_l$ padded out with zeroes. Some specific examples will be presented in Sect. 4 below. Expression (12) makes it clear that the present work deals with the situation of generalised linear mixed models with a precision matrix for the random component that is linear in the precision parameters σ_l^{-2} .

The next section presents the proposed estimation method that we call SOP.

3.1 The SOP method

Regardless of the structure of \mathbf{G} , estimation of $\boldsymbol{\beta}$ and $\boldsymbol{\alpha}$, and, when necessary, the dispersion parameter ϕ , does not pose a problem and can be done as discussed in Sect. 2 above (see also “Appendix C” for a detailed description of the estimating algorithm). Recall that estimates for $\boldsymbol{\alpha}_k$ are obtained as $\hat{\boldsymbol{\alpha}}_k = \mathbf{G}_k \mathbf{Z}_k^\top \mathbf{P} \mathbf{z}$. The hat matrix associated with the k -th random component $\boldsymbol{\alpha}_k$ is, once again, $\mathbf{H}_k = \mathbf{Z}_k \mathbf{G}_k \mathbf{Z}_k^\top \mathbf{P}$, and the effective degrees of freedom of this component are $\text{ED}_k = \text{trace}(\mathbf{H}_k)$. The variance parameters $\sigma_{k_l}^2$ ($l = 1, \dots, p_k$ and $k = 1, \dots, c$), however, cannot be estimated by means of Harville’s approach. This is a consequence of \mathbf{G} not being linear in the variance parameters.

The key of our approach is to work with \mathbf{G}^{-1} instead of with \mathbf{G} . Given that \mathbf{G}^{-1} is linear in the precision parameters $\sigma_{k_l}^{-2}$, the first-order partial derivatives of the (approximate) REML log-likelihood function can be explicitly obtained, as well as the REML-based estimates of the variance parameters. We state the result in the following Theorem, whose proof is given in “Appendix A”.

Theorem 1 Let $\mathbf{G} = \bigoplus_{k=1}^c \mathbf{G}_k$ be a symmetric positive definite matrix, with $\mathbf{G}_k^{-1} = \sum_{l=1}^{p_k} \sigma_{k_l}^{-2} \boldsymbol{\Lambda}_{k_l}$ symmetric positive definite and $\boldsymbol{\Lambda}_{k_l}$ known symmetric positive semi-definite. Then, the REML-based estimate updates of the variance parameters $\sigma_{k_l}^2$ ($l = 1, \dots, p_k$ and $k = 1, \dots, c$) are given by

$$\hat{\sigma}_{k_l}^2 = \frac{\hat{\boldsymbol{\alpha}}_k^{[t]\top} \boldsymbol{\Lambda}_{k_l} \hat{\boldsymbol{\alpha}}_k^{[t]}}{\text{ED}_{k_l}^{[t]}}, \quad (13)$$

where

$$\text{ED}_{k_l}^{[t]} = \text{trace} \left(\mathbf{Z}_k^\top \mathbf{P}^{[t]} \mathbf{Z}_k \mathbf{G}_k^{[t]} \frac{\boldsymbol{\Lambda}_{k_l}}{\hat{\sigma}_{k_l}^{2[t]}} \mathbf{G}_k^{[t]} \right), \quad (14)$$

with $\hat{\boldsymbol{\alpha}}^{[t]}$, $\mathbf{P}^{[t]}$ and $\mathbf{G}^{[t]}$ evaluated at the current estimates $\hat{\sigma}_{k_l}^{2[t]}$ ($l = 1, \dots, p_k$ and $k = 1, \dots, c$), and, when necessary, $\hat{\phi}^{[t]}$.

Note that when $\mathbf{G}_k = \sigma_k^2 \mathbf{I}_{q_k}$, expressions (13) and (14) reduce to those of Harville [expressions (4) and (5), respectively].

An important and desirable property of the updates given in expression (13) is that they are always non-negative, provided that the previous estimates of the variance parameters are non-negative. In addition, under rather weak conditions, these updates are strictly positive (although it is possible to obtain values very close to zero).

Theorem 2 If $\hat{\sigma}_{k_l}^{2[t]} > 0$, then the REML-based estimate updates of the variance parameters given in expression (13) are larger or equal than zero, with strict inequality holding if: (i) $\text{rank}(\mathbf{X}, \mathbf{Z}_k \mathbf{G}_k^{[t]} \boldsymbol{\Lambda}_{k_l}) > \text{rank}(\mathbf{X})$; and (ii) \mathbf{z} (the working response vector) is not in the space spanned by the columns of \mathbf{X} .

The proof is provided in “Appendix B”. We note that condition (i) is needed for both the numerator and denominator of expression (13) to be strictly positive, while (ii) is only needed for the numerator. From an applied point of view, it would undoubtedly important to be able to check whether the conditions are fulfilled before fitting the model. This may not be an easy task, since they depend on $\mathbf{G}_k^{[t]}$ and thus may vary from iteration to iteration. There are, however, common situations where condition (i) could be checked in advance:

- If $\boldsymbol{\Lambda}_{k_l}$ is of full rank, then condition (i) simplifies to $\text{rank}(\mathbf{X}, \mathbf{Z}_k) > \text{rank}(\mathbf{X})$. We note that this condition is the same as that discussed by Harville (1977) in Lemma 1.
- If $\mathbf{G}_k^{[t]}$ and $\boldsymbol{\Lambda}_{k_l}$ commute (i.e., $\mathbf{G}_k^{[t]} \boldsymbol{\Lambda}_{k_l} = \boldsymbol{\Lambda}_{k_l} \mathbf{G}_k^{[t]}$), then condition (i) simplifies to $\text{rank}(\mathbf{X}, \mathbf{Z}_k \boldsymbol{\Lambda}_{k_l}) > \text{rank}(\mathbf{X})$. Examples when $\mathbf{G}_k^{[t]}$ and $\boldsymbol{\Lambda}_{k_l}$ commute include, for instance, when both are diagonal.

We discuss these situations in more detail in Sect. 4, where some examples of application of the SOP method are presented.

3.2 Effective degrees of freedom in the SOP method

In line with the Harville method discussed in Sect. 2, the denominator of expression (13) has been denoted as $\text{ED}_{\{ \cdot \}}$, from effective degrees of freedom. Result (14) makes it easy to show that the sum of the ED_{k_l} over the p_k variance parameters involved in \mathbf{G}_k^{-1} [see Eq. (11)] corresponds to ED_k (the effective degrees of freedom of $\boldsymbol{\alpha}_k$)

$$\begin{aligned} \sum_{l=1}^{p_k} \text{ED}_{k_l} &= \sum_{l=1}^{p_k} \text{trace} \left(\mathbf{Z}_k^\top \mathbf{P} \mathbf{Z}_k \mathbf{G}_k \frac{\boldsymbol{\Lambda}_{k_l}}{\sigma_{k_l}^2} \mathbf{G}_k \right) \\ &= \text{trace} \left(\mathbf{Z}_k^\top \mathbf{P} \mathbf{Z}_k \mathbf{G}_k \right) = \text{trace}(\mathbf{H}_k) = \text{ED}_k. \end{aligned}$$

As a consequence, at convergence, the estimated effective degrees of freedom associated with each random component in model (10) are obtained as a by-product of the SOP method.

To finish this part, we would like to point out an interesting link between the upper bound for ED_{k_l} (denoted with $\text{ED}_{k_l}^{ub}$) and condition (i) in Theorem 2. It can be shown that

$$\begin{aligned} \text{ED}_{k_l} &\leq \text{rank}(X, Z_k G_k \Lambda_{k_l}) - \text{rank}(X) \\ &= \text{rank}((I_n - P_X) Z_k G_k \Lambda_{k_l}) = \text{ED}_{k_l}^{ub}, \end{aligned}$$

where $P_X = X(X^\top X)^{-1} X^\top$. Thus, if condition (i) in Theorem 2 is not verified, ED_{k_l} would be exactly zero. We omit the proof of the previous result. It can be obtained in a similar fashion as in the paper by Cui et al. (2010) (see Web Appendix (f) in that paper), by noting that

$$\begin{aligned} \text{ED}_{k_l} &= \text{trace} \left(Z_k^\top P Z_k G_k \frac{\Lambda_{k_l}}{\sigma_{k_l}^2} G_k \right) \\ &= \text{trace} \left(Z_k G_k \frac{\Lambda_{k_l}}{\sigma_{k_l}^2} G_k Z_k^\top [(I_n - P_X) V (I_n - P_X)]^+ \right), \end{aligned}$$

where Γ^+ denotes the Moore–Penrose pseudoinverse of Γ . This equivalence has been proved, in a less general situation, by Rodríguez-Álvarez et al. (2018) (see Web Appendix D in that paper). Following a similar reasoning, we obtain the upper bound for the effective degrees of freedom of the k -th random component

$$\begin{aligned} \text{ED}_k &\leq \text{rank}(X, Z_k) - \text{rank}(X) \\ &= \text{rank}((I_n - P_X) Z_k) = \text{ED}_k^{ub}. \end{aligned}$$

Using the well-known result that the rank of a matrix sum cannot exceed the sum of the ranks of the summand matrices, we have that $\text{ED}_k^{ub} \leq \sum_{l=1}^{p_k} \text{ED}_{k_l}^{ub}$. In general, however, this is a strict inequality, since most of the cases

$$\mathcal{C}(Z_k G_k \Lambda_{k_u}) \cap \mathcal{C}(Z_k G_k \Lambda_{k_v}) \neq \{0\} \quad (u \neq v),$$

where $\mathcal{C}(A)$ denotes the linear space spanned by the columns of A (see, e.g., Theorem 18.5.7 in Harville 1997). Intuitively, we can interpret this as a sort of competition among the p_k “elements” associated with the k -th random component. The ED_{k_l} cannot vary “free” between 0 and $\text{ED}_{k_l}^{ub}$, but they have to fulfil that their sum does not exceed ED_k^{ub} .

3.3 Computational aspects

From a computational point of view, the evaluation of the expression given in (14) can be very costly. However, this computation can be relaxed by using the result given in (9). For our purpose, it is easy to show that

$$\begin{aligned} \text{trace} \left(Z_k^\top P Z_k G_k \Lambda_{k_l} G_k \right) &= \text{trace} \left(G_k Z_k^\top P Z_k G_k \Lambda_{k_l} \right) \\ &= \text{trace} \left((G_k - C_{kk}^*) \Lambda_{k_l} \right). \end{aligned}$$

This result is exploited in the estimating algorithm described in “Appendix C”. We note that C^* [i.e., the inverse of C in

(2)] is computed in order to estimate $\hat{\beta}$ and $\hat{\alpha}$. Thus, the computation of the traces needed to implement the algorithm reduces to the computation of the diagonal elements of $(G_k - C_{kk}^*) \Lambda_{k_l}$. In addition, in those cases where Λ_{k_l} is diagonal, only the element-wise product of the diagonals of $(G_k - C_{kk}^*)$ and Λ_{k_l} is needed. This will considerably reduce the number of operations required and therefore the computing time.

4 Penalised smoothing and the SOP method

This section discusses several situations in the P-spline framework where estimation can be approached using the SOP method. As it will be seen, the method can be used whenever there are multiple penalties acting on the same coefficients. Anisotropic tensor-product P-splines is an example of overlapping penalties, and it has been extensively discussed in the paper by Rodríguez-Álvarez et al. (2015b). However, multiple penalties arise in a broader class of situations. We describe here two of those: spatially adaptive P-splines and P-splines for hierarchical curve data.

4.1 Spatially adaptive P-splines

Consider a generalised regression problem

$$g(\mu_i) = f(x_i), \quad i = 1, \dots, n, \quad (15)$$

where $\mu_i = \mathbb{E}(y_i)$, $g(\cdot)$ is the link function and $f(\cdot)$ is a smooth and unknown function. We assume further that $\text{Var}(y_i) = \phi v(\mu_i)$. In the P-spline framework (Eilers and Marx 1996), the unknown function $f(x)$ is approximated by a linear combination of d B-splines basis functions, i.e., $f(x) = \sum_{j=1}^d \theta_j B_j(x)$. In matrix notation, model (15) is thus expressed as

$$g(\mu) = B\theta, \quad (16)$$

where $\mu = (\mu_1, \mu_2, \dots, \mu_n)^\top$, $\theta = (\theta_1, \theta_2, \dots, \theta_d)^\top$ and B is a B-spline regression matrix of dimension $n \times d$, i.e., $b_{ij} = B_j(x_i)$ is the j -th B-spline evaluated at x_i . Smoothness is achieved by imposing a penalty on the regression coefficients θ in the form

$$\lambda \sum_{k=q+1}^d (\Delta^q \theta_k)^2 = \lambda \theta^\top D_q^\top D_q \theta, \quad (17)$$

where λ is the smoothing parameter and Δ^q forms differences of order q on adjacent coefficients, i.e., $\Delta \theta_k = \theta_k - \theta_{k-1}$, $\Delta^2 \theta_k = \Delta(\Delta \theta_k) = \theta_k - \theta_{k-1} - (\theta_{k-1} - \theta_{k-2}) = \theta_k - 2\theta_{k-1} + \theta_{k-2}$, and so on for higher q . Finally, D_q is simply the matrix representation of Δ^q .

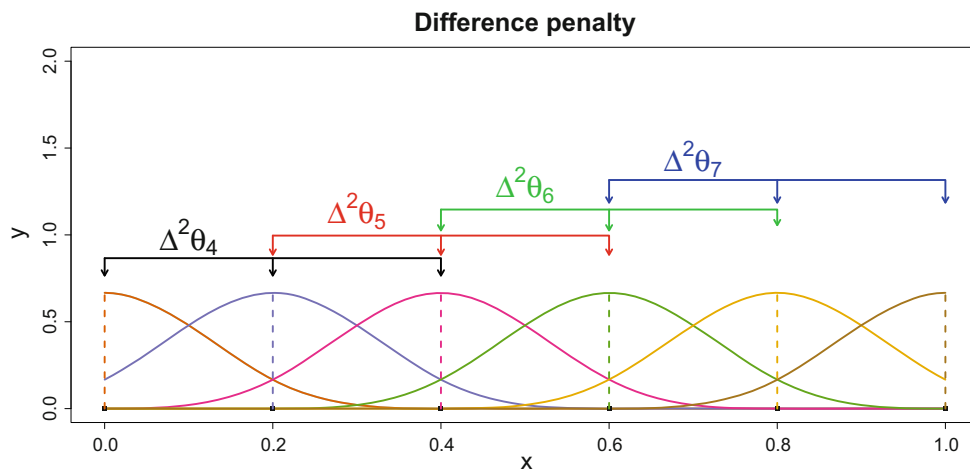


Fig. 1 Graphical representation of differences of order 2 on adjacent coefficients of cubic B-splines basis functions. Note the local and ordered nature of these differences

As can be observed in Eqn. (17), the same smoothing parameter λ applies to all coefficient differences, irrespective of their location (see Fig. 1). Thus, the model assumes that the same amount of smoothing is needed across the whole domain of the covariate. Adaptive P-splines (see, e.g., Krivobokova et al. 2008; Ruppert and Carroll 2000 among others) relax this assumption. The idea is simple, to replace the global smoothing parameter by smoothing parameters that vary locally according to the covariate value. This can be accomplished by specifying a different smoothing parameter for each coefficient difference (Ruppert and Carroll 2000; Wood 2011)

$$\sum_{k=q+1}^d \lambda_{k-q} (\Delta^q \theta_k)^2 = \boldsymbol{\theta}^\top \mathbf{D}_q^\top \text{diag}(\boldsymbol{\lambda}) \mathbf{D}_q \boldsymbol{\theta}, \quad (18)$$

where $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_{d-q})^\top$. Note that this approach would imply as many smoothing parameters as coefficient differences (i.e., $d - q$), which could lead to undersmoothing and unstable computations. Given the local and ordered nature of the coefficient differences (see Fig. 1), we may model the smoothing parameters λ_k as a smooth function of k (its position) and use B-splines for this purpose (here no penalty is assumed)

$$\boldsymbol{\lambda} = \boldsymbol{\Psi} \boldsymbol{\xi}, \quad (19)$$

where $\boldsymbol{\Psi}$ is a B-spline regression matrix of dimension $(d - q) \times p$ with $p < (d - q)$ and $\boldsymbol{\xi} = (\xi_1, \dots, \xi_p)^\top$ is the new vector of smoothing parameters. Performing some simple algebraic operations, it can be shown that the *adaptive* penalty (18) is

$$\boldsymbol{\theta}^\top \left(\sum_{l=1}^p \xi_l \mathbf{D}_q^\top \text{diag}(\boldsymbol{\psi}_l) \mathbf{D}_q \right) \boldsymbol{\theta}, \quad (20)$$

where $\boldsymbol{\psi}_l$ denotes the column l of $\boldsymbol{\Psi}$. Note that under this adaptive penalty, all coefficients are penalised by multiple smoothing parameters, i.e., there are overlapping penalties.

4.1.1 Mixed model reparametrisation

Estimation of the P-spline model (16) subject to the adaptive penalty defined in (20) can be carried out based on the connection between P-splines and mixed models (e.g., Currie and Durban 2002; Wand 2003). It is easy to show that the null space (i.e., the unpenalised function space) of the adaptive penalty matrix $\mathbf{P}_{Ad} = \sum_{l=1}^p \xi_l \mathbf{D}_q^\top \text{diag}(\boldsymbol{\psi}_l) \mathbf{D}_q$ is independent of $\boldsymbol{\xi}$. In addition, note that when $\xi_u = \xi_v = \lambda$ ($\forall u, v$) then

$$\begin{aligned} \mathbf{P}_{Ad} &= \sum_{l=1}^p \xi_l \mathbf{D}_q^\top \text{diag}(\boldsymbol{\psi}_l) \mathbf{D}_q = \lambda \mathbf{D}_q^\top \left(\sum_{l=1}^p \text{diag}(\boldsymbol{\psi}_l) \right) \mathbf{D}_q \\ &= \lambda \mathbf{D}_q^\top \mathbf{D}_q = \mathbf{P}. \end{aligned}$$

This is a consequence of the rows of a B-spline regression matrix adding up to 1. Thus, the null space of \mathbf{P}_{Ad} is the same as that of \mathbf{P} . Different reparametrisations of P-spline models have been suggested in the literature (see, e.g., Currie and Durban 2002; Eilers 1999), all aiming to decompose the model into the unpenalised and the penalised part. The consequence of this decomposition is that the penalty matrix of the reparametrised P-spline model is of full rank and so is the precision matrix of the corresponding mixed model. For our application, we use the proposal given in Eilers (1999). As will be seen, this approach gives rise to diagonal precision matrices. As discussed in Sect. 3.3, this is very convenient

from a computational point of view. Using Eilers (1999)'s transformation, model (16) is re-expressed as

$$g(\boldsymbol{\mu}) = \mathbf{B}\boldsymbol{\theta} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\alpha},$$

where $\mathbf{X} = [\mathbf{1}_n | \mathbf{x} | \dots | \mathbf{x}^{(q-1)}]$ and $\mathbf{Z} = \mathbf{B}\mathbf{D}_q^\top (\mathbf{D}_q \mathbf{D}_q^\top)^{-1}$, and the precision (penalty) matrix of the vector of random (penalised) coefficients $\boldsymbol{\alpha}$ becomes

$$\mathbf{G}^{-1} = \frac{1}{\phi} \mathbf{F}^\top \mathbf{P}_{Ad} \mathbf{F} = \sum_{l=1}^p \sigma_l^{-2} \text{diag}(\boldsymbol{\psi}_l) = \sum_{l=1}^p \sigma_l^{-2} \tilde{\mathbf{\Lambda}}_l, \quad (21)$$

where $\mathbf{F} = \mathbf{D}_q^\top (\mathbf{D}_q \mathbf{D}_q^\top)^{-1}$, $\tilde{\mathbf{\Lambda}}_l = \text{diag}(\boldsymbol{\psi}_l)$, and $\sigma_l^2 = \phi/\xi_l$. Thus, the precision matrix is linear in the precision parameters σ_l^{-2} , and the SOP method can therefore be used.

We note that the model has a single random component ($c = 1$). The reparametrisation ensures that $\text{rank}(\mathbf{X}, \mathbf{Z}) = \text{rank}(\mathbf{X}) + \text{rank}(\mathbf{Z})$. Thus, $\text{rank}(\mathbf{X}, \mathbf{Z}\tilde{\mathbf{\Lambda}}_l) = \text{rank}(\mathbf{X}) + \text{rank}(\mathbf{Z}\tilde{\mathbf{\Lambda}}_l) > \text{rank}(\mathbf{X})$. Exploiting the fact that \mathbf{G} and $\tilde{\mathbf{\Lambda}}_l$ are diagonal, and thus commute, condition (i) of Theorem 2 is satisfied.

4.2 P-splines for hierarchical curve data

For simplicity, let's assume balanced hierarchical curve data. Our data consist on m individuals each with s different measurements at times $\mathbf{t} = (t_1, t_2, \dots, t_s)$. Our interest is focused on

$$g(\mu_{ij}) = f(t_i) + g_j(t_i), \quad 1 \leq i \leq s, \quad 1 \leq j \leq m, \quad (22)$$

where $\mu_{ij} \equiv \mathbb{E}(y_{ij})$, y_{ij} is the response variable on the j -th subject at time t_i , $f(\cdot)$ is a function describing the population effect, and $g_j(\cdot)$ are random functions measuring the deviation of the j -th subject from the population effect. As before, $\text{Var}(y_{ij}) \equiv \phi v(\mu_{ij})$. A simple model would consist in a parametric specification for $f(\cdot)$ and $g_j(\cdot)$, e.g., $f(t) = \beta_0 + \beta_1 t$ and $g_i(t) = \alpha_{0j} + \alpha_{1j} t$, with $\boldsymbol{\alpha}_0 = (\alpha_{01}, \dots, \alpha_{0m})^\top \sim N(\mathbf{0}, \sigma_0^2 \mathbf{I}_m)$ and $\boldsymbol{\alpha}_1 = (\alpha_{11}, \dots, \alpha_{1m})^\top \sim N(\mathbf{0}, \sigma_1^2 \mathbf{I}_m)$.

A more flexible approach consists in assuming $f(\cdot)$ and $g_j(\cdot)$ smooth and unknown. Important contributions in the P-spline framework can be found in Durban et al. (2005) and Ruppert et al. (2003). Both approaches are based on modelling $f(\cdot)$ and $g_j(\cdot)$ using truncated line basis, and estimation is based on penalised methods and linear mixed model techniques. More recently, Djeundje and Currie (2010) extend those models and propose the inclusion of an extra penalty for the individual curve coefficients. The authors argue that this extra penalty is needed to address identifiability issues when estimating the population effect. Under the

so-called M0 model in Djeundje and Currie (2010)'s paper, model (22) is expressed in matrix notation as

$$g(\boldsymbol{\mu}_{\cdot j}) = \underbrace{\mathbf{B}\boldsymbol{\theta}}_{f(\mathbf{t})} + \underbrace{\check{\mathbf{B}}\check{\boldsymbol{\theta}}_j}_{g_j(\mathbf{t})},$$

where $\boldsymbol{\mu}_{\cdot j} = (\mu_{1j}, \dots, \mu_{sj})^\top$ is the mean response vector for the j -th individual and \mathbf{B} and $\check{\mathbf{B}}$ are B-spline regression matrices of possibly different size for, respectively, the population curve and the individual deviation. The vector $\boldsymbol{\theta}$ is assumed fixed, but subject to a q -th order penalty of the form $\mathbf{P} = \lambda_1 \mathbf{D}_q^\top \mathbf{D}_q$, and $\check{\boldsymbol{\theta}}_j$ is a random vector with distribution $N(\mathbf{0}, \phi \check{\mathbf{P}}^{-1})$ where

$$\check{\mathbf{P}} = \lambda_2 \check{\mathbf{D}}_q^\top \check{\mathbf{D}}_q + \lambda_3 \mathbf{I}_{\check{d}}. \quad (23)$$

The first term $\lambda_2 \check{\mathbf{D}}_q^\top \check{\mathbf{D}}_q$ is responsible for the smoothness of the individuals curves, whereas $\lambda_3 \mathbf{I}_{\check{d}}$ addresses the identifiability aspect (see Djeundje and Currie 2010 for more details). Note that each random effect is shrunk (penalised) by both smoothing parameters λ_2 and λ_3 , and thus the precision matrix $\check{\mathbf{G}}^{-1} = 1/\phi \check{\mathbf{P}}$ is linear in the precision parameters

$$\check{\mathbf{G}}^{-1} = \sum_{l=2}^3 \sigma_l^{-2} \check{\mathbf{\Lambda}}_l, \quad (24)$$

where $\sigma_2^2 = \phi/\lambda_2$ and $\sigma_3^2 = \phi/\lambda_3$ and $\check{\mathbf{\Lambda}}_2 = \check{\mathbf{D}}_q^\top \check{\mathbf{D}}_q$ and $\check{\mathbf{\Lambda}}_3 = \mathbf{I}_{\check{d}}$.

In more compact way, we express the model for the whole sample as

$$g(\boldsymbol{\mu}) = [\mathbf{1}_m \otimes \mathbf{B}]\boldsymbol{\theta} + [\mathbf{I}_m \otimes \check{\mathbf{B}}]\check{\boldsymbol{\theta}}, \quad (25)$$

where \otimes denotes the Kronecker product, $\boldsymbol{\mu} = (\boldsymbol{\mu}_{\cdot 1}^\top, \dots, \boldsymbol{\mu}_{\cdot m}^\top)^\top$, $\check{\boldsymbol{\theta}} = (\check{\boldsymbol{\theta}}_1^\top, \dots, \check{\boldsymbol{\theta}}_m^\top)^\top$ and

$$\check{\boldsymbol{\theta}} \sim N(\mathbf{0}, \mathbf{I}_m \otimes \check{\mathbf{G}}).$$

One last step is needed in order to apply the SOP method for the estimation of model (25): the decomposition of the population effect into the unpenalised and the penalised part. We use here the approach based on the eigenvalue decomposition (EVD) of the penalty. Let $\mathbf{D}_q^\top \mathbf{D}_q = \mathbf{U}\boldsymbol{\Sigma}\mathbf{U}^\top$ be the EVD of $\mathbf{D}_q^\top \mathbf{D}_q$. Here \mathbf{U} denotes the matrix of eigenvectors and $\boldsymbol{\Sigma}$ the diagonal matrix of eigenvalues. Let us also denote by $\mathbf{U}_+(\boldsymbol{\Sigma}_+)$ and $\mathbf{U}_0(\boldsymbol{\Sigma}_0)$ the sub-matrices corresponding to the nonzero and zero eigenvalues, respectively. In this case, model (25) is re-expressed as

$$g(\boldsymbol{\mu}) = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\alpha},$$

where $\beta = U_+^\top \theta$, $\alpha = ((U_0^\top \theta)^\top, \check{\theta}^\top)^\top$, $X = [\mathbf{1}_m \otimes BU_0]$ and $Z = [\mathbf{1}_m \otimes BU_+ : \mathbf{I}_m \otimes \check{B}]$. Finally, $\alpha \sim N(\mathbf{0}, G)$, where

$$G^{-1} = \begin{pmatrix} \frac{1}{\sigma_1^2} \Sigma_+ & \mathbf{0}_{d \times (m\check{d})} \\ \mathbf{0}_{(m\check{d}) \times d} & \mathbf{I}_m \otimes \check{G}^{-1} \end{pmatrix} = \sum_{l=1}^3 \sigma_l^{-2} \tilde{\Lambda}_l,$$

with $\sigma_l^2 = \frac{\phi}{\lambda_l}$ and

$$\begin{aligned} \tilde{\Lambda}_1 &= \begin{pmatrix} \Sigma_+ & \mathbf{0}_{d \times (m\check{d})} \\ \mathbf{0}_{(m\check{d}) \times d} & \mathbf{0}_{(m\check{d}) \times (m\check{d})} \end{pmatrix} \\ \tilde{\Lambda}_2 &= \begin{pmatrix} \mathbf{0}_{d \times d} & \mathbf{0}_{d \times (m\check{d})} \\ \mathbf{0}_{(m\check{d}) \times d} & \mathbf{I}_m \otimes \check{\Lambda}_2 \end{pmatrix} \\ \tilde{\Lambda}_3 &= \begin{pmatrix} \mathbf{0}_{d \times d} & \mathbf{0}_{d \times (m\check{d})} \\ \mathbf{0}_{(m\check{d}) \times d} & \mathbf{I}_m \otimes \check{\Lambda}_3 \end{pmatrix}. \end{aligned}$$

The precision matrix is linear in the precision parameters and thus appropriate for the SOP method. There are, in this case, two random components, modelling, respectively, the population curve and the individual deviations, with

$$Z_1 = \mathbf{1}_m \otimes BU_+ \quad \text{and} \quad G_1^{-1} = \sigma_1^{-2} \Sigma_+,$$

$$Z_2 = \mathbf{I}_m \otimes \check{B} \quad \text{and}$$

$$\begin{aligned} G_2^{-1} &= \sigma_2^{-2} \mathbf{I}_m \otimes \check{\Lambda}_2 + \sigma_3^{-2} \mathbf{I}_m \otimes \check{\Lambda}_3 \\ &= \sigma_2^{-2} \mathbf{I}_m \otimes \check{D}_q^\top \check{D}_q + \sigma_3^{-2} \mathbf{I}_m \otimes \check{I}_{\check{d}}. \end{aligned}$$

Recall that $X = \mathbf{1}_m \otimes BU_0$, and note that $\text{rank}(\mathbf{1}_m \otimes BU_0) = \text{rank}(BU_0)$. We now show that condition (i) of Theorem 2 is satisfied for all variance parameters:

σ_1^2 : By construction, Σ_+ is positive definite and of full rank, and $\text{rank}(X, Z_1) = \text{rank}(X) + \text{rank}(Z_1)$. Noting that $\text{rank}(X, Z_1 G_1 \Sigma_+) = \text{rank}(X, Z_1) > \text{rank}(X)$, the condition is verified.

σ_2^2 : By, e.g., Corollary 18.2.2 in Harville (1997), G_2 and $\mathbf{I}_m \otimes \check{D}_q^\top \check{D}_q$ commute. Thus, as long as $m > 1$, it is easy to show that

$$\begin{aligned} \text{rank}(X, Z_2 \mathbf{I}_m \otimes \check{D}_q^\top \check{D}_q) \\ = \text{rank}(\mathbf{1}_m \otimes BU_0, \mathbf{I}_m \otimes \check{B} \check{D}_q^\top \check{D}_q) > \text{rank}(BU_0). \end{aligned}$$

σ_3^2 : Note that G_2 and $\mathbf{I}_m \otimes \check{I}_{\check{d}}$ commute. As before, as long as $m > 1$, it is easy to show that

$$\begin{aligned} \text{rank}(X, Z_2 \mathbf{I}_m \otimes \check{I}_{\check{d}}) &= \text{rank}(\mathbf{1}_m \otimes BU_0, \mathbf{I}_m \otimes \check{B}) \\ &> \text{rank}(BU_0). \end{aligned}$$

5 Examples

This section presents several data examples where the SOP method represents a powerful alternative to existing estimation procedures. We discuss three different analyses: the first two examples are concerned with spatially adaptive P-splines, but each of them deals with a different situation regarding complexity and aim; the last example is devoted to illustrating our method for the analysis of hierarchical curve data. All computations were performed in (64-bit) R 3.4.4 (Core Team 2018), and a 2.30GHz \times 4 Intel® Core™i5 processor computer with 15.6GB of RAM and Ubuntu 16.04 LTS operating system.

5.1 Doppler function

For our first example, we consider the Doppler function. This is a common example in the adaptive smoothing literature and has been discussed by Ruppert and Carroll (2000), Krivobokova et al. (2008) and Tibshirani (2014), among others. Data are generated according to

$$y_i = \sin(4/x_i) + 1.5 + \varepsilon_i, \quad i = 1, \dots, n,$$

where $x_i \sim U[0, 1]$, $\varepsilon_i \sim N(0, 0.2^2)$, and $n = 1000$. For fitting the data, we assume the spatially adaptive P-spline model discussed in Sect. 4.1. We compare the performance of the SOP method with that implemented in the R-package `mgcv`, version 1.8-23, and described in Wood (2011). It is worth noticing that both approaches implement in essence the same adaptive P-spline model; the only difference is the estimation procedure (and, possibly, the reparametrisation). In addition, we also fit the model without assuming an adaptive penalty. In this case, the SOP method reduces to Harville's approach (see Sect. 2). In all cases, we use 200 cubic B-splines to represent the smooth function, jointly with second-order differences. For the adaptive approaches, 15 equally spaced cubic B-splines are used for the smoothing parameters [see Eq. (19)]. These values are chosen to provide enough flexibility to the model. Under this configuration, there are a total of 15 variance parameters. Figure 2a shows the true simulated Doppler function. Figure 2b, c shows, respectively, the estimated curves based on the SOP method without and with an adaptive penalty. Results using `mgcv` are depicted in Fig. 2d. As expected, both adaptive approaches perform similarly. With the specified configuration, they are able to capture 7 cycles of the Doppler function. On the other hand, the non-adaptive approach is able to capture only 4 cycles and presents very wiggly estimates, especially on the right-hand side of the covariate domain. In terms of EDs, for the SOP model without adaptive penalty, we obtain a total ED of 95.8 (out of 200). For models with adaptive smoothing, we obtained identical results, i.e., 50.2 (with SOP) and 50.0 (with

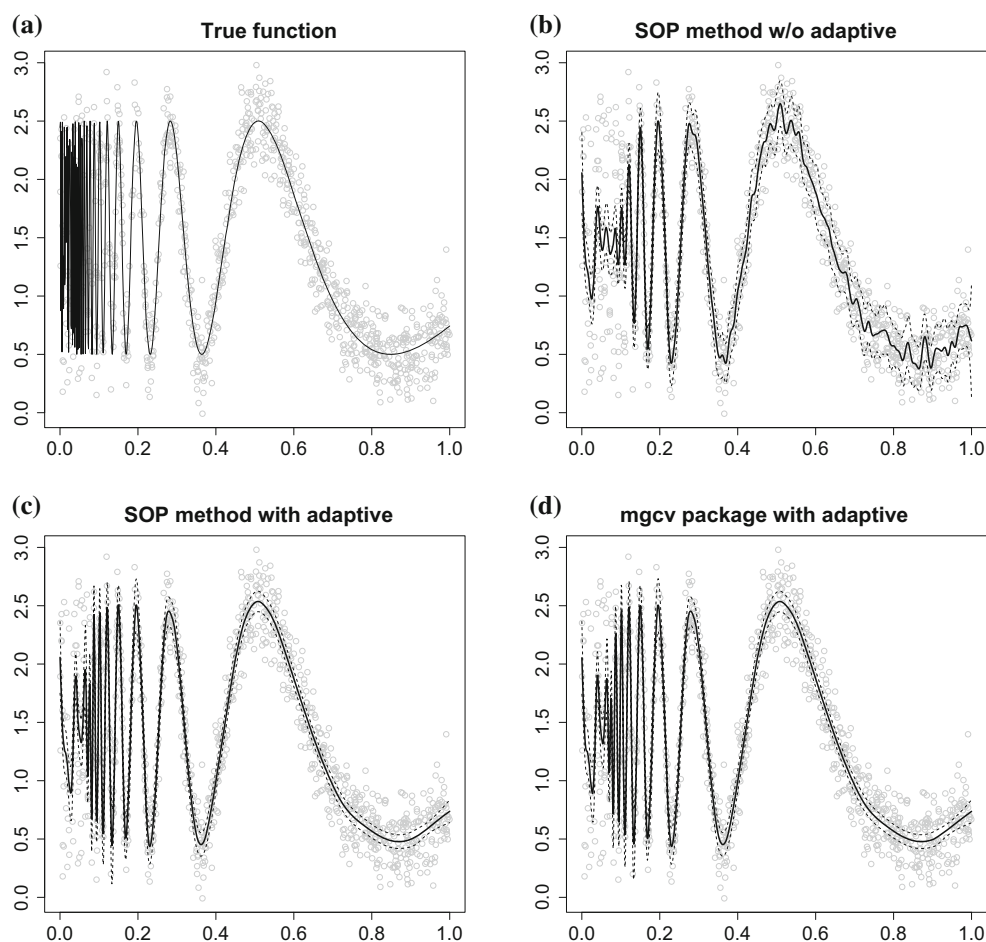


Fig. 2 For the Doppler function: **a** true function (solid line) and simulated data points (grey points), **b** estimated curve using the SOP method without adaptive penalty (solid line) jointly with 95% approximate confidence intervals (dotted lines), **c** estimated curve using the SOP method

and adaptive penalty (solid line) jointly with 95% approximate confidence intervals (dotted lines); and **d** estimated curve using the *mgcv* package (solid line) jointly with 95% approximate confidence intervals (dotted lines)

mgcv). It is worth remembering that using the SOP method the total ED is obtained by adding up the EDs associated with each variance parameter in the model (plus the dimension of the fixed part). These EDs are the denominator of the estimate update expressions of the variance parameters. The gain of the SOP method is clear when we compare the computing times: 1.0 second with our approach (0.4 if we do not consider adaptive), and 45 seconds using *mgcv*. To gain more insights into the performance of both adaptive approaches beyond the visual similarity shown in Fig. 2 for a single dataset, we run a small simulation study. We perform a total of 100 replicates with the same specification described above and make comparisons regarding root-mean-square error (RMSE) and computing time. In terms of RMSE, the mean (standard deviation) over the 100 replicates is 0.265 (0.0113) for SOP and 0.267 (0.0121) using *mgcv*. Concerning computing times, the mean (standard deviation)—in seconds—is 0.84 (0.22) and 52.79 (30.77) for, respectively, SOP and *mgcv*.

5.2 X-ray diffraction data

For this example, we use data from a X-ray crystallography radiation scan of a thin layer of indium tin oxide. X-ray crystallography allows the exploration of the molecular and atomic structure of crystals. Crystallographers precisely rotate the crystal by entering its desired orientation, while it is illuminated by a beam of X-rays. Depending on the angle, the number of diffracted photons varies and they are detected and counted by an electronic detector. The dataset was analysed by Davies et al. (2008, 2013) and can be found in the R-package *diffractometry* as *indiumoxide*. Figure 3 shows such an X-ray diffraction scan (grey lines). The aim of X-ray diffraction analysis is to determine (a) the signal baseline (and remove it); and (b) the number of peaks (and isolate them to further analysis of their position, height, symmetry, and so forth). This example is solely included to illustrate the potentiality of the method presented in this paper for the anal-

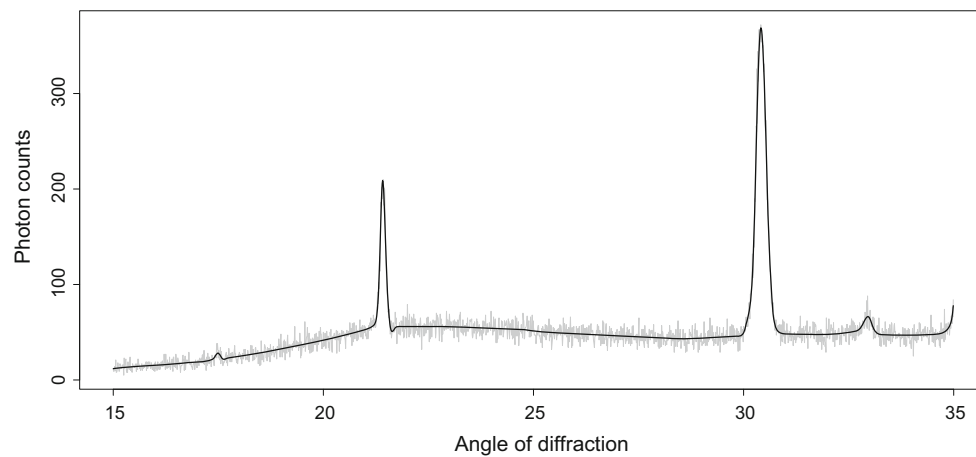


Fig. 3 For the X-ray radiation example: estimated smooth effect of the angle of diffraction on the X-ray radiation using the SOP method (solid black line). The grey lines represent the raw data

ysis of very complex data. For a different modelling approach applied to this dataset, see Camarda et al. (2016). Given that the outcome variable represents count data, a Poisson model is adopted

$$\log(\mathbb{E}[y_i]) = \log(\mu_i) = f(x_i),$$

where y_i and x_i ($i = 1, \dots, 2000$) denote, respectively, the photon counts and the angle of diffraction. To provide enough flexibility to the model in order to make it able to capture the peaks (see Fig. 3), we use 200 cubic B-splines and second-order differences for the function, and 80 cubic B-splines for the adaptive penalty. Results are shown in Fig. 3. The results using the `mgcv` package are almost identical to our proposal and are not depicted. In this case, our method takes less than 3 s, whereas `mgcv` is around 750 times slower (33 min). Regarding the EDs, we obtain a total of 29.5 and 24.9 using SOP and `mgcv`, respectively.

5.3 Diffusion tensor imaging scan data

Our last example deals with hierarchical curve data. We analyse the DTI dataset that can be found in the R-package `refund` (Goldsmith et al. 2016). A detailed description of the study and data can be found in Goldsmith et al. (2011), Goldsmith et al. (2012) and Greven and Scheipl (2017). In brief, the study aimed at comparing the white matter tracts in patients affected by multiple sclerosis (MS) and healthy individuals. Multiple sclerosis is a disease of the central nervous system that causes lesions in white matter tracts, thus interrupting the travel of nerve impulses to and from the brain and spinal cord. Diffusion tensor imaging (DTI) is a magnetic resonance imaging technique which makes it possible to study the white matter tracts by tracing the diffusion of water on the brain. From DTI scans, fractional anisotropy (FA) mea-

surements can be extracted. FA is related to water diffusion and thus to MS diagnosis and progression. The DTI dataset contains FA measurements of healthy and diseased individuals, recorded at several locations along the callosal fibre tract (CCA) and the right corticospinal tract in the brain. In Fig. 4a, the observed FA measurements at different tract locations in the CAA are shown, separately for cases and controls, i.e., individuals affected and non-affected with MS. Each line in these plots represents an individual, and only the first visit is considered. Note that, in general, MS patients present lower FA measurements than healthy individuals.

For illustration purposes, we present two different analyses and comparisons. We first focus our interest on the subgroup of individuals affected with MS. In this group, there are a total of $m = 99$ individuals, each with $s = 93$ FA measurements at different CAA tract locations. The SOP method is used to estimate the model described in Sect. 4.2 (for Gaussian homoskedastic errors) and presented in Djeundje and Currie (2010). To compare results and computing times, the code associated with the paper by Djeundje and Currie (2010) is also tested. For this example, both implementations take the advantage of the array structure of the data: generalised linear array models (GLAM, Currie et al. 2006) are used to efficiently compute the inner products for the Henderson equations [Eq. (27)]. Here, 43 cubic B-spline basis are used for the population curve, and 23 for the individual curves. This configuration gives rise to a model with 2320 ($= 43 + 23 \times 99$) coefficients (both random and fixed). SOP method needs about 150 seconds to fit the model, and Djeundje and Currie (2010)'s code is 14 times slower. We note that the computational time can be further improved if the sparse structure of the matrices involved in the model is exploited. Using the R-package `Matrix`, we are able to reduce the computing time using SOP to 35 seconds. Figure 4b shows the estimated population effect

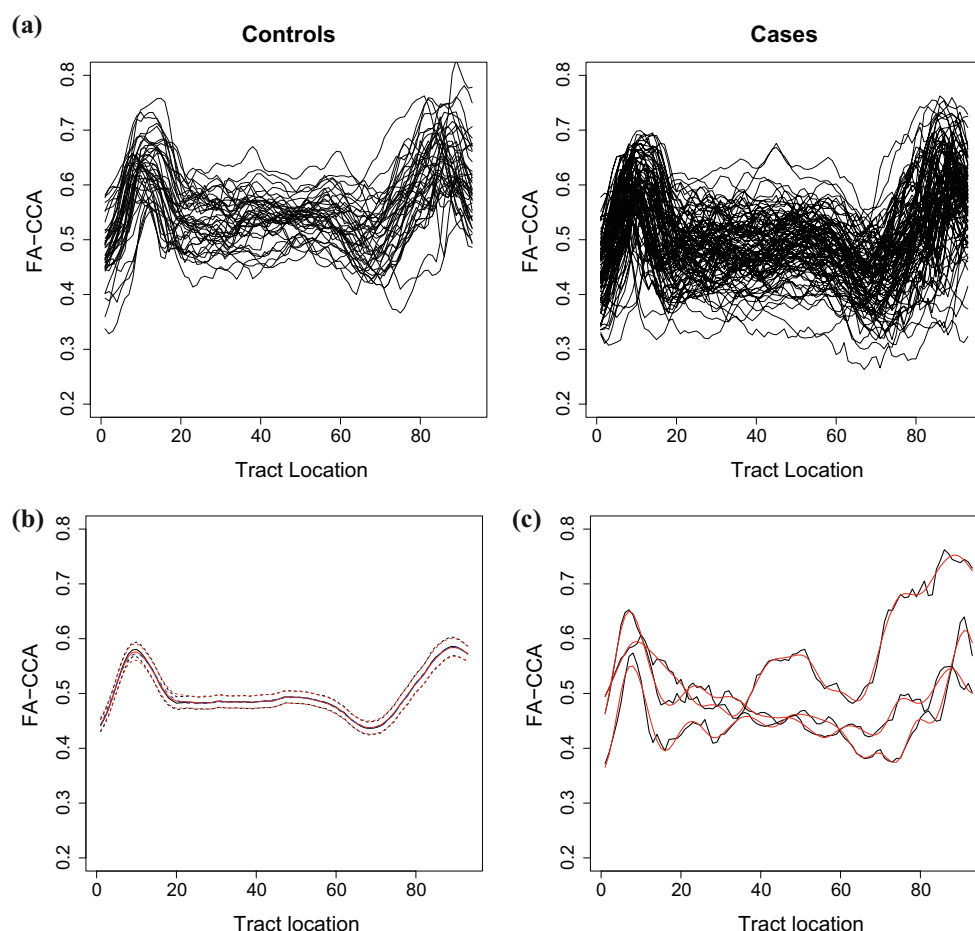


Fig. 4 For the diffusion tensor imaging scan data: **a** observed FA values along the CCA tract. Left: healthy controls. Right: MS patients; **b** estimated population (group) FA profile for individuals affected by MS. Solid red line: SOP method. Dotted blue line: Djeundje and Currie (2010)’s code. The dotted red lines are the pointwise 95% confidence

intervals based on the full-sandwich standard errors proposed by Heckman et al. (2013). The back line is the observed mean, and the dotted black lines are the empirical confidence intervals; and **c** estimated (red lines) and observed (black lines) individual FA profiles for 3 selected individuals. (Color figure online)

using both approaches, which provide very similar results. 95% pointwise confidence intervals are calculated using the full-sandwich standard errors proposed by Heckman et al. (2013) but adapted to our case. Figure 4c shows, for several MS patients, the estimated (and observed) FA profiles. In terms of EDs, we obtain 35.03 (out of 43) for the population curve (including the unpenalised or fixed part), and a total of 2025.78 (out of 2275 ($= 23 \times 99 - 2$)) for all individual curves. We note that this total corresponds to the sum of the EDs associated with the variances σ_2^2 and σ_3^2 involved in the modelling of these individual curves [see expressions (23) and (24) for details]. More precisely, using the SOP method we obtain an ED of 870.44 for σ_2^2 and of 1155.34 for σ_3^2 .

Our second analysis considers all individuals, cases and controls. The interest here is to compare the FA profiles at the first visit between these two groups. To that aim, the following factor-by-curve interaction model is considered.

$$y_{ij} = f_{z_j}(t_i) + g_j(t_i) + \varepsilon_{ij} \quad 1 \leq i \leq s, \quad 1 \leq j \leq m,$$

where $z_j = 1$ if the j -th individual is affected by MS (case) and $z_j = 0$ otherwise (control), and $\varepsilon_{ij} \sim N(0, \phi)$. Note that this model assumes a different FA profile for each group. For this analysis, there are a total of $m_0 = 42$ controls and $m_1 = 99$ MS patients ($m = m_0 + m_1 = 141$), and $s = 93$ different tract locations. A detailed description of the model can be found in “Appendix D”. As for the first analysis 43 cubic B-spline basis are used for the population curves (FA profiles), and 23 for the individual curves, yielding a total of 3329 ($= 43 \times 2 + 23 \times 141$) coefficients. Using GLAM and sparse matrix techniques (R-package *Matrix*), the fit takes 65 seconds. Figure 5a shows the estimated FA profiles for both cases and controls, jointly with 95% pointwise confidence intervals (Heckman et al. 2013). The ED for the FA profile in controls is 32.21 and in MS patients is 35.55. In both cases, we include the fixed part. Regarding the indi-

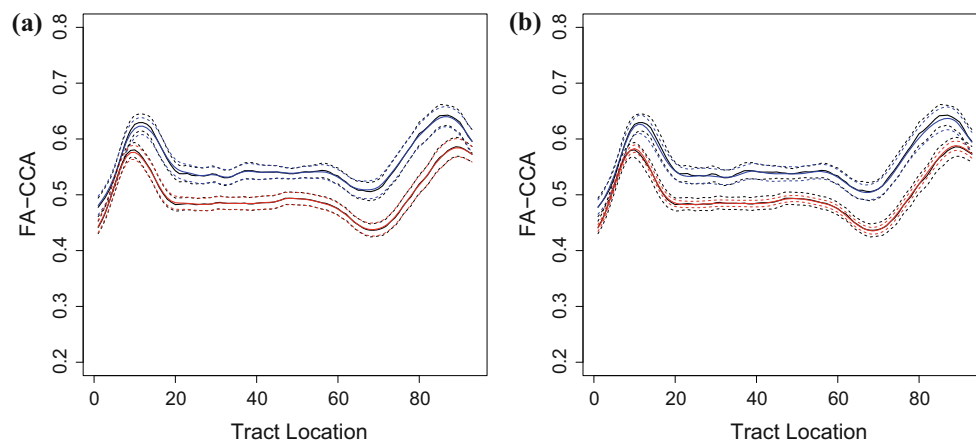


Fig. 5 For the diffusion tensor imaging scan data: estimated population (group) FA profiles. From left to right: results using the SOP method and the functional regression approach by Greven and Scheipl (2017). Solid red line: MS patients. Solid blue line: controls. The dotted blue and red lines are the pointwise 95% confidence intervals. For the SOP

method confidence intervals are based on the full-sandwich standard errors proposed by Heckman et al. (2013). The solid back line is the observed mean, and the dotted black lines are the empirical confidence intervals. (Color figure online)

vidual curves, we obtain a total ED of 2863.46, the sum of 1263.26 and 1600.20.

We compare the results with the functional regression model presented in the paper by Greven and Scheipl (2017) [see model (1.1) and Fig. 2 in that paper]. We would like to note that the P-spline model used in this paper was discussed in Durban and Aguilera-Morillo (2017) as a competitive alternative to the functional approach by Greven and Scheipl (2017). For the functional approach, we consider Gaussian homoskedastic errors, and, as suggested by the authors, 25 cubic B-spline basis and a first-order difference penalty, as well as 8 functional principal components (FPC) functions. The code for fitting the model was kindly provided by the authors. Results are depicted in Fig. 5b. As can be observed, both approaches provide very similar results. However, note that the pointwise 95% confidence interval for the estimated FA profile in MS patients is narrower than the empirical confidence interval. This results can be explained by the (possibly wrong) assumption of Gaussian homoskedastic errors. In terms of computing times, the functional regression model needs 895 s to be fitted (in contrast to 65 s using SOP). We are aware that the computing times of both approaches (the SOP method and functional approach) are not fully comparable since they assume different model specifications.

6 Discussion

This paper presents a new estimation method, called SOP, for (generalised) linear mixed models. The method is an extension of a previous proposal by Harville (1977), and generalised by Schall (1991) among others. In contrast to those

previous approaches, the SOP method is suitable for models where the precision matrix is linear in the precision parameters. These precision matrices are common when penalised smooth models are reformulated as (generalised) linear mixed models. They appear when there are multiple quadratic penalties acting on the same regression coefficients. One special case is anisotropic tensor-product P-splines. This situation was discussed in the paper by Rodríguez-Álvarez et al. (2015b) where the SAP algorithm was proposed. The present paper goes one step further and generalises SAP to a more general case. SOP is, as far as we know, the only method in this context where the variance parameters estimates involve “partial” effective degrees of freedom. As a consequence, the method provides, at convergence, the estimated effective partial degrees of freedom associated with each smooth/random component in the model. This is especially relevant when working with smoothing models, where the effective degrees of freedom of each smooth term give important insights on the complexity of the fitted function. Furthermore, we show in the paper that the SOP method ensures non-negative estimates of the variance parameters and discuss the conditions under which these estimates are strictly positive.

We present in the paper several situations in the context of penalised spline regression in which SOP represents a powerful alternative to existing estimation methods. In particular, we show the use of SOP in the case of spatially adaptive P-splines and for the estimation of subject-specific curves in longitudinal data. We discuss several real data analyses dealing with these situations and show the outperformance of SOP in terms of computing times. We use simple modelling situations with the aim of describing the method, models and examples in a detailed way, avoiding generalisations

that could complicate the reading of the paper and obscure the simplicity of the proposal. However, there are several other fields of application of SOP. For instance, overlapping penalties appear in brain imaging research (Karas et al. 2017) and derivative curve estimation (Simpkin and Newell 2013). Also, the method can be used for more complex models including linear effects (multidimensional) smooth functions, random Gaussian effects, etc. The use of other basis functions beyond B-splines is also possible as long as quadratic penalties are combined.

The proposal and examples presented in this paper pave the way for further research efforts. For instance, the approach discussed for adaptive P-splines is based on smoothing the locally varying smoothing parameter by means of B-splines. This implies that smoothness is solely controlled by the number of B-spline basis. The selection of the appropriate basis dimension may not be an easy task, with the undesirable consequence that if a large basis dimension is chosen (larger than needed), we may end up with a local linear fit. To reduce the impact of the basis dimension, we will explore the inclusion of a penalty on the coefficients (variance parameters) associated with the B-spline basis. This can be accomplished by means of a hierarchical structure for the random effects (see, e.g., Krivobokova et al. 2008). Another challenging field is the study of suitable penalties and efficient estimation methods for adaptive P-splines in more than one dimension. Whereas some attempts have been done in two dimensions (see, e.g., Crainiceanu et al. 2007; Krivobokova et al. 2008), to the best of our knowledge the literature is lacking in three dimensional approaches (e.g., space and time). Some preliminary results using SOP are available at Rodríguez-Álvarez et al. (2016), but further work still needs to be done. Also, the use and extension of the SOP method to the analysis of multilevel and longitudinal functional data (see, e.g., Greven and Scheipl 2017) represents an exciting field of research. Finally, for variable selection there exists some works that propose sparse regression models using local quadratic approximations to the L1-norm adopting a penalised likelihood approach (see Fan and Li 2001; Hunter and Li 2005; Zou and Li 2008). More recently, (generalised) linear mixed effects approaches have also been proposed (see Taylor et al. 2012; Groll and Tutz 2014) allowing for the penalty to be estimated simultaneously with the variance parameters using REML. We intend to extend the SOP method in this direction. In conclusion, this paper opens up a pathway for a general estimating method allowing for both smoothing and variable selection in reasonable computing times.

The R-code used for the real data examples presented in Sect. 5 as well as an R-package implementing the SOP method for generic generalised linear mixed models, spatially adaptive P-spline models and P-spline models for

hierarchical curve data can be downloaded from <https://bitbucket.org/mxrodriguez/sop>.

Acknowledgements This research was supported by the Basque Government through the BERC 2018-2021 program and by Spanish Ministry of Economy and Competitiveness MINECO through BCAM Severo Ochoa excellence accreditation SEV-2013-0323 and through projects MTM2017-82379-R funded by (AEI/FEDER, UE) and acronym “AFTERAM”, MTM2014-52184-P and MTM2014-55966-P. The MRI/DTI data were collected at Johns Hopkins University and the Kennedy-Krieger Institute. We are grateful to Pedro Caro and Iain Currie for useful discussions, to Martin Boer and Cajo ter Braak for the detailed reading of the paper and their many suggestions, and to Bas Engel for sharing with us his knowledge. We are also grateful to the two peer referees for their constructive comments of the paper.

A Proof of Theorem 1

Proof We first note that the first-order partial derivatives of the (approximate) REML log-likelihood function can be expressed as (see, e.g., Rodríguez-Álvarez et al. 2015b)

$$\frac{\partial l}{\partial \sigma_{kl}^2} = -\frac{1}{2} \text{trace} \left(\mathbf{Z}^\top \mathbf{P} \mathbf{Z} \frac{\partial \mathbf{G}}{\partial \sigma_{kl}^2} \right) + \frac{1}{2} \hat{\boldsymbol{\alpha}}^\top \mathbf{G}^{-1} \frac{\partial \mathbf{G}}{\partial \sigma_{kl}^2} \mathbf{G}^{-1} \hat{\boldsymbol{\alpha}}.$$

Given that \mathbf{G} is a positive definite matrix, we have the identity

$$\frac{\partial \mathbf{G}}{\partial \sigma_{kl}^2} = -\mathbf{G} \frac{\partial \mathbf{G}^{-1}}{\partial \sigma_{kl}^2} \mathbf{G},$$

and thus

$$\frac{\partial \mathbf{G}}{\partial \sigma_{kl}^2} = \frac{1}{\sigma_{kl}^4} \text{diag} \left(\mathbf{0}^{(1)}, \mathbf{G}_k \boldsymbol{\Lambda}_{kl} \mathbf{G}_k, \mathbf{0}^{(2)} \right),$$

where $\mathbf{0}^{(1)}$ and $\mathbf{0}^{(2)}$ are zero square matrices of appropriate dimensions.

The first-order partial derivatives of the REML log-likelihood function are then expressed as

$$2 \frac{\partial l}{\partial \sigma_{kl}^2} = -\frac{1}{\sigma_{kl}^4} \text{trace} \left(\mathbf{Z}_k^\top \mathbf{P} \mathbf{Z}_k \mathbf{G}_k \boldsymbol{\Lambda}_{kl} \mathbf{G}_k \right) + \frac{1}{\sigma_{kl}^4} \hat{\boldsymbol{\alpha}}_k^\top \boldsymbol{\Lambda}_{kl} \hat{\boldsymbol{\alpha}}_k.$$

When the REML estimates are positive, they are obtained by equating the former expression to zero (see, e.g., Engel 1990)

$$\frac{\hat{\boldsymbol{\alpha}}_k^\top \boldsymbol{\Lambda}_{kl} \hat{\boldsymbol{\alpha}}_k}{\text{trace} \left(\mathbf{Z}_k^\top \mathbf{P} \mathbf{Z}_k \mathbf{G}_k \boldsymbol{\Lambda}_{kl} \mathbf{G}_k \right)} = 1.$$

We now multiply both sides with σ_{kl}^2 and evaluate the left-hand side for the previous iterates and the right-hand side for the update, obtaining

$$\begin{aligned}\hat{\sigma}_{kl}^2 &= \frac{\hat{\alpha}_k^{[t]\top} \Lambda_{kl} \hat{\alpha}_k^{[t]}}{\text{trace} \left(\mathbf{Z}_k^\top \mathbf{P}^{[t]} \mathbf{Z}_k \mathbf{G}_k^{[t]} \Lambda_{kl} \mathbf{G}_k^{[t]} \right)} \hat{\sigma}_{kl}^{2[t]} \\ &= \frac{\hat{\alpha}_k^{[t]\top} \Lambda_{kl} \hat{\alpha}_k^{[t]}}{\text{trace} \left(\mathbf{Z}_k^\top \mathbf{P}^{[t]} \mathbf{Z}_k \mathbf{G}_k^{[t]} \frac{\Lambda_{kl}}{\hat{\sigma}_{kl}^{2[t]}} \mathbf{G}_k^{[t]} \right)}.\end{aligned}$$

□

B Proof of Theorem 2

Proof First let us recall some notation and introduce some needed results. We denote as $\mathbf{P} = \mathbf{V}^{-1} - \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{V}^{-1}$, where $\mathbf{V} = \mathbf{R} + \mathbf{Z} \mathbf{G} \mathbf{Z}^\top$, $\mathbf{R} = \phi \mathbf{W}^{-1}$ with \mathbf{W} being the diagonal weight matrix involved in the Fisher scoring algorithm.

Denote as $\mathcal{C}(\mathbf{A})$ the linear space spanned by the columns of \mathbf{A} , and let $\mathbf{P}_{\mathbf{XV}^{-1}} = \mathbf{X} (\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{V}^{-1}$ be the orthogonal projection matrix for $\mathcal{C}(\mathbf{X})$ with respect to \mathbf{V}^{-1} . It is easy to show that

$$\begin{aligned}\mathbf{P} &= \mathbf{V}^{-1} (\mathbf{I}_n - \mathbf{P}_{\mathbf{XV}^{-1}}) \\ &= (\mathbf{I}_n - \mathbf{P}_{\mathbf{XV}^{-1}}) \mathbf{V}^{-1} (\mathbf{I}_n - \mathbf{P}_{\mathbf{XV}^{-1}}).\end{aligned}$$

By Theorem 14.2.9 in Harville (1997), \mathbf{P} is a (symmetric) positive semi-definite matrix. In addition,

$$\mathbf{P} \mathbf{X} = \mathbf{0},$$

and

$$\begin{aligned}\text{rank}(\mathbf{P}) &= \text{rank} \left(\mathbf{V}^{-1} (\mathbf{I}_n - \mathbf{P}_{\mathbf{XV}^{-1}}) \right) \\ &= \text{rank} ((\mathbf{I}_n - \mathbf{P}_{\mathbf{XV}^{-1}})) \\ &= n - \text{rank} (\mathbf{P}_{\mathbf{XV}^{-1}}) \\ &= n - \text{rank} (\mathbf{X}).\end{aligned}$$

Thus,

$$\ker(\mathbf{P}) = \mathcal{C}(\mathbf{X}), \quad (26)$$

i.e., $\mathbf{P} \mathbf{x} = \mathbf{0}$ if and only if \mathbf{x} is in $\mathcal{C}(\mathbf{X})$.

Let $\Lambda_{kl} = \mathbf{U} \Sigma \mathbf{U}^\top$ be the eigenvalue decomposition of Λ_{kl} . Note that $\Lambda_{kl} = \mathbf{U}_+ \Sigma_+ \mathbf{U}_+^\top$, where \mathbf{U}_+ and Σ_+ are the sub-matrices corresponding to the nonzero eigenvalues. Then

$$\begin{aligned}\hat{\alpha}_k^\top \Lambda_{kl} \hat{\alpha}_k &= \hat{\alpha}_k^\top \mathbf{U}_+ \Sigma_+ \mathbf{U}_+^\top \hat{\alpha}_k \\ &= \mathbf{z}^\top \mathbf{P} \mathbf{Z}_k \mathbf{G}_k \mathbf{U}_+ \Sigma_+ \mathbf{U}_+^\top \mathbf{G}_k \mathbf{Z}_k^\top \mathbf{P} \mathbf{z} \geq 0,\end{aligned}$$

with equality holding if and only if $\mathbf{U}_+^\top \mathbf{G}_k \mathbf{Z}_k^\top \mathbf{P} \mathbf{z} = \mathbf{0}$ (since Σ_+ is positive definite). Thus, using result (26), the equality

holds if \mathbf{z} is in $\mathcal{C}(\mathbf{X})$ or $\mathcal{C}(\mathbf{Z}_k \mathbf{G}_k \mathbf{U}_+) \subset \mathcal{C}(\mathbf{X})$. By Lemma 4.2.2 and Corollary 4.5.2 in Harville (1997), we have

$$\begin{aligned}\mathcal{C}(\mathbf{Z}_k \mathbf{G}_k \mathbf{U}_+) &= \mathcal{C}(\mathbf{Z}_k \mathbf{G}_k \Lambda_{kl}) \subset \mathcal{C}(\mathbf{X}) \\ &\iff \text{rank}(\mathbf{Z}_k \mathbf{G}_k \Lambda_{kl}, \mathbf{X}) = \text{rank}(\mathbf{X}).\end{aligned}$$

Regarding the denominator of the REML-based estimates updates, we follow a similar reasoning. Using Corollary 14.7.5 (and Theorem 14.2.9) in Harville (1997), we have

$$\begin{aligned}\text{trace} \left(\mathbf{Z}_k^\top \mathbf{P} \mathbf{Z}_k \mathbf{G}_k \Lambda_{kl} \mathbf{G}_k \right) \\ = \text{trace} \left(\mathbf{U}_+^\top \mathbf{G}_k \mathbf{Z}_k^\top \mathbf{P} \mathbf{Z}_k \mathbf{G}_k \mathbf{U}_+ \Sigma_+ \right) \geq 0,\end{aligned}$$

with equality holding if and only if $\mathbf{U}_+^\top \mathbf{G}_k \mathbf{Z}_k^\top \mathbf{P} \mathbf{Z}_k \mathbf{G}_k \mathbf{U}_+ = \mathbf{0}$. Again, this equality holds if and only if $\mathcal{C}(\mathbf{Z}_k \mathbf{G}_k \mathbf{U}_+) \subset \mathcal{C}(\mathbf{X})$ (i.e., $\iff \text{rank}(\mathbf{Z}_k \mathbf{G}_k \Lambda_{kl}, \mathbf{X}) = \text{rank}(\mathbf{X})$). □

C Estimating algorithm

This appendix summarises the steps of the estimating algorithm for model (10) based on the SOP method. Recall that interest lies in estimating model

$$g(\boldsymbol{\mu}) = \mathbf{X} \boldsymbol{\beta} + \mathbf{Z} \boldsymbol{\alpha} = \mathbf{X} \boldsymbol{\beta} + \sum_{k=1}^c \mathbf{Z}_k \boldsymbol{\alpha}_k,$$

where $\mathbf{Z} = [\mathbf{Z}_1, \dots, \mathbf{Z}_c]$, $\boldsymbol{\alpha} = (\boldsymbol{\alpha}_1^\top, \dots, \boldsymbol{\alpha}_c^\top)^\top$, $\boldsymbol{\alpha}_k \sim N(\mathbf{0}, \mathbf{G}_k)$ with $\mathbf{G}_k^{-1} = \sum_{l=1}^{p_k} \sigma_{kl}^{-2} \Lambda_{kl}$, and $\boldsymbol{\alpha} \sim N(\mathbf{0}, \mathbf{G})$ with $\mathbf{G} = \bigoplus_{k=1}^c \mathbf{G}_k$.

Initialise Set initial values for $\hat{\boldsymbol{\mu}}^{[0]}$ and the variance parameters $\hat{\sigma}_{kl}^{2[0]}$ ($l = 1, \dots, p_k$ and $k = 1, \dots, c$). In those situations where ϕ is unknown, establish an initial value for this parameter, $\hat{\phi}^{[0]}$. Set $t = 0$.

Step 1 Construct the *working* response vector \mathbf{z} and the matrix of weights \mathbf{W} as follows

$$\begin{aligned}z_i &= g(\hat{\mu}_i^{[t]}) + (y_i - \hat{\mu}_i^{[t]}) g'(\hat{\mu}_i^{[t]}), \\ w_{ii} &= \left\{ g'(\hat{\mu}_i^{[t]})^2 v(\hat{\mu}_i^{[t]}) \right\}^{-1}.\end{aligned}$$

Step 1.1 Given the initial *estimates* of variance parameters, estimate $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ by solving

$$\underbrace{\begin{bmatrix} \mathbf{X}^\top \mathbf{R}^{[t]-1} \mathbf{X} & \mathbf{X}^\top \mathbf{R}^{[t]-1} \mathbf{Z} \\ \mathbf{Z}^\top \mathbf{R}^{[t]-1} \mathbf{X} & \mathbf{Z}^\top \mathbf{R}^{[t]-1} \mathbf{Z} + \mathbf{G}^{[t]-1} \end{bmatrix}}_{\mathbf{C}^{[t]}} \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\alpha}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}^\top \mathbf{R}^{[t]-1} \mathbf{z} \\ \mathbf{Z}^\top \mathbf{R}^{[t]-1} \mathbf{z} \end{bmatrix}, \quad (27)$$

where $\mathbf{R}^{[t]} = \hat{\phi}^{[t]} \mathbf{W}^{-1}$. Let $\hat{\alpha}^{[t]}$ and $\hat{\beta}^{[t]}$ be these estimates.

Step 1.2 Update the variance parameters as follows

$$\hat{\sigma}_{k_l}^2 = \frac{\hat{\alpha}_k^{[t]\top} \Lambda_{k_l} \hat{\alpha}_k^{[t]}}{\text{ED}_{k_l}^{[t]}},$$

and, when necessary,

$$\hat{\phi} = \frac{(z - X\hat{\beta}^{[t]} - Z\hat{\alpha}^{[t]})^\top \mathbf{W} (z - X\hat{\beta}^{[t]} - Z\hat{\alpha}^{[t]})}{n - \text{rank}(X) - \sum_{k=1}^c \sum_{l=1}^{p_k} \text{ED}_{k_l}^{[t]}},$$

with

$$\text{ED}_{k_l}^{[t]} = \text{trace} \left(\left(\mathbf{G}_k^{[t]} - \mathbf{C}_{kk}^{[t]*} \right) \frac{\Lambda_{k_l}}{\hat{\sigma}_{k_l}^{2[t]}} \right).$$

Recall that $\mathbf{C}^{[t]*}$ denotes the inverse of $\mathbf{C}^{[t]}$ in (27), and $\mathbf{C}_{kk}^{[t]*}$ is that partition of $\mathbf{C}^{[t]*}$ corresponding to the k -th random component α_k .

Step 1.3 Repeat Step 1.1 and Step 1.2, replacing $\hat{\sigma}_{k_l}^{2[t]}$, and, if updated, $\hat{\phi}^{[t]}$, by $\hat{\sigma}_{k_l}^2$ and $\hat{\phi}$, until convergence. For the examples presented in Sect. 5, we use the REML deviance as the convergence criterion.

Step 2 Repeat Step 1, replacing the variance parameters and the model's fixed and random effects (and thus $\hat{\mu}^{[t]}$) by those obtained in the last iteration of Steps 1.1–Step 1.3, until convergence.

It is worth noting that for notational convenience, in the examples described in Sect. 4, the precision matrix was rewritten as $\mathbf{G}^{-1} = \sum_{l=1}^p \sigma_l^{-2} \tilde{\Lambda}_l$, where $p = \sum_{k=1}^c p_k$ is the number of variance parameters, and $\tilde{\Lambda}_l$ are the matrices Λ_{k_l} padded out with zeroes. Here, $\tilde{\Lambda}_l$ are matrices of dimension $q \times q$, where $q = \sum_{k=1}^c q_k$ is the number of random effects coefficients. For this specification, the estimating algorithm discussed above remains essentially the same, but in Step 1.2, $\alpha_k^{[t]}$, $\mathbf{G}_k^{[t]}$, and Λ_{k_l} are replaced by, respectively, $\alpha^{[t]}$, $\mathbf{G}^{[t]}$ and $\tilde{\Lambda}_l$, and $\mathbf{C}_{kk}^{[t]*}$ would be that partition of $\mathbf{C}^{[t]*}$ corresponding to the random vector α (and thus the same for all variance parameters).

D Factor-by-curve hierarchical curve model

This appendix describes in detail the factor-by-curve interaction model discussed in Sect. 5.3, i.e.,

$$y_{ij} = f_{z_j}(t_i) + g_j(t_i) + \varepsilon_{ij} \quad 1 \leq i \leq s, \quad 1 \leq j \leq m,$$

where $z_j = 1$ if the j -th individual is affected by MS (case) and $z_j = 0$ otherwise (control). Let's order the data with the observations on controls first, followed by observations on MS patients. In matrix notation, the model can be expressed as

$$\mathbf{y} = [\mathbf{Q} \otimes \mathbf{B}] \boldsymbol{\theta} + [\mathbf{I}_m \otimes \check{\mathbf{B}}] \check{\boldsymbol{\theta}} + \boldsymbol{\varepsilon}, \quad (28)$$

with \mathbf{B} , $\check{\mathbf{B}}$, $\check{\boldsymbol{\theta}}$ as defined in Sect. 4.2, and $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}_{\cdot,1}^\top, \dots, \boldsymbol{\varepsilon}_{\cdot,m}^\top)^\top$, with $\boldsymbol{\varepsilon}_{\cdot,j} = (\varepsilon_{1j}, \dots, \varepsilon_{sj})^\top$. Matrix \mathbf{Q} is any suitable contrast matrix of dimension $m \times 2$, where $m = m_0 + m_1$, with m_0 being the number of controls and m_1 the number of MS patients. For our application, we consider

$$\mathbf{Q} = \begin{pmatrix} \mathbf{1}_{m_0} & \mathbf{0}_{m_0} \\ \mathbf{0}_{m_1} & \mathbf{1}_{m_1} \end{pmatrix},$$

and a different amount of smoothing is assumed for f_0 and f_1 , i.e., the penalty matrix acting over the vector of coefficients $\boldsymbol{\theta}$ is of the form

$$\mathbf{P} = \begin{pmatrix} \lambda_1 \mathbf{D}_q^\top \mathbf{D}_q & \mathbf{0}_{c \times c} \\ \mathbf{0}_{c \times c} & \lambda_2 \mathbf{D}_q^\top \mathbf{D}_q \end{pmatrix}.$$

The reformulation as a mixed model can be done in a similar fashion to that described in Sect. 4.2, with, in this case

$$\mathbf{X} = [\mathbf{Q} \otimes \mathbf{B} \mathbf{U}_0],$$

$$\mathbf{Z} = [\mathbf{Q} \otimes \mathbf{B} \mathbf{U}_+ : \mathbf{I}_m \otimes \check{\mathbf{B}}],$$

and

$$\mathbf{G}^{-1} = \begin{pmatrix} \sigma_1^{-2} \boldsymbol{\Sigma}_+ & \mathbf{0}_{d \times d} & \mathbf{0}_{d \times (m\check{d})} \\ \mathbf{0}_{d \times d} & \sigma_2^{-2} \boldsymbol{\Sigma}_+ & \mathbf{0}_{d \times (m\check{d})} \\ \mathbf{0}_{(m\check{d}) \times d} & \mathbf{0}_{(m\check{d}) \times c} & \mathbf{I}_m \otimes \check{\mathbf{G}}^{-1} \end{pmatrix},$$

where

$$\check{\mathbf{G}}^{-1} = \sigma_3^{-2} \check{\mathbf{D}}_q^\top \check{\mathbf{D}}_q + \sigma_4^{-2} \mathbf{I}_{\check{d}}.$$

References

- Breslow, N.E., Clayton, D.G.: Approximate inference in generalized linear mixed models. *J. Am. Stat. Assoc.* **88**(421), 9–25 (1993)
- Camarda, C.G., Eilers, P.H., Gampe, J.: Sums of smooth exponentials to decompose complex series of counts. *Stat. Model.* **16**(4), 279–296 (2016)
- Crainiceanu, C.M., Ruppert, D., Carroll, R.J., Joshi, A., Goodner, B.: Spatially adaptive Bayesian penalized splines with heteroscedastic errors. *J. Comput. Graph. Stat.* **16**(2), 265–288 (2007)
- Crump, S.L.: The present status of variance component analysis. *Biometrics* **7**(1), 1–16 (1951)

- Cui, Y., Hodges, J.S., Kong, X., Carlin, B.P.: Partitioning degrees of freedom in hierarchical and other richly-parameterized models. *Technometrics* **52**, 124–136 (2010)
- Currie, I.D., Durban, M.: Flexible smoothing with P-splines: a unified approach. *Stat. Model.* **2**(4), 333–349 (2002)
- Currie, I.D., Durban, M., Eilers, P.H.C.: Generalized linear array models with applications to multidimensional smoothing. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **68**(2), 259–280 (2006)
- Davies, P.L., Gather, U., Meise, M., Mergel, D., Mildenerberger, T.: Residual-based localization and quantification of peaks in X-ray diffractograms. *Ann. Appl. Stat.* **2**(3), 861–886 (2008)
- Davies, P.L., Gather, U., Meise, M., Mergel, D., Mildenerberger, T., Bernholt, T., Hofmeister, T.: diffractometry: baseline identification and peak decomposition for x-ray diffractograms. R package version 0.1-10 (2018)
- Djeundje, V.A., Currie, I.D.: Appropriate covariance-specification via penalties for penalized splines in mixed models for longitudinal data. *Electron. J. Stat.* **4**, 1202–1224 (2010)
- Durban, M., Aguilera-Morillo, M.C.: On the estimation of functional random effects. *Stat. Model.* **17**(1–2), 50–58 (2017)
- Durban, M., Harezlak, J., Wand, M.P., Carroll, R.J.: Simple fitting of subject-specific curves for longitudinal data. *Stat. Med.* **24**(8), 1153–1167 (2005)
- Eilers, P.H.C.: Discussion of Verbyla et al. *J. R. Stat. Soc. Ser. C (Appl. Stat.)* **48**, 300–311 (1999)
- Eilers, P.H.C., Marx, B.D.: Flexible smoothing with B-splines and penalties. *Stat. Sci.* **11**(2), 89–121 (1996)
- Engel, B.: The analysis of unbalanced linear models with variance components. *Stat. Neerl.* **44**, 195–219 (1990)
- Engel, B., Buist, W.: Analysis of a generalized linear mixed model: a case study and simulation results. *Biom. J.* **38**(1), 61–80 (1996)
- Engel, B., Keen, A.: A simple approach for the analysis of generalized linear mixed models. *Stat. Neerl.* **48**(1), 1–22 (1994)
- Fan, J., Li, R.: Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Am. Stat. Assoc.* **96**(456), 1348–1360 (2001)
- Gilmour, A.R., Thompson, R., Cullis, B.R.: Average information REML: an efficient algorithm for variance parameter estimation in linear mixed models. *Biometrics* **51**(4), 1440–1450 (1995)
- Goldsmith, J., Bobb, J., Crainiceanu, C.M., Caffo, B., Reich, D.: Penalized functional regression. *J. Comput. Graph. Stat.* **20**(4), 830–851 (2011)
- Goldsmith, J., Crainiceanu, C.M., Caffo, B., Reich, D.: Longitudinal penalized functional regression for cognitive outcomes on neuronal tract measurements. *J. R. Stat. Soc. Ser. C (Appl. Stat.)* **61**(3), 453–469 (2012)
- Goldsmith, J., Scheipl, F., Huang, L., Wrobel, J., Gellar, J., Harezlak, J., McLean, M.W., Swihart, B., Xiao, L., Crainiceanu, C., Reiss, P.T.: refund: Regression with Functional Data. R package version 0.1-16 (2016)
- Graser, H.-U., Smith, S.P., Tier, B.: A derivative-free approach for estimating variance components in animal models by restricted maximum likelihood. *J. Anim. Sci.* **2**(64), 1362–1373 (1987)
- Green, P.J.: Penalized likelihood for general semi-parametric regression models. *Int. Stat. Rev./Revue Internationale de Statistique* **55**(3), 245–259 (1987)
- Greven, S., Scheipl, F.: A general framework for functional regression modelling. *Stat. Model.* **17**(1–2), 1–35 (2017)
- Groll, A., Tutz, G.: Variable selection for generalized linear mixed models by L1-penalized estimation. *Stat. Comput.* **24**(2), 137–154 (2014)
- Harville, D.A.: Maximum likelihood approaches to variance component estimation and to related problems. *J. Am. Stat. Assoc.* **72**(358), 320–338 (1977)
- Harville, D.A.: Matrix Algebra from a Statistician’s Perspective. Springer, Berlin (1997)
- Hastie, T.J., Tibshirani, R.J.: Generalized Additive Models. Chapman & Hall, London (1990)
- Heckman, N., Lockhart, R., Nielsen, J.D.: Penalized regression, mixed effects models and appropriate modelling. *Electron. J. Stat.* **7**, 1517–1552 (2013)
- Henderson, C.R.: Selection index and expected genetic advance. *Stat. Genet. Plant Breed.* **982**, 141–163 (1963)
- Hunter, D.R., Li, R.: Variable selection using MM algorithms. *Ann. Stat.* **33**(4), 1617–1642 (2005)
- Johnson, D.L., Thompson, R.: Restricted maximum likelihood estimation of variance components for univariate animal models using sparse matrix techniques and average information. *J. Dairy Sci.* **78**, 449–456 (1995)
- Karas, M., Brzyski, D., Dziedzic, M., Goñi, J., Kareken, D.A., Randolph, T.W., Harezlak, J.: Brain connectivity-informed regularization methods for regression. *Stat. Biosci.* (2017). <https://doi.org/10.1007/s12561-017-9208-x>
- Krivobokova, T.: Smoothing parameter selection in two frameworks for penalized splines. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **75**(4), 725–741. <https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/rssb.12010> (2009)
- Krivobokova, T., Crainiceanu, C.M., Kauermann, G.: Fast adaptive penalized splines. *J. Comput. Graph. Stat.* **17**(1), 1–20 (2008)
- Lee, D.-J.: Smoothing mixed model for spatial and spatio-temporal data. PhD thesis. Department of Statistics, Universidad Carlos III de Madrid, Spain (2010)
- McCullagh, P., Nelder, J.: Generalized Linear Models. Chapman and Hall/CRC Monographs on Statistics and Applied Probability Series, 2nd edn. Chapman & Hall, London (1989)
- Patterson, H.D., Thompson, R.: Recovery of inter-block information when block sizes are unequal. *Biometrika* **58**(3), 545–554 (1971)
- R Core Team: R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria (2018)
- Reiss, P.T., Ogden, R.T.: Smoothing parameter selection for a class of semiparametric linear models. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **71**(2), 505–523 (2009)
- Rodríguez-Álvarez, M.X., Durban, M., Lee, D.-J., Eilers, P.H.C.: Fast estimation of multidimensional adaptive p-spline models. In: Friedl, H., Wagner, H. (eds.) Proceedings of the 30th International Workshop on Statistical Modelling, pp 330 – 335. [arXiv:1610.06861](https://arxiv.org/abs/1610.06861) (2015a)
- Rodríguez-Álvarez, M.X., Lee, D.-J., Kneib, T., Durban, M., Eilers, P.H.C.: Fast smoothing parameter separation in multidimensional generalized P-splines: the sap algorithm. *Stat. Comput.* **25**, 941–957 (2015b)
- Rodríguez-Álvarez, M.X., Durban, M., Lee, D.-J., Eilers, P.H.C., Gonzalez, F.: Spatio-temporal adaptive penalized splines with application to neuroscience. In: Dupuy, J.-F., Josse, J. (eds.) Proceedings of the 31th International Workshop on Statistical Modelling, pp 267–272. [arXiv:1610.06860](https://arxiv.org/abs/1610.06860) (2016)
- Rodríguez-Álvarez, M.X., Boer, M.P., van Eeuwijk, F.A., Eilers, P.H.: Correcting for spatial heterogeneity in plant breeding experiments with p-splines. *Spat. Stat.* **23**, 52–71 (2018)
- Ruppert, D., Carroll, R.J.: Spatially-adaptive penalties for spline fitting. *Aust. N. Z. J. Stat.* **42**(2), 205–223 (2000)
- Ruppert, D., Wand, M.P., Carroll, R.: Semiparametric Regression. Cambridge University Press, Cambridge (2003)
- Schall, R.: Estimation in generalized linear models with random effects. *Biometrika* **78**(4), 719–727 (1991)
- Simpkin, A., Newell, J.: An additive penalty p-spline approach to derivative estimation. *Comput. Stat. Data Anal.* **68**, 30–43 (2013)
- Smith S.P.: Estimation of genetic parameters in non-linear models. In: Gianola, D., Hammond, K. (eds.) Advances in Statistical Methods for Genetic Improvement of Livestock. Advanced Series in Agricultural Sciences, vol. 18. Springer, Berlin, Heidelberg (1990)

- Taylor, J.D., Verbyla, A.P., Cavanagh, C., Newberry, M.: Variable selection in linear mixed models using an extended class of penalties. *Aust. N. Z. J. Stat.* **54**(4), 427–449 (2012)
- Tibshirani, R.J.: Adaptive piecewise polynomial estimation via trend filtering. *Ann. Stat.* **42**(1), 285–323 (2014)
- Wand, M.P.: Smoothing and mixed models. *Comput. Stat.* **18**(2), 223–249 (2003)
- Wood, S.N.: Fast stable direct fitting and smoothness selection for generalized additive models. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **70**(3), 495–518 (2008)
- Wood, S.N.: Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **73**(1), 2–36 (2011)
- Wood, S.N.: *Generalized Additive Models: An Introduction with R*, 2nd edn. Chapman & Hall CRC, London (2017)
- Wood, S.N., Fasiolo, M.: A generalized Fellner-Schall method for smoothing parameter optimization with application to Tweedie location, scale and shape models. *Biometrics* **73**, 1071–1081 (2017)
- Wood, S.N., Pya, N., Säfken, B.: Smoothing parameter and model selection for general smooth models. *J. Am. Stat. Assoc.* **111**(516), 1548–1563 (2016)
- Zou, H., Li, R.: One-step sparse estimates in nonconcave penalized likelihood models. *Ann. Stat.* **36**(4), 1509–1533 (2008)