



Published in final edited form as:

*Biometrics*. 2017 December ; 73(4): 1092–1101. doi:10.1111/biom.12697.

## Pointwise Influence Matrices for Functional-Response Regression

Philip T. Reiss<sup>1,2,3</sup>, Lei Huang<sup>4</sup>, Pei-Shien Wu<sup>1</sup>, Huaihou Chen<sup>5</sup>, and Stan Colcombe<sup>6</sup>

<sup>1</sup>Department of Child & Adolescent Psychiatry, New York University School of Medicine, USA

<sup>2</sup>Department of Population Health, New York University School of Medicine, USA

<sup>3</sup>Department of Statistics, University of Haifa, Israel

<sup>4</sup>Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, USA

<sup>5</sup>Department of Biostatistics, University of Florida, USA

<sup>6</sup>Nathan S. Kline Institute for Psychiatric Research, USA

### SUMMARY

We extend the notion of an influence or hat matrix to regression with functional responses and scalar predictors. For responses depending linearly on a set of predictors, our definition is shown to reduce to the conventional influence matrix for linear models. The pointwise degrees of freedom, the trace of the pointwise influence matrix, are shown to have an adaptivity property that motivates a two-step bivariate smoother for modeling nonlinear dependence on a single predictor. This procedure adapts to varying complexity of the nonlinear model at different locations along the function, and thereby achieves better performance than competing tensor product smoothers in an analysis of the development of white matter microstructure in the brain.

### Keywords

Bivariate smoothing; Degrees of freedom; Fractional anisotropy; Function-on-scalar regression; Functional nonlinear regression; Neurodevelopmental trajectory; Tensor product spline

### 1. Introduction

The influence or hat matrix is a central concept in parametric regression (Hoaglin and Welsch, 1978), as well as in semiparametric regression (Ruppert et al., 2003), for which the term *smoother matrix* is often used. Given the growing importance of functional data analysis (Ramsay and Silverman, 2005), there is a need to extend the idea of a influence matrix to regression with functional responses. In this article we define influence matrices in a pointwise sense for functional-response regression, and to show that this definition, and a

Supplementary Materials

Web Appendix A, referenced in Sections 4, 5.2 and 7, and Web Appendix B, referenced in Section 5.1, are available with this paper at the *Biometrics* website on Wiley Online Library, along with the data and R code for the corpus callosum analyses.

related notion of pointwise degrees of freedom, can be useful tools for studying the performance of functional regression estimators.

Letting  $y(s)$  denote a real-valued functional response, defined for  $s$  in an interval  $\mathcal{S} \in \mathbb{R}$ , we assume throughout that

$$E[y(s) | \mathbf{x}] = \mathbf{x}^T \boldsymbol{\beta}(s), \quad (1)$$

where  $\mathbf{x} \in \mathbb{R}^p$  is a predictor vector and  $\boldsymbol{\beta}(s)$  is a vector-valued function on  $\mathcal{S}$ . This formulation will allow a unified presentation of linear and nonlinear dependence on scalar predictors, or more specifically, of the following two special cases of (1).

### Problem 1: Linear dependence on several predictors

If  $\mathbf{x} = (1, x_1, \dots, x_{p-1})^T$ , where  $x_1, \dots, x_{p-1}$  are scalar predictors, we have the varying-coefficient model

$$E[y(s) | x_1, \dots, x_{p-1}] = \beta_0(s) + \sum_{j=1}^{p-1} x_j \beta_j(s) \quad (2)$$

where for each  $j$ ,  $\beta_j(\cdot)$  is a smooth function defined on  $\mathcal{S}$ . Model (2) is often called a functional linear model (e.g., Ramsay and Silverman, 2005); Reiss et al. (2010) coined the term *function-on-scalar* regression to distinguish (2) from linear models with functional predictors.

### Problem 2: Nonlinear dependence on a single predictor

Let  $p = K_t$  and  $\mathbf{x} = \mathbf{b}_t(t)$ , where  $t \in \mathbb{R}$  is a predictor variable and  $\mathbf{b}_t(t) = [b_{t1}(t), \dots, b_{tK_t}(t)]^T$  is a vector of smooth basis functions defined on the domain of  $t$ ; our main example will be a  $B$ -spline basis that is sufficiently rich to have small approximation error. Here the components  $\beta_1(s), \dots, \beta_{K_t}(s)$  of  $\boldsymbol{\beta}(s)$  are again assumed to be smooth, but may not be individually of primary interest; instead, we are interested in estimating the bivariate smooth function

$$f(t, s) = E[y(s) | t] = \mathbf{b}_t(t)^T \boldsymbol{\beta}(s). \quad (3)$$

As suggested by the arguments  $t$  and  $s$ , in some settings (such as our motivating application), the predictor refers to a point in time, while the argument of the functional response denotes a location in space. Thus we sometimes refer to the temporal or spatial dimension in what follows.

For either of the above cases we shall assume that each component of  $\boldsymbol{\beta}(s)$  can be represented with negligible approximation error as a linear combination of basis functions such as  $B$ -splines, so that (1) becomes  $E[y(s) | \mathbf{x}] = \mathbf{x}^T \boldsymbol{\Theta} \mathbf{b}_s(s)$  where  $\mathbf{b}_s(s) = [b_{s1}(s), \dots, b_{sK_s}(s)]^T$

is a basis function vector and  $\Theta$  is a  $p \times K_s$  coefficient matrix. Consequently our estimators for model (1) will have the form

$$\hat{E}[y(s) | x] = x^T \hat{\Theta} b_s(s). \quad (4)$$

The present work was motivated by Problem 2 in the context of human brain development studies, in which the family of functions  $\{f(\cdot, s) : s \in \mathcal{S}\}$  may be of particular interest. When  $t$  represents age,  $f(\cdot, s)$  is the mean, as a function of age, of the quantity measured by  $y$  at point  $s$ . In the specific example that will be presented below,  $\mathcal{S}$  is a set of locations in the brain, and  $y(s)$  denotes fractional anisotropy, a measure of white matter integrity, at location  $s$ . Thus  $f(t, s)$  denotes the mean fractional anisotropy for that location at age  $t$ , and the function  $f(\cdot, s)$  is what neuroscientists often refer to as the developmental trajectory of fractional anisotropy at location  $s$ . As a brief example of the scientific meaning of such trajectories, suppose that for given  $s$ ,  $f(t, s)$  characteristically increases with  $t$  up to some point  $t_s \equiv \operatorname{argmax}_t f(t, s)$ , then decreases. Then the peak age  $t_s$  can provide information about typical maturation for location  $s$ , and can be compared between diagnostic groups to study the links between psychiatric disorders and brain development (Shaw et al., 2007).

In our data set, fractional anisotropy was measured at 107 voxels, or  $1 \times 1 \times 1$  mm volume units, in 146 individuals age 7–48. These voxels, based on registration to a standard image from the FMRIB Software Library, trace a path along a midsagittal cross-section of the corpus callosum (see Fig. 1). We take  $s$  to represent arc length along this path, which ranges from  $s_1 = 0$  mm (the leftmost point in the figure, toward the back of the brain) to  $s_{107} = 110.55$  mm. At right in Fig. 1, a rainbow plot (Hyndman and Shang, 2010), with fractional anisotropy curves color-coded by age, is used to visualize the relationship between age and the functional response. This relationship is not easily discerned from the plot, and may be non-monotonic (and hence nonlinear) in some locations.

In ordinary nonparametric estimation of a function  $g(t) = E(y|t)$  based on data pairs  $(t_1, y_1), \dots, (t_n, y_n)$ , an important role is played by the  $n \times n$  influence or smoother matrix  $H$  such that the response vector  $y = (y_1, \dots, y_n)^T$  is related to a fitted value vector

$$\hat{y} = (\hat{y}_1, \dots, \hat{y}_n)^T = [\hat{g}(t_1), \dots, \hat{g}(t_n)]^T \text{ by} \\ \hat{y} = Hy. \quad (5)$$

(Since its effect is to add a “hat” to  $y$ ,  $H$  is also known as the *hat matrix*.) In particular, the trace of  $H$  serves as a popular definition of the effective degrees of freedom, a measure of the complexity of a smoother (Wahba, 1983; Buja et al., 1989). When trajectory estimates  $t \mapsto \hat{f}(t, s)$  are of interest, it would be very useful to extend these concepts to the curve estimates  $\hat{f}(\cdot, s)$  for each point  $s \in \mathcal{S}$  at which the functional responses are observed. However, for estimators of the type considered here, which share information among points  $s \in \mathcal{S}$ , it is not clear how to extend the definition of  $H$  to a *pointwise* influence matrix for a given  $s \in \mathcal{S}$ .

In this paper we develop a notion of pointwise influence matrices for estimators of the form (4), and demonstrate its conceptual and methodological utility. After introducing some needed notation and assumptions in Section 2, we define the pointwise influence matrix, and a corresponding notion of pointwise effective degrees of freedom, in Section 3. In Section 4 we provide some support for our definition by showing that, for Problem 1 (the linear case), it reduces to the influence matrix as ordinarily defined. Section 5 presents a result about pointwise influence matrices for “post-smoothed” estimators which suggests that a simple but apparently novel method should work well for the motivating class of applications; and this is corroborated by the brain imaging data analysis in Section 6. Section 7 offers some concluding remarks.

## 2. Further notation and assumptions

In what follows,  $\mathbf{1}_m$  denotes a vector of  $m$  1's, while  $\mathbf{e}_{r,m}$  is a vector in  $\mathbb{R}^m$  with 1 in the  $r$ th position and 0 elsewhere. We assume independent data pairs  $[\mathbf{x}_i, y_i(\cdot)]$ ,  $i = 1, \dots, n$ , where  $n > p$ , with each of the functional responses  $y_i(\cdot)$  observed at a common set of points  $s_1, \dots, s_L$ , giving rise to an  $n \times L$  response matrix

$$\mathbf{Y} = \begin{pmatrix} y_{11} & \cdots & y_{1L} \\ \vdots & \ddots & \vdots \\ y_{n1} & \cdots & y_{nL} \end{pmatrix} = \begin{pmatrix} y_1(s_1) & \cdots & y_1(s_L) \\ \vdots & \ddots & \vdots \\ y_n(s_1) & \cdots & y_n(s_L) \end{pmatrix}.$$

Irregularly sampled functions can be accommodated by adding a presmoothing step (cf.

Chiou et al., 2003). Let  $\mathbf{y}_i^T$  and  $\mathbf{y}_\ell$  denote the  $i$ th row and  $\ell$ th column of  $\mathbf{Y}$ , respectively, and let  $\mathbf{y} = \text{vec}(\mathbf{Y}) \in \mathbb{R}^{nL}$ .

An estimator of the general type considered here, i.e., having the form (4), yields an  $n \times L$  fitted value matrix

$$\hat{\mathbf{Y}} = \mathbf{X} \hat{\boldsymbol{\Theta}} \mathbf{B}_s^T, \quad (6)$$

where  $\mathbf{X}$  is the matrix with  $i$ th row  $\mathbf{x}_i^T$  ( $i = 1, \dots, n$ ) and  $\mathbf{B}_s$  is the  $L \times K_s$  matrix with  $\ell$ th row  $[b_{s1}(s_\ell), \dots, b_{sK_s}(s_\ell)]$ . A standard result on the vec operator and Kronecker products yields

$$\hat{\mathbf{y}} = (\mathbf{B}_s \otimes \mathbf{X}) \hat{\boldsymbol{\theta}} \quad (7)$$

where  $\hat{\mathbf{y}} = \text{vec}(\hat{\mathbf{Y}})$  and  $\hat{\boldsymbol{\theta}} = \text{vec}(\hat{\boldsymbol{\Theta}})$ .

When the above  $K_s$  basis functions are  $B$ -splines, these are typically used in conjunction with a roughness penalty (Green and Silverman, 1994), which is implemented by means of a symmetric positive semidefinite  $K_s \times K_s$  matrix  $\mathbf{P}_s$  such that, for a given function  $g(s) = \boldsymbol{\gamma}^T \mathbf{b}_s(s)$ , we have  $\boldsymbol{\gamma}^T \mathbf{P}_s \boldsymbol{\gamma} = r_s(g)$  where  $r_s(g)$  is some measure of the roughness of  $g$ . Popular choices include the second-derivative penalty matrix  $\mathbf{P}_s = [\int b_{si}''(s) b_{sj}''(s) ds]_{i,j=1, \dots, K_s}$  for

which  $\gamma^T P_s \gamma = r_s(g) \equiv \int_{\mathcal{S}} g''(s)^2 ds$ , or difference penalties as in Eilers and Marx (1996).

We can then define the  $L \times L$  smoother matrix  $H_s^{\lambda_s} = B_s (B_s^T B_s + \lambda_s P_s)^{-1} B_s^T$ , where  $\lambda_s \geq 0$  is the tuning parameter that governs the extent of the smoothing. For Problem 2, we have  $X = B_b$  defined analogously to  $B_s$  but with temporal-dimension smooth basis functions  $b_b(t)$  as in (3); we likewise define the  $K_t \times K_t$  penalty matrix  $P_b$  and the  $n \times n$  smoother matrix  $H_t^{\lambda_t} = B_t (B_t^T B_t + \lambda_t P_t)^{-1} B_t^T$ .

We shall require two assumptions regarding smoothing in the spatial direction:

$$B_s \mathbf{1}_{K_s} = \mathbf{1}_L, \quad (8)$$

$$P_s \mathbf{1}_{K_s} = 0. \quad (9)$$

In particular, assumption (8) holds for a  $B$ -spline basis in one dimension (Schumaker, 2007), while (9) holds for a derivative or difference penalty (Eilers and Marx, 1996).

### 3. The pointwise influence matrix

All of the estimation procedures presented below (in Sections 4 and 5) produce fitted values that can be written as  $\hat{y} = \mathcal{H} y$  for some  $nL \times nL$  matrix

$$\mathcal{H} = \begin{pmatrix} \mathcal{H}_{11} & \cdots & \mathcal{H}_{1L} \\ \vdots & \ddots & \vdots \\ \mathcal{H}_{L1} & \cdots & \mathcal{H}_{LL} \end{pmatrix}, \quad (10)$$

where each of the blocks  $\mathcal{H}_{\ell_1 \ell_2}$  is  $n \times n$ . Moreover, for all of the procedures we consider one can write  $\hat{\theta} = C y$ , for some  $pK_s \times nL$  matrix  $C$ , in (7), so the influence matrix has the form

$$\mathcal{H} = (B_s \otimes X) C. \quad (11)$$

While an explicit estimate  $\hat{\theta}$  is required to estimate the conditional mean of  $y(s)$  for arbitrary  $s$ , some procedures are more conveniently described by merely specifying  $\hat{y}$ .

Some intuition regarding the block influence matrix (10) can be gained from Fig. 2. Subfigure (a) displays a portion of  $\mathcal{H}$  obtained by fitting a tensor product smooth (as in Section 5) to a subset of the corpus callosum data. Had we fitted separate models at each voxel, the nonzero entries in the influence matrix, representing influence of the responses on the fitted values, would be confined to diagonal blocks such as those outlined in black. The sharing of information across locations is expressed as a blockwise blurring in the horizontal

direction. Formally, whereas with separate models the fitted values at the  $\ell$ th location would be  $\hat{\mathbf{y}}_{\cdot\ell} = \mathcal{H}_{\ell\ell} \mathbf{y}_{\cdot\ell}$ , here we have

$$\hat{\mathbf{y}}_{\cdot\ell} = \sum_{\ell^*=1}^L \mathcal{H}_{\ell\ell^*} \mathbf{y}_{\cdot\ell^*}, \quad (12)$$

implying that the fitted values are influenced not just by the  $\ell$ th location data via  $\mathcal{H}_{\ell\ell}$ , but by all of the response data via  $\mathcal{H}_{\ell 1}, \dots, \mathcal{H}_{\ell L}$ . This motivates the following two definitions.

### Definition 1

The *pointwise influence matrix* or *pointwise hat matrix* at location  $\ell$  is the  $n \times n$  matrix

$$\mathcal{H}_{\ell\cdot} = \sum_{\ell^*=1}^L \mathcal{H}_{\ell\ell^*} \quad (13)$$

$$= (\mathbf{e}_{\ell:L}^T \otimes \mathbf{I}_n) \mathcal{H} (\mathbf{1}_L \otimes \mathbf{I}_n). \quad (14)$$

### Definition 2

The *pointwise effective degrees of freedom*, or simply *pointwise degrees of freedom*, at location  $\ell$  is

$$\text{DF}_{\ell} = \text{tr}(\mathcal{H}_{\ell\cdot}) = \sum_{\ell^*=1}^L \text{tr}(\mathcal{H}_{\ell\ell^*}) = \sum_{i=1}^n \sum_{\ell^*=1}^L \frac{\partial \hat{y}_{i\ell}}{\partial y_{i\ell^*}}.$$

As an aid to intuition regarding Definition 2, consider a toy example with  $n = 5$  observations and  $L = 4$  locations, so that the overall influence matrix (10) comprises a  $4 \times 4$  grid of  $5 \times 5$  blocks. The conventional degrees of freedom for the 2nd-location model, which would be appropriate if separate models were fitted at each location, is the sum of the shaded values in Fig. 2(b). The pointwise degrees of freedom  $\text{DF}_2$ , which takes into account the influence of neighboring locations, is the sum of the shaded values in Fig. 2(c).

A formal justification for Definition 1 is provided below by Theorem 1, but for now we offer a heuristic interpretation that is valid when  $\hat{\mathbf{y}}_{\cdot\ell}$  lies in the column space of  $\mathcal{H}_{\ell\cdot}$ . On that assumption we can rewrite (12) in a form resembling (5), namely as

$$\hat{\mathbf{y}}_{\cdot\ell} = \mathcal{H}_{\ell\cdot} \mathbf{y}_{\cdot\ell}^{\dagger} \quad (15)$$

where

$$\mathbf{y}_{\cdot\ell}^\dagger = \left( \sum_{\ell^*=1}^L \mathcal{H}_{\ell\ell^*} \right)^- \sum_{\ell^*=1}^L \mathcal{H}_{\ell\ell^*} \mathbf{y}_{\cdot\ell^*}, \quad (16)$$

with  $(\cdot)^-$  denoting a generalized inverse. The  $\ell$ -location fitted value vector  $\hat{\mathbf{y}}_{\cdot\ell}$  can thus be viewed as resulting from a two-step process: first, in (16), we share information among the  $\mathbf{y}_1, \dots, \mathbf{y}_L$  by forming their weighted average  $\mathbf{y}_{\cdot\ell}^\dagger$  with matrix-valued weights  $\mathcal{H}_{\ell 1}, \dots, \mathcal{H}_{\ell L}$ ; then, in (15), we premultiply by the pointwise influence matrix to obtain  $\hat{\mathbf{y}}_{\cdot\ell}$ .

#### 4. Application to varying-coefficient models

In this section we study the pointwise influence matrix in the context of Problem 1, i.e., the varying coefficient linear model (2). We show that for two natural approaches to estimating this model, the pointwise influence matrix reduces to the usual hat matrix for linear regression. We refer to these two approaches as *contemporaneous smoothing* and *post-smoothing*. In contemporaneous smoothing one minimizes a penalized criterion and thereby simultaneously estimates the linear relationship at each point  $s$  and maintains smoothness across points, as opposed to the post-smoothing approach of first fitting separate initial models at each point and then smoothing the results. Post-smoothing can be more computationally efficient than contemporaneous smoothing, but its favorable performance relies on being able to fit the initial models with a modicum of accuracy.

A general expression for the contemporaneous-smoothing coefficient estimate is

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \left[ \{ \mathbf{y} - (\mathbf{B}_s \otimes \mathbf{X}) \boldsymbol{\theta} \}^T \left( \hat{\boldsymbol{\Sigma}}^{-1} \otimes \mathbf{I}_n \right) \{ \mathbf{y} - (\mathbf{B}_s \otimes \mathbf{X}) \boldsymbol{\theta} \} + \boldsymbol{\theta}^T (\mathbf{P}_s \otimes \boldsymbol{\Lambda}_s) \boldsymbol{\theta} \right], \quad (17)$$

where  $\hat{\boldsymbol{\Sigma}}$  is an  $L \times L$  matrix and  $\boldsymbol{\Lambda}_s = \text{Diag}(\lambda_{s,1}, \dots, \lambda_{s,p})$ ; more transparently, the penalty  $\boldsymbol{\theta}^T (\mathbf{P}_s \otimes \boldsymbol{\Lambda}_s) \boldsymbol{\theta}$  in (17) can be written as  $\sum_{j=1}^p \lambda_{s,j} \boldsymbol{\theta}_j^T \mathbf{P}_s \boldsymbol{\theta}_j$  where  $\boldsymbol{\theta}_j^T$  is the  $j$ th row of  $\boldsymbol{\Theta}$ , i.e., as the sum of separate penalties for the  $p$  coefficient functions  $\beta_k(s) = \boldsymbol{\theta}_k^T \mathbf{b}_s(s)$ ,  $k = 1, \dots, p$ .

If we take  $\hat{\boldsymbol{\Sigma}} = \mathbf{I}_L$  then the first term on the right side of (17) reduces to the squared Frobenius norm  $\| \mathbf{Y} - \mathbf{X} \boldsymbol{\Theta} \mathbf{B}_s^T \|_F^2$ , and  $\hat{\boldsymbol{\theta}}$  is a penalized ordinary least squares (OLS) estimate as in Chapter 13 of Ramsay and Silverman (2005). On the other hand, if  $\hat{\boldsymbol{\Sigma}}^{-1}$  is an estimate of the inverse of the covariance matrix  $\boldsymbol{\Sigma} = [\text{Cov}\{\mathbf{y}(s), \mathbf{y}(s')\} | \mathbf{x}]$ ,  $\ell, \ell' = 1, \dots, L$ , typically derived from the residuals of an OLS fit, then  $\hat{\boldsymbol{\theta}}$  is a penalized generalized least squares (GLS) estimate (Reiss et al., 2010).

A post-smoothing estimator via local polynomial smoothing is studied by Fan and Zhang (2000) (see also Zhu et al., 2014). An analogous penalized approach proceeds as follows:

1. Obtain a  $p \times L$  matrix of initial pointwise least squares estimates

$$\tilde{\beta} = \begin{pmatrix} \tilde{\beta}_1(s_1) & \cdots & \tilde{\beta}_1(s_L) \\ \vdots & \ddots & \vdots \\ \tilde{\beta}_p(s_1) & \cdots & \tilde{\beta}_p(s_L) \end{pmatrix} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

2. For  $j = 1, \dots, p$ , obtain final estimates  $\hat{\beta}_j(s_1), \dots, \hat{\beta}_j(s_L)$  by applying a smoother matrix  $\mathbf{H}_s^{\lambda_{s,j}} = \mathbf{B}_s (\mathbf{B}_s^T \mathbf{B}_s + \lambda_{s,j} \mathbf{P}_s)^{-1} \mathbf{B}_s^T$  to the  $j$ th row of  $\tilde{\beta}$ .

This heuristic description suggests the more formal specification

$$\hat{\theta}_{j\cdot}^T = \mathbf{e}_{j:p}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \mathbf{B}_s (\mathbf{B}_s^T \mathbf{B}_s + \lambda_{s,j} \mathbf{P}_s)^{-1} \quad (18)$$

for the  $j$ th row of the  $p \times K_s$  coefficient estimate matrix  $\hat{\Theta}$ .

We can now state the following result, which is proved in Web Appendix A.

### Theorem 1

Assume that  $\mathbf{X}$  is of rank  $p$  and that (8) and (9) hold. For either the contemporaneous smoothing estimator (17) or the post-smoothing estimator (18), the pointwise influence matrix (13) equals  $\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  for  $\ell=1, \dots, L$ .

In short, Theorem 1 says that the pointwise influence matrix for the varying-coefficient model (2) directly generalizes the usual influence matrix for linear regression. This property serves as a form of validation for Definition 1; indeed, we can conceive of no other property that would more cogently validate the definition.

## 5. Post-smoothing and nonlinear trajectories

### 5.1 Pointwise influence matrices for a class of post-smoothing procedures

In the above post-smoothed estimator for model (2) (Problem 1), if  $\lambda_{s,1} = \dots = \lambda_{s,p} = \lambda_s$ , i.e., a common spatial smoother matrix  $\mathbf{H}_s^{\lambda_s}$  is used for each of the  $p$  coefficients, then the final fitted value matrix is

$$\hat{\mathbf{Y}} = \tilde{\mathbf{Y}} \mathbf{H}_s^{\lambda_s} \quad (19)$$

where  $\tilde{\mathbf{Y}}$  is the initial fitted value matrix  $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ . More generally, consider the class of post-smoothed model fits of the form (19), i.e., those obtainable by forming an initial fit  $\tilde{\mathbf{Y}}$  and then postmultiplying by a spatial smoother matrix  $\mathbf{H}_s^{\lambda_s}$ . For example, the fast sandwich smoother

$$\hat{\mathbf{Y}} = \mathbf{H}_t^{\lambda_t} \mathbf{Y} \mathbf{H}_s^{\lambda_s} \quad (20)$$



(Xiao et al., 2013) is of this form, with  $\tilde{Y} = H_t^{\lambda_t} Y$ . Given a fitted value matrix (19) such that  $\tilde{y} = \text{vec}(\tilde{Y}) = \mathcal{H} y$  for some matrix  $\mathcal{H}$ , we can define the initial pointwise influence matrix  $\tilde{\mathcal{H}}_\ell$  analogously to Definition 1, but with  $\mathcal{H}$  replacing  $\mathcal{H}$  in (13), as well as the initial pointwise degrees of freedom  $\tilde{\text{DF}}_\ell = \text{tr}(\tilde{\mathcal{H}}_\ell)$ . We then have the following result, which is proved in Web Appendix B.

**Theorem 2**—Using the subscript  $ij$  to denote the  $(i, j)$  entry of a matrix, we have

$$\begin{pmatrix} \mathcal{H}_{1,ij} \\ \vdots \\ \mathcal{H}_{L,ij} \end{pmatrix} = H_s^{\lambda_s} \begin{pmatrix} \tilde{\mathcal{H}}_{1,ij} \\ \vdots \\ \tilde{\mathcal{H}}_{L,ij} \end{pmatrix}, \quad (21)$$

for  $i, j = 1, \dots, n$ . Moreover,

$$\begin{pmatrix} \text{DF}_1 \\ \vdots \\ \text{DF}_L \end{pmatrix} = H_s^{\lambda_s} \begin{pmatrix} \tilde{\text{DF}}_1 \\ \vdots \\ \tilde{\text{DF}}_L \end{pmatrix}. \quad (22)$$

Theorem 2 says, first, that the pointwise influence matrices  $\mathcal{H}_1, \dots, \mathcal{H}_L$  can be obtained by elementwise smoothing of the initial pointwise influence matrices  $\tilde{\mathcal{H}}_1, \dots, \tilde{\mathcal{H}}_L$  with smoother matrix  $H_s^{\lambda_s}$ ; and second, that applying that smoother to the initial pointwise degrees of freedom vector yields the final pointwise degrees of freedom vector.

## 5.2 An adaptive post-smoothing procedure for nonlinear trajectories

Theorem 2 will turn out to have interesting implications for Problem 2, estimating the bivariate smooth function  $f(t, s)$  (3) given a set of functional responses associated with a scalar predictor. Problem 2 admits of the same two types of approaches to spatial smoothness that we described above for Problem 1 (see Table 1). A *contemporaneous smoothing* approach is the tensor product smoother  $\hat{f}(t, s) = b_t(t)^T \hat{\Theta} b_s(s)$  with

$$\hat{\Theta} = \arg \min_{\Theta} \left[ \| (Y - B_t \Theta B_s^T) \sum^{-1/2} \|^2_{\mathbf{F}} + \lambda_t p_t(\Theta) + \lambda_s p_s(\Theta) \right], \quad (23)$$

where  $p_t, p_s$  penalize roughness in the  $t$ -direction and the  $s$ -direction, respectively (e.g.,

Eilers and Marx, 2003; Wood, 2006b). As for Problem 1, one can set  $\hat{\sum} = I$  in (23),

resulting in an OLS smoother, or let  $\hat{\sum}^{-1}$  be an estimate of the inverse of

$[\text{Cov}\{y(s_\ell), y(s_{\ell^*})\}|t]_{\ell, \ell^*=1, \dots, L}$ , giving a GLS smoother. An example of a *post-smoothing*

approach is the sandwich smoother (20), although that method was not conceived as such by Xiao et al. (2013).

For the sandwich smoother, the initial pointwise influence matrix  $\tilde{\mathcal{H}}_{\ell}$  equals  $H_t^{\lambda_t}$  for each  $\ell$ . This, combined with Theorem 2 and Lemma 1 in Web Appendix A, implies that the sandwich smoother has final pointwise influence matrix  $H_t^{\lambda_t}$ , and thus  $DF_{\ell} = \text{tr}(H_t^{\lambda_t})$ , for each  $t$ . Recall, however, that our motivating class of applications involves estimating trajectories  $f(\cdot, s)$  with possibly varying complexity. For this purpose it seems suboptimal to have the same pointwise degrees of freedom at each location. We therefore propose an alternative post-smoothing procedure, which generalizes the sandwich smoother to allow for separate temporal smoothing parameters at each location:

1. Derive the initial estimate  $\tilde{Y} = (\tilde{y}_{\cdot 1} \dots \tilde{y}_{\cdot L})$  with  $\tilde{y}_{\cdot \ell} = H_t^{\lambda_{t,\ell}} y_{\cdot \ell}$  where, for each  $\ell$ ,  $\lambda_{t,\ell}$  is chosen by restricted maximum likelihood (Ruppert et al., 2003; Reiss and Ogden, 2009) for an optimal  $\mathcal{A}$ -location smooth. Reiss et al. (2014) propose a fast procedure for this.
2. Obtain the final estimate  $\hat{Y} = \tilde{Y} H_s^{\lambda_s}$ , where the spatial smoothing parameter  $\lambda_s$  may be chosen by multifold cross-validation.

This procedure avoids inversion of  $nL \times nL$  matrices and is thus much faster than the contemporaneous methods, as will be seen below. In this two-step procedure,

$\tilde{DF}_{\ell} = \text{tr}(H_t^{\lambda_{t,\ell}})$ , for  $\ell = 1, \dots, L$ . If we think of  $(\tilde{DF}_1, \dots, \tilde{DF}_L)^T$  as a vector of noisily measured complexity indices for  $\{f(\cdot, s): \ell = 1, \dots, L\}$ , then by Theorem 2, the final pointwise degrees of freedom vector  $(DF_1, \dots, DF_L)^T$  can be viewed as a denoised version of these measurements, with the denoising performed by application of the smoother matrix  $H_s^{\lambda_s}$ . In this sense the final fit adapts in a smooth manner to the varying complexity of  $f(\cdot, s)$ . The application described next illustrates how such adaptivity can be beneficial in practice.

## 6. Application: Development of corpus callosum microstructure

We now return to the corpus callosum fractional anisotropy data described in the introduction. A full description of the image data processing, as well as an illuminating set of analyses, can be found in Imperati et al. (2011). Here we aim to estimate the mean fractional anisotropy  $f(t, s)$  where  $t$  denotes age and  $s$  denotes location, expressed as arc length along the path depicted in Fig. 1.

Before presenting our modeling results, let us consider some evidence for nonlinear change in mean fractional anisotropy with respect to age. For  $\ell = 1, \dots, 107$  we performed a restricted likelihood ratio test (Crainiceanu and Ruppert, 2004), implemented by the `vows` package (Reiss et al., 2016) for R (R Core Team, 2016), to test the null hypothesis that  $f(\cdot, s)$  is linear—that is, mean fractional anisotropy for the  $\mathcal{A}$ h voxel changes linearly with age—against the alternative of nonlinear change. Figure 3 shows the resulting  $p$ -values for each voxel. These  $p$ -values are not adjusted for multiple testing, since our aim here is descriptive

rather than to test the global null hypothesis that  $f(\cdot, s)$  is linear for each  $\ell$ . Also shown are separate penalized spline smooths at three voxels, with the smoothing parameter chosen by restricted maximum likelihood. The first of these voxels is located in the sensorimotor portion of the brain; the second, in the posterior (rear) portion of the prefrontal lobes, with projections into the inferior and middle frontal gyri; and the third, in the anterior (front) portion of the prefrontal lobes. In the first and third voxels, it appears that mean fractional anisotropy attains a peak in young adulthood, and then declines. These nonlinear trajectories are consistent with the  $p$ -values of .019 and .001, respectively, for the two voxels. But for the second voxel, there is no strong evidence for nonlinear change. Features of these curves, such as peaks, are of biological interest, as they may provide insight into characteristic patterns of development for different cognitive abilities. However, the data are quite noisy, suggesting that estimation might be improved by sharing information across voxels.

We applied the post-smoothing procedure of Section 5.2 to this data set, and compare its performance to that of the OLS and GLS smoothers based on (23) with automatic selection of  $\lambda_t$  and  $\lambda_s$  provided by the R package mgcv (Wood, 2006a, 2011). For OLS, we chose these tuning parameters by restricted maximum likelihood, ignoring within-function dependence. While this approach may be suboptimal, it finds some support in recent work (e.g., Krivobokova and Kauermann, 2007; Crainiceanu et al., 2012). For GLS, we treated the

prewhitened response vectors  $\sum_{i=1}^n y_i \cdot (i=1, \dots, n)$  as the data, and again chose  $\lambda_t$  and  $\lambda_s$  by restricted maximum likelihood, as in Reiss et al. (2010). We used 15-dimensional  $B$ -spline bases in the temporal direction, but since the amount of detail to capture is much greater in the spatial direction (i.e., differences among brain regions), we used more  $B$ -splines in that direction: 30 for the post-smoothing method, and slightly fewer, 25, for the contemporaneous methods to mitigate their much higher computation time. The postsmoothing procedure requires a grid search for  $\lambda_s$ ; we selected  $\log(\lambda_s)$  by five-fold crossvalidation from the candidate values  $-20, -19, \dots, 10$ . The post-smoothing fit required 2.1 seconds, versus 23.4 seconds for OLS and 42.6 seconds for GLS.

The rainbow plots (Hyndman and Shang, 2010) in the upper panels of Fig. 4 make it easier to examine the shape of the estimates  $\hat{f}(\cdot, s)$  for different locations  $s$ . Vertical lines are drawn at the same three voxels as in Fig. 3, and vertical progression from blue to red at a given arc length indicates that the estimate  $\hat{f}(\cdot, s)$  is monotonic for that  $s$ . Recall that separate nonparametric regressions suggested that fractional anisotropy attained a peak for the first and third voxels, but not for the second. The OLS estimates  $\hat{f}(\cdot, s)$  are very erratic, with apparently spurious fluctuations with respect to age. The GLS estimates, on the other hand, indicate that mean fractional anisotropy changes monotonically in all three regions, suggesting oversmoothing with respect to age. Only the post-smoothed estimates appear consistent with the scatterplots in Fig. 3: they capture the peak during young adulthood in the sensorimotor and anterior prefrontal regions, whereas in the posterior prefrontal region, the rainbow pattern suggests linear decrease with age.

These impressions are borne out by the lower left panel of Fig. 4, in which the pointwise degrees of freedom for OLS is seen to be uniformly high while that for GLS is uniformly

low. In line with Theorem 2, the pointwise degrees of freedom for the post-smoothed method is a smoothed version of the degrees of freedom for the separate smooths, so this method adapts to spatial variation in temporal-direction complexity. The lower right panel of Fig. 4 presents functional  $R^2$  values (using a definition similar to that of Müller and Yao, 2008) for the three methods against prediction error estimates, based on repeated five-fold crossvalidation (Burman, 1989) with 10 different splits of the observations into five validation sets. This figure provides further evidence that OLS and GLS overfit and underfit the data, respectively, but suggests that in this case underfitting is not very detrimental to predictive performance. It might be objected that these methods' over- and underfitting is attributable to poor automatic selection of the smoothing parameter  $\lambda_t$  in (23). What these results suggest, however, is that the use of a fixed  $\lambda_b$  in itself, severely limits the ability of these contemporaneous smoothers to adapt to the varying complexity of  $f(\cdot, s)$  for different  $s$ .

Figure 5 presents the estimates  $\hat{f}(\cdot, s_\ell)$ , for  $\ell = 91, \dots, 98$ , of the mean fractional anisotropy as a smooth function of age. These eight voxels correspond to arc lengths 94.1–101.5, a range whose right terminus corresponds to the rightmost peak observed in Fig. 4, and to the most prominent trough in the  $p$ -value plot of Fig. 3. The upper left subfigure displays separate curve estimates for the eight voxels. In line with the results shown in Figures 4, the OLS and GLS curves clearly overfit and underfit, respectively, whereas the post-smoothed estimates appear to do a reasonable job of denoising the separate curve estimates, as Theorem 2 might lead us to expect.

## 7. Discussion

In introducing the pointwise influence matrix, we have to a large extent focused on the trace of that matrix, i.e., the pointwise degrees of freedom. We offer here a few remarks on some related (and not-so-related) work on the notion of degrees of freedom.

For linear modeling procedures in which the fitted value matrix cannot be expressed in the linear form (5) for fixed  $\mathbf{H}$ , Ye (1998) proposed the generalized degrees of freedom

$\text{GDF} = \sum_{i=1}^n \text{Cov}(\hat{y}_i, y_i) / \sigma^2$ , where we assume uncorrelated mean-zero errors with common variance  $\sigma^2$ . Following Efron (2004) and Efron et al. (2004), GDF is now often referred to as simply the degrees of freedom. While GDF indeed generalizes the degrees of freedom of Wahba (1983), i.e.  $\text{tr}(\mathbf{H})$  in (5), to a much broader class of fitting procedures (e.g., Tibshirani and Taylor, 2012), it is *less* general than  $\text{tr}(\mathbf{H})$  in that the latter entails no assumptions about the error structure. Notwithstanding the caveats recently issued by Kaufman and Rosset (2014) and Janson et al. (2015) about GDF, we believe that the degrees of freedom of Wahba (1983) remains useful as a description of the complexity of smooth model fits, and our concern here has been to extend it in a completely different direction: from scalar to functional responses. Extending GDF to functional responses would be a worthy goal for future research.

The ordinary degrees of freedom are sometimes partitioned among the parameters in the model (Hodges and Sargent, 2001; Cui et al., 2010). By (14), influence matrices of the

general form (11) imply the expression  $\mathcal{H}_\ell = \left[ \mathbf{b}_s(s_\ell)^T \otimes \mathbf{X} \right] \mathbf{C} (\mathbf{1}_L \otimes \mathbf{I}_n)$  for the point-wise influence matrix, which suggests dividing the degrees of freedom into contributions of the individual columns of  $\mathbf{X}$ , i.e.,

$\text{DF}_\ell(j) = \text{tr} \left[ \{ \mathbf{b}_s(s_\ell)^T \otimes \mathbf{X} \mathbf{e}_{j;p} \mathbf{e}_{j;p}^T \} \mathbf{C} (\mathbf{1}_L \otimes \mathbf{I}_n) \right] = \left[ \mathbf{b}_s(s_\ell)^T \otimes \mathbf{e}_{j;p}^T \right] \mathbf{C} (\mathbf{1}_L \otimes \mathbf{X}_{\cdot j})$  for  $j = 1, \dots, p$ , where  $\mathbf{X}_{\cdot j} = \mathbf{X} \mathbf{e}_{j;p}$  is the  $j$ th column of  $\mathbf{X}$ . This directly extends the degrees of freedom per parameter in section 4.4 of Wood (2006a). By plugging in the particular values of  $\mathbf{C}$  that are given implicitly in Web Appendix A for the functional linear model estimators of Theorem 1, one can readily show that  $\text{DF}_\ell(j) \equiv 1$  for each  $j$  in the linear case. For nonlinear models,  $\text{DF}_\ell(j)$  can be interpreted as a pointwise parameter-specific shrinkage factor.

More straightforwardly,  $\text{DF}_\ell = \text{tr}(\mathcal{H}_\ell)$  can also be partitioned into contributions of the individual observations, given by the  $n$  elements of the main diagonal of  $\mathcal{H}_\ell$ . These elements provide a novel definition of pointwise leverage that could be used to examine which of the functional observations are influential in a pointwise sense.

While we have focused on independent functional observations and the simple settings of Problems 1 and 2, the pointwise influence matrix is readily extended to more complex cases such as longitudinal functional responses and pointwise additive models, as well as to scalar responses observed on a grid or array. Finally, the Bayesian notion of priors on degrees of freedom (Hodges and Sargent, 2001; Hodges, 2013) might be extended to the pointwise setting. We hope to pursue these developments in future work.

An earlier unpublished report by the authors (Reiss et al., 2014) includes additional estimators, comparative simulations, and details not covered here. An R package implementing the methods described here is available at <https://github.com/philreiss/vsm>.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments

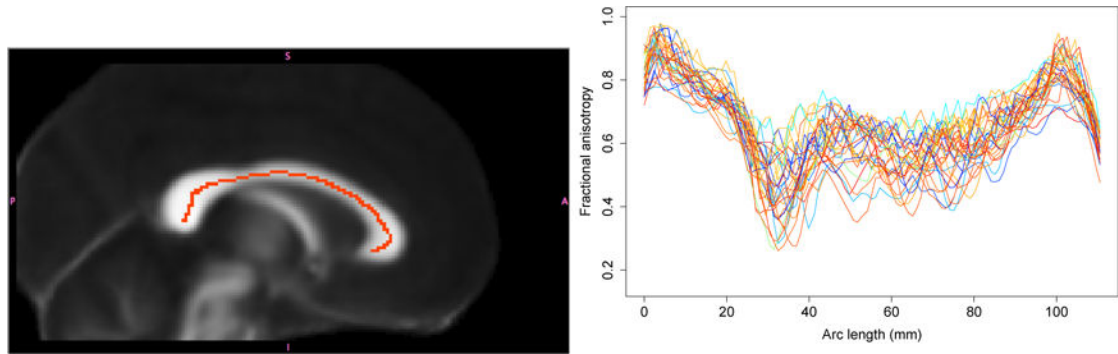
The authors thank the reviewers as well as Yin-Hsiu Chen, Ciprian Crainiceanu, Yair Goldberg, Jeff Goldsmith, Lan Huo, Mike Milham, Todd Ogden, Eva Petkova, David Ruppert, Fabian Scheipl, and Simon Wood for very helpful advice and feedback. The first author's work was supported by U.S. National Science Foundation grant DMS-0907017, and the work of the first four authors was supported by U.S. National Institutes of Health grant 1R01MH095836.

## References

- Buja A, Hastie T, Tibshirani R. Linear smoothers and additive models. *Ann Statist.* 1989; 17:453–510.
- Burman P. A comparative study of ordinary cross-validation,  $v$ -fold cross-validation and the repeated learning-testing methods. *Biometrika.* 1989; 76:503–514.
- Chiou JM, Müller HG, Wang JL. Functional quasi-likelihood regression models with smooth random effects. *J R Statist Soc B.* 2003; 65:405–423.
- Crainiceanu CM, Ruppert D. Likelihood ratio tests in linear mixed models with one variance component. *J R Statist Soc B.* 2004; 66:165–185.

- Crainiceanu CM, Staicu AM, Ray S, Punjabi N. Bootstrap-based inference on the difference in the means of two correlated functional processes. *Statistics in Medicine*. 2012; 31:3223–3240. [PubMed: 22855258]
- Cui Y, Hodges JS, Kong X, Carlin BP. Partitioning degrees of freedom in hierarchical and other richly parameterized models. *Technometrics*. 2010; 52:124–136. [PubMed: 20559456]
- Efron B. The estimation of prediction error (with Discussion). *J Am Statist Assoc*. 2004; 99:619–642.
- Efron B, Hastie T, Johnstone I, Tibshirani R. Least angle regression (with Discussion). *Ann Statist*. 2004; 32:407–499.
- Eilers PHC, Marx BD. Flexible smoothing with B-splines and penalties. *Stat Sci*. 1996; 11:89–102.
- Eilers PHC, Marx BD. Multivariate calibration with temperature interaction using two-dimensional penalized signal regression. *Chemometr Intell Lab*. 2003; 66:159–174.
- Fan J, Zhang JT. Two-step estimation of functional linear models with applications to longitudinal data. *J R Statist Soc B*. 2000; 62:303–322.
- Green, PJ., Silverman, BW. *Nonparametric Regression and Generalized Linear Models: A Roughness Penalty Approach*. Chapman & Hall; Boca Raton, Florida: 1994.
- Hoaglin DC, Welsch RE. The hat matrix in regression and ANOVA. *The American Statistician*. 1978; 32:17–22.
- Hodges, JS. *Richly Parameterized Linear Models: Additive, Time Series, and Spatial Models Using Random Effects*. CRC Press; Boca Raton, Florida: 2013.
- Hodges JS, Sargent DJ. Counting degrees of freedom in hierarchical and other richly-parameterised models. *Biometrika*. 2001; 88:367–379.
- Hyndman RJ, Shang HL. Rainbow plots, bagplots, and boxplots for functional data. *J Comp Graph Statist*. 2010; 19:29–45.
- Imperati D, Colcombe S, Kelly C, Di Martino A, Zhou J, Castellanos FX, Milham MP. Differential development of human brain white matter tracts. *PLoS ONE*. 2011; 6:e23437. [PubMed: 21909351]
- Janson L, Fithian W, Hastie TJ. Effective degrees of freedom: a flawed metaphor. *Biometrika*. 2015; 102:479–485. [PubMed: 26977114]
- Kaufman S, Rosset S. When does more regularization imply fewer degrees of freedom? Sufficient conditions and counterexamples *Biometrika*. 2014; 101:771–784.
- Krivobokova T, Kauermann G. A note on penalized spline smoothing with correlated errors. *J Am Statist Assoc*. 2007; 102:1328–1337.
- Müller HG, Yao F. Functional additive models. *J Am Statist Assoc*. 2008; 103:1534–1544.
- R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing; Vienna, Austria: 2016.
- Ramsay, JO., Silverman, BW. *Functional Data Analysis*. 2nd. Springer; New York: 2005.
- Reiss, P., Chen, YH., Huang, L., Huo, L., Tan, R., Jiao, R. *vows: Voxelwise Semiparametrics*. 2016. R package version 0.5
- Reiss, PT., Huang, L., Chen, H., Colcombe, S. Varying-smoother models for functional responses. arXiv:1412.0778 [stat.ME]. 2014. available at <http://arxiv.org/abs/1412.0778>
- Reiss PT, Huang L, Chen YH, Huo L, Tarpey T, Mennes M. Massively parallel nonparametric regression, with an application to developmental brain mapping. *J Comp Graph Statist*. 2014; 23:232–248.
- Reiss PT, Huang L, Mennes M. Fast function-on-scalar regression with penalized basis expansions. *Int J Biostat*. 2010; 6 article 28.
- Reiss PT, Ogden RT. Smoothing parameter selection for a class of semiparametric linear models. *J R Statist Soc B*. 2009; 71:505–523.
- Ruppert, D., Wand, MP., Carroll, RJ. *Semiparametric Regression*. Cambridge University Press; New York: 2003.
- Schumaker, L. *Spline Functions: Basic Theory*. 3rd. Cambridge University Press; Cambridge, UK: 2007.

- Shaw P, Eckstrand K, Sharp W, Blumenthal J, Lerch JP, Greenstein D, Clasen L, Evans A, Giedd J, Rapoport JL. Attention-deficit/hyperactivity disorder is characterized by a delay in cortical maturation. *Proc Natl Acad Sci USA*. 2007; 104:19649–19654. [PubMed: 18024590]
- Tibshirani RJ, Taylor J. Degrees of freedom in lasso problems. *Ann Statist*. 2012; 40:1198–1232.
- Wahba G. Bayesian “confidence intervals” for the cross-validated smoothing spline. *J R Statist Soc B*. 1983; 45:133–150.
- Wood, SN. *Generalized Additive Models: An Introduction with R*. Chapman & Hall; Boca Raton, Florida: 2006a.
- Wood SN. Low-rank scale-invariant tensor product smooths for generalized additive mixed models. *Biometrics*. 2006b; 62:1025–1036. [PubMed: 17156276]
- Wood SN. Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models. *J R Statist Soc B*. 2011; 73:3–36.
- Xiao L, Li Y, Ruppert D. Fast bivariate P-splines: the sandwich smoother. *J R Statist Soc B*. 2013; 75:577–599.
- Ye J. On measuring and correcting the effects of data mining and model selection. *J Am Statist Assoc*. 1998; 93:120–131.
- Zhu H, Fan J, Kong L. Spatially varying coefficient model for neuroimaging data with jump discontinuities. *J Am Statist Assoc*. 2014; 109:1084–1098.

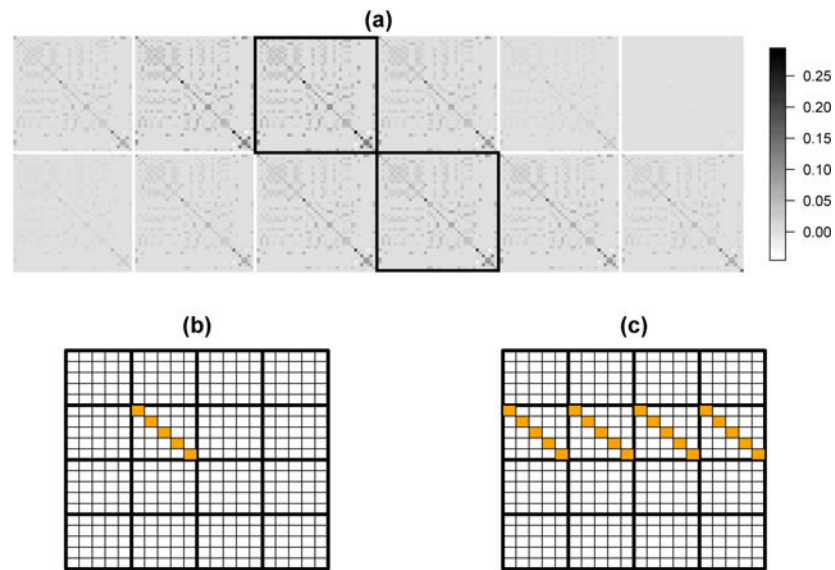


**Figure 1.**

Left: Sequence of 107 corpus callosum voxels at which fractional anisotropy was recorded.

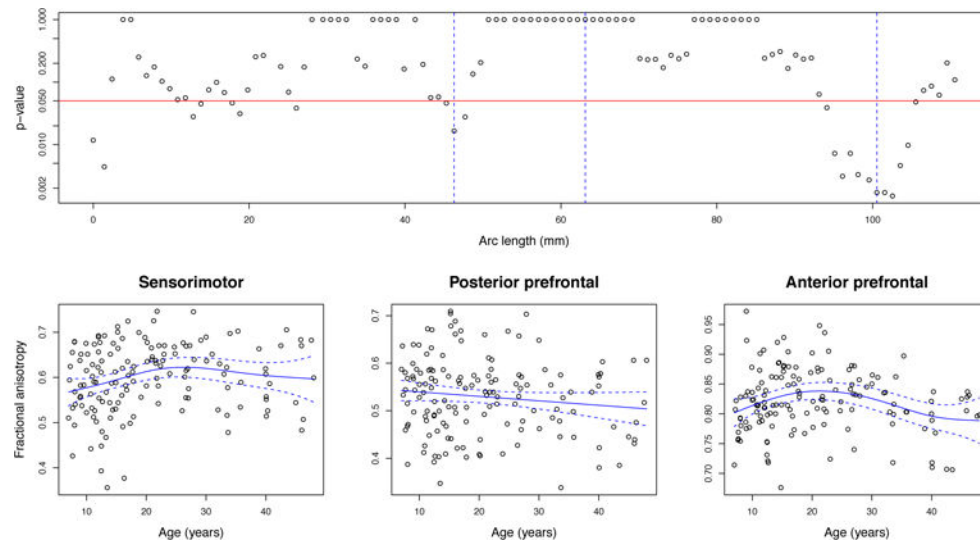
Right: Rainbow plots displaying how the resulting fractional anisotropy profiles (functional responses) vary with age, which ranges from 7 (dark blue) to 48 (red).





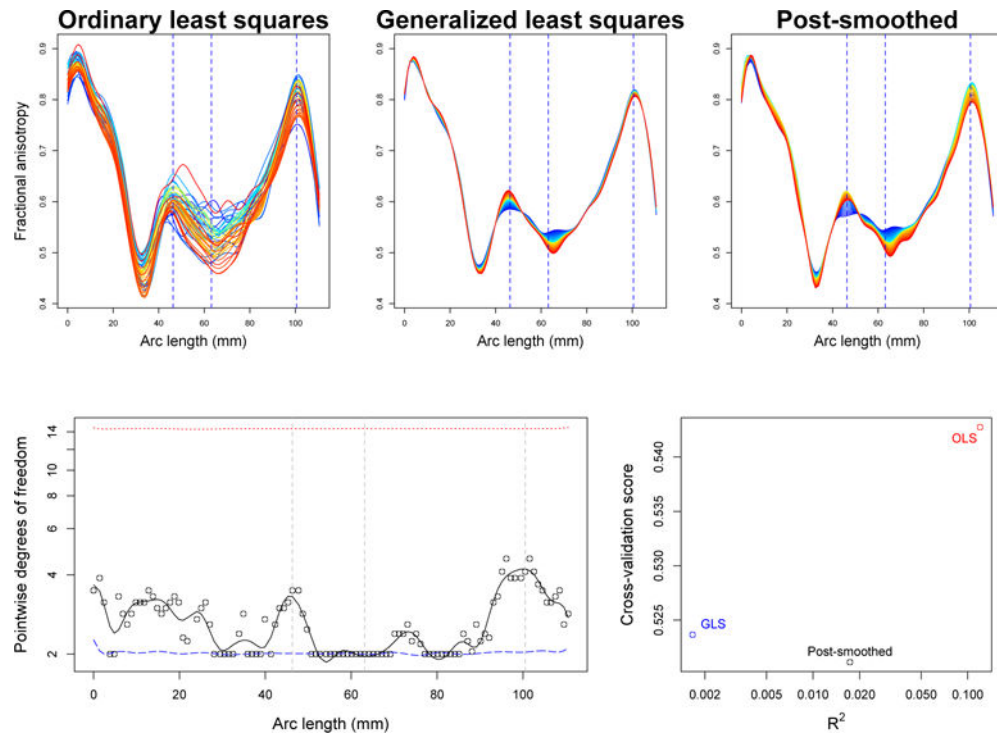
**Figure 2.**

(a) Excerpt from the block influence matrix (10) for a fit to the corpus callosum data, with diagonal blocks outlined in black. (b) Schematic illustration of the usual definition of degrees of freedom for the 2nd location, in a toy example with separate models at each of  $L = 4$  locations. (c) Proposed pointwise degrees of freedom  $DF_2$  for the 2nd location. This figure appears in color in the electronic version of this article.



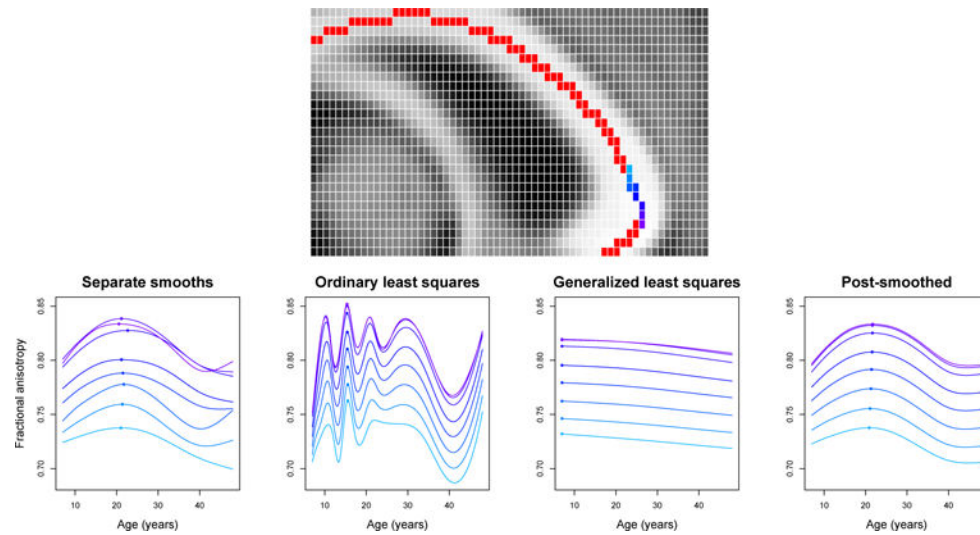
**Figure 3.**

Above,  $p$ -values from restricted likelihood ratio tests of the null hypothesis that mean fractional anisotropy changes linearly with age. Below, curve estimates,  $\pm 2$  approximate standard errors, for the three voxels indicated above by dashed vertical lines; headings refer to the brain regions where these voxels are located. This figure appears in color in the electronic version of this article.



**Figure 4.**

Above: Rainbow plots of fitted value functions  $\hat{f}(t, \cdot)$  with  $t = 7, 8, \dots, 48$ , with 7 in dark blue and 48 in red. Lower left: Pointwise degrees of freedom for the corpus callosum data—separate smooths at each of the 107 voxels (circles), ordinary least squares (dotted curve), generalized least squares (dashed curve), and post-smoothed (solid curve). The vertical axis is on a log scale for clarity. Lower right: Functional  $R^2$  plotted against prediction error estimate from repeated five-fold cross-validation.



**Figure 5.**

Estimated mean fractional anisotropy as a function of age for eight voxels in the prefrontal cortex, which are displayed in the image above. Dots indicate estimated age of peak fractional anisotropy. This figure appears in color in the electronic version of this article.

**Table 1**

Representative references for some existing approaches to fitting function-on-scalar regression models of the general form (1).

<i>Problem</i>	<i>Dependence of y on predictor(s)</i>	<i>Contemporaneous smoothing</i>	<i>Post-smoothing</i>
1	Linear	Ramsay and Silverman (2005), Reiss et al. (2010)	Fan and Zhang (2000), Zhu et al. (2014)
2	Nonlinear	Wood (2006b)	Xiao et al. (2013)