



Simple and multiple P-splines regression with shape constraints

Kaatje Bollaerts^{1,3*}, Paul H. C. Eilers² and Iven van Mechelen¹

¹Katholieke Universiteit Leuven, Belgium

²Leiden University Medical Centre, The Netherlands

³Universiteit Hasselt, Belgium

In many research areas, especially within social and behavioural sciences, the relationship between predictor and criterion variables is often assumed to have a particular shape, such as monotone, single-peaked or U-shaped. Such assumptions can be transformed into (local or global) constraints on the sign of the n th-order derivative of the functional form. To check for such assumptions, we present a non-parametric regression method, P-splines regression, with additional asymmetric discrete penalties enforcing the constraints. We show that the corresponding loss function is convex and present a Newton–Raphson algorithm to optimize. Constrained P-splines are illustrated with an application on monotonicity-constrained regression with both one and two predictor variables using data from research on cognitive development of children.

1. Introduction

In many research areas, much effort is made to investigate the relationship between a set of predictor variables (also called independent variables) X_i and a criterion or dependent variable Y . Such a relationship is mostly denoted as a function f in $Y = f(X_1, \dots, X_i, \dots, X_n)$. Examples include the subjective value of money as a function of a person's wealth or a child's task performance as a function of the child's motivation. Often in social and behavioural sciences, theoretical hypotheses regarding the shape of the relationship between predictor variables and a criterion variable are made. For instance, the subjective value of money is assumed to be a decreasing function of a person's wealth. Or a child's performance is assumed to increase first and then decrease as a function of increasing motivation. Such hypotheses, which are commonly formulated within social and behavioural sciences, are mostly not easily translated in terms of a specific parametric relationship. Instead,

*Correspondence should be addressed to Kaatje Bollaerts, Universiteit Hasselt, Center for Statistics, Agoralaan 1 Gebouw D, B-3590 Diepenbeek, Belgium (e-mail: kaatje.bollaerts@uhasselt.be).

they typically imply a non-parametric functional relationship such as monotonicity, U-shapedness or single-peakedness. For instance, it is more realistic to assume that the subjective value of money is a monotone decreasing function of a person's wealth rather than assuming, say, a logarithmic or a linear function. In the following, we will discuss in more detail some non-parametric shapes, together with substantive psychological theories in which these non-parametric shapes play a central role.

Monotonicity means that f is either monotone decreasing or monotone increasing. As an example, consider a function f which expresses the relationship between some measure of cognitive performance of children and age. In this case, it is common to assume that f is a monotone increasing function; note, in this respect, that assuming monotonicity is more plausible than assuming linearity or an exponential function. Another example stems from the trait approach to personality (McCrae & Costa, 1987). Within this approach, it is assumed that the ordering of persons with respect to a particular behaviour, such as fighting in a specific situation, corresponds to the ordering of persons with respect to some relevant underlying trait (e.g. aggressiveness). Hence, people's behaviour is assumed to be a monotone increasing function of the relevant underlying trait, even though no specific parametric assumption is made.

A second non-parametric functional form we may consider is *U-shapedness*, as put forward in many theories concerning growth and development. Strauss (1982) defines U-shaped behavioural growth curves as curves indicating the initial appearance of a behaviour, a later dropping out of that behaviour and its subsequent reappearance.

A third non-parametric functional form we consider is *single-peakedness*. A single-peaked function is a function that increases up to some point and then decreases. The importance of non-parametric single-peaked functions has been underscored by Coombs (1977). In particular, preference and psychophysical functions are frequently observed to be single-peaked.

So far, we have only considered examples of assumptions on non-parametric regression on a single predictor variable. However, non-parametric functional forms with respect to two or more predictor variables can be assumed as well. For instance, one may consider the assumed monotone increasing functional relationship between some measure of cognitive performance of children and both age and amount of training, or between aggressiveness of persons and both trait anger and the anger-eliciting power of the situation (McCrae & Costa, 1987). Yet another example concerns the single-peaked preference functions for options varying along two dimensions, such as preference for beer as a function of its alcohol level and its bitterness (Coombs, 1977).

The examples above illustrate the versatility of assumptions that imply non-parametric functional forms. In this paper, we will introduce a statistical tool by which such assumptions can be dealt with, namely constrained P-splines regression. P-splines regression as such is introduced in a target article by Eilers and Marx (1996) with illustrations on data containing a single predictor variable. An illustration of P-splines regression with two predictor variables is given in Durban, Currie, and Eilers (2002). Computational issues of smoothing with different predictor variables are discussed in Eilers, Currie, and Durban (2006). In the present paper, the main focus is on P-splines regression with shape constraints, which is briefly introduced in Eilers (1994). Illustrations with both one and two predictor variables will be given.

The remainder of this paper is organized as follows. In Section 2 we will discuss simple as well as multivariate unconstrained and constrained P-splines regression. In Section 3 an application of simple and multiple P-splines regression with monotonicity constraints is given, with data from research on cognitive development of children. In Section 4 we present some concluding remarks. The condition and optimization of the loss function are discussed in Appendices A and B, respectively.

2. Method

In this section, we successively discuss unconstrained and constrained P-splines regression. In each case, we will discuss regression with one and two predictor variables. The discussion of unconstrained P-splines regression with a single predictor variable is mainly based on Eilers and Marx (1996).

2.1. Unconstrained P-splines regression

2.1.1. Simple regression

Eilers and Marx (1996) introduced non-parametric regression with P-splines, which is essentially least squares regression with an excessive number of univariate B-splines (De Boor, 1978; Dierckx, 1993) and an additional discrete penalty to correct for overfitting.

Univariate B-splines are piecewise linear functions with local support. A B-spline of degree q consists of $q + 1$ polynomial pieces of degree q joined smoothly (i.e. differentiable of order $q - 1$) at q points λ_i (called *interior knots*) between boundaries λ_{\min} and λ_{\max} (called *exterior knots*) and with a positive value between and a value of zero outside these boundaries. An example of a B-spline of the first degree is given in Figure 1a; it is clear that this B-spline consists of two linear pieces joined at one interior knot. An example of a B-spline of the third degree is given in Figure 1b; this B-spline consists of four cubic pieces joined smoothly at three interior knots. Note that the B-splines shown in Figure 1 both have equally spaced knots. B-splines with unequally spaced knots exist as well, but are not considered in this paper.

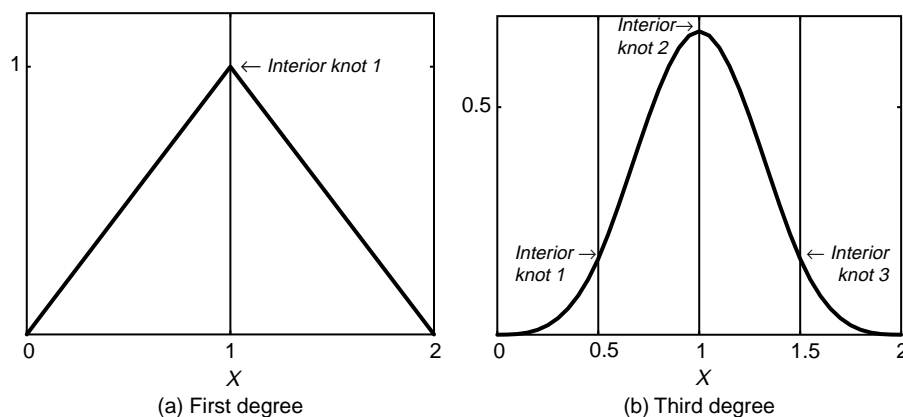


Figure 1. Single B-splines of first and third degree.

In order to use B-splines for non-parametric regression, a basis of r overlapping B-splines is constructed, which is such that

$$\forall x : \sum_{j=1}^r B_j(x, q) = 1 \quad (1)$$

with $B_j(x, q)$ denoting a B-spline of degree q with leftmost knot j . In Figure 2a one can see an example of a basis of B-splines of the third degree, which is the most commonly used degree in B-splines regression.

The B-splines of a B-spline basis act as the predictors in spline regression. Given m observations (x_i, y_i) , least squares regression with B-splines of Y on the basis of X comes down to minimizing the following loss function:

$$S = \sum_{i=1}^m (y_i - \hat{y}_{(\alpha)_i})^2 \quad (2)$$

with

$$\hat{y}_{(\alpha)_i} = \sum_{j=1}^r \alpha_j B_j(x_i, q) \quad (3)$$

and the α_j s being the coefficients (or amplitudes) of the corresponding B-splines. In Figure 2b, spline regression with B-splines of the third degree is illustrated.

A major problem in B-spline regression is the choice of the optimal number of B-splines. An insufficient number of B-splines leads to underfitting such that the fitted curve is too smooth and, hence, relevant information is neglected. On the other hand, too many B-splines lead to overfitting such that the fitted curve is too flexible and, hence, random fluctuations are modelled. To overcome this problem, O'Sullivan (1988) suggested using an excessive number of B-splines with a smoothness penalty consisting of the integrated squared second-order derivative of the fitted curve, in order to correct for overfitting; this approach has become standard in spline literature. Eilers and Marx (1996) propose using an excessive number of equally spaced B-splines together with a discrete smoothness penalty based on second- (or higher-) order differences of the coefficients of adjacent B-splines. They call this approach *P-splines regression*, and it is very similar to O'Sullivan's approach. Furthermore, P-splines regression is easy to implement, has no

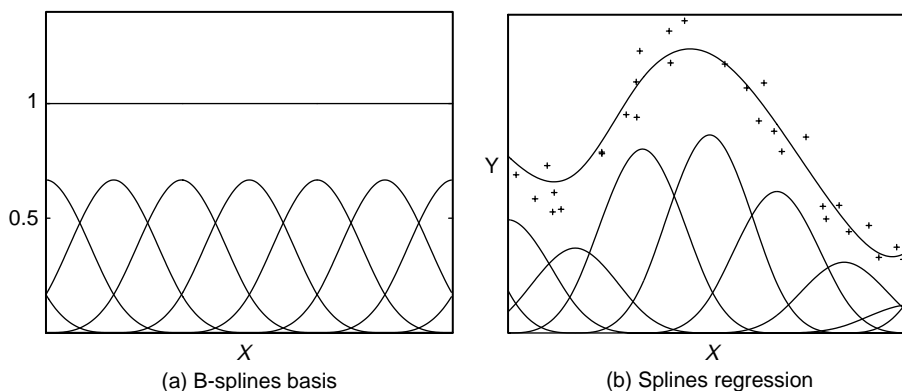


Figure 2. Spline regression with B-splines of third degree.

boundary effects, conserves moments and has polynomial curve fits as limits. For a penalty on second-order differences, the corresponding least squares loss function equals

$$S = \sum_{i=1}^m (y_i - \hat{y}_{(\alpha)_i})^2 + \lambda \sum_{j=3}^r (\Delta^2 \alpha_j)^2, \quad (4)$$

with λ being the smoothness parameter and $\Delta^2 \alpha_j$ being the second-order differences, that is,

$$\Delta^2 \alpha_j = \Delta^1(\Delta^1 \alpha_j) = \alpha_j - 2\alpha_{j-1} + \alpha_{j-2}.$$

The system of equations that follows from minimization of the loss function given in (4) can be written in matrix notation as

$$\mathbf{B}^T \mathbf{y} = (\mathbf{B}^T \mathbf{B} + \lambda (\mathbf{D}^2)^T \mathbf{D}^2) \hat{\alpha}, \quad (5)$$

with \mathbf{B} being the design matrix consisting of B-splines and \mathbf{D}^k the matrix representation of the difference operator Δ^2 .

The amount of smoothness can be controlled for by means of λ , which is a user-defined smoothness parameter. If $\lambda \rightarrow \infty$, then, for regression with a smoothness penalty on m th-order differences, the fitted function will approach a polynomial of degree $m - 1$. To choose an optimal value for λ , Eilers and Marx (1996) propose using Akaike's information criterion:

$$AIC(\lambda) = -2L(\alpha, \lambda) + 2\dim(\alpha, \lambda), \quad (6)$$

with $L(\alpha, \lambda)$ denoting the log-likelihood of the data and $\dim(\alpha, \lambda)$ the effective dimension of the vector of parameters. The determination of the latter requires some extra attention. Indeed, since the rationale behind P-splines regression is to use an excessive number of B-splines with a penalty to correct for overfitting, the total number of parameters of the P-splines model is an overestimation of the effective dimension of the vector of parameters. This problem can be solved by using the trace of the hat matrix \mathbf{H} as an approximation of the effective dimension of the vector of parameters (Hastie & Tibshirani, 1990). For P-splines regression, the trace of the hat matrix equals

$$\text{tr}(\mathbf{H}) = \text{tr}(\mathbf{B}^T \mathbf{B} + \lambda (\mathbf{D}^k)^T \mathbf{D}^k)^{-1} \mathbf{B}^T \mathbf{B}. \quad (7)$$

Then, under the assumption of normally distributed homoscedastic errors, $\tilde{y}_i \mathcal{N}(\hat{y}_i, \sigma^2)$, Akaike's information criterion equals

$$AIC(\lambda) = \sum_{i=1}^m \frac{(y_i - \hat{y}_i)^2}{\hat{\sigma}^2} + 2m \ln(\hat{\sigma}) + m \ln(2\pi) + 2\text{tr}(\mathbf{H}). \quad (8)$$

As an estimate of the nuisance parameter $\hat{\sigma}^2$, Eilers and Marx (1996) propose using the variance of the residuals computed for an optimal value for λ chosen on the basis of (generalized) cross-validation.

2.1.2. Multiple regression with two predictor variables

P-splines regression with two predictor variables is a straightforward extension of P-splines regression with one predictor variable as introduced in the previous section. The constitutive elements of P-splines regression with two predictor variables are bivariate B-splines, illustrated in Figure 3a. A bivariate B-spline of degree q is the product

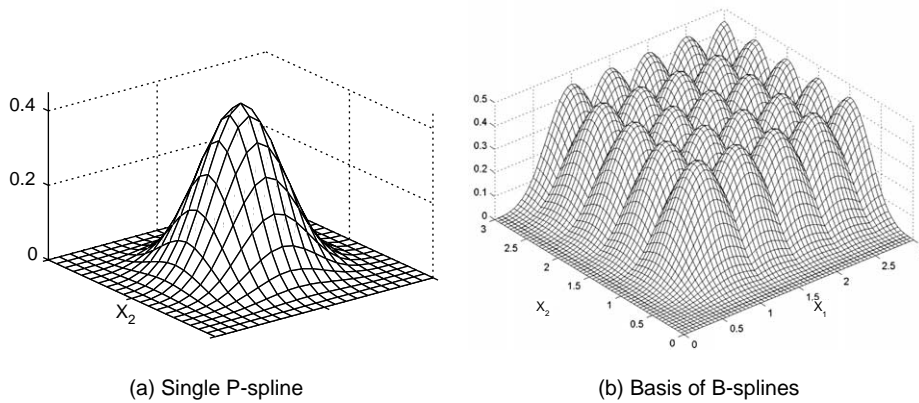


Figure 3. Bivariate B-splines of third degree.

of two univariate B-splines of degree q , that is,

$$B_{jk}(\mathbf{x}_1, \mathbf{x}_2, q) = B_j^{(1)}(\mathbf{x}_1, q) \times B_k^{(2)}(\mathbf{x}_2, q), \quad (9)$$

with $B^{(1)}$ and $B^{(2)}$ denoting two distinct univariate B-splines.

A basis of bivariate B-splines computed as the tensor product of two vectors of B-splines of the third degree is displayed in Figure 3b. Then, for m observations $(\mathbf{x}_{1i}, \mathbf{x}_{2i}, y_i)$ regression of y on \mathbf{x}_1 and \mathbf{x}_2 with a basis of $r \times r'$ overlapping bivariate B-splines comes down to minimizing the following least squares loss function:

$$S = \sum_{i=1}^m (y_i - \hat{y}_{(\alpha)_i})^2, \quad (10)$$

with

$$\hat{y}_{(\alpha)_i} = \sum_{j=1}^r \sum_{k=1}^{r'} \alpha_{jk} B_{jk}(\mathbf{x}_1, \mathbf{x}_2, q), \quad (11)$$

where the α_{jk} s are the coefficients of the corresponding bivariate B-splines.

In P-splines regression with two predictor variables, two smoothness penalties are used, one for each predictor variable. For penalties on second-order differences, the loss function to be minimized is

$$S = \sum_{i=1}^m (y_i - \hat{y}_{(\alpha)_i})^2 + \lambda_1 \sum_{j=3}^r \sum_{k=1}^{r'} (\Delta_1^2 \alpha_{jk})^2 + \lambda_2 \sum_{k=1}^{r'} \sum_{j=3}^r (\Delta_2^2 \alpha_{jk})^2 \quad (12)$$

with Δ_1^2 being a column-wise and Δ_2^2 a row-wise smoothness penalty on the matrix $\mathbf{A} = [\alpha_{jk}]$, and with λ_1 and λ_2 being the smoothness parameters of \mathbf{x}_1 and \mathbf{x}_2 , respectively.

2.2. Constrained P-splines regression

2.2.1. Simple regression

As indicated in the Introduction, assumptions of non-parametric functional forms that can be transformed into local or global constraints on the sign of the n th-order derivative

can be checked using constrained P-splines regression. This is P-splines regression with an asymmetric discrete penalty on n -th order differences, reflecting the assumed non-parametric functional form. This penalty is asymmetric since it differentially penalizes positive and negative n -th order differences, in order to restrict the sign of the n -th order differences and as such restrict the sign of the n -th order derivative of the fitted function. Indeed, the first-order derivative of a B-splines function with equally spaced knots equals

$$f^{(1)}(x) = \frac{\partial f(x)}{\partial x} = \frac{\partial}{\partial x} \sum_{j=1}^r \alpha_j B_j(x, q) = (qb)^{-1} q \sum_{j=1}^{r+1} \Delta_{\alpha_j}^1 B_j(x, q-1), \quad (13)$$

with b denoting the distance between two adjacent knots (De Boor, 1978). By induction, the n -th order derivative of a B-splines function is

$$f^{(n)}(x) = \frac{\partial^n f(x)}{\partial x^n} = \prod_{l=1}^n [(q+1-l)b]^{-1} (q+1-l) \sum_{j=1}^{r+n} \Delta_{\alpha_j}^n B_j(x, q-n). \quad (14)$$

Then, since b , $q+1-l$ and $B_j(x, q-n)$ are all positive by definition, restricting $\Delta_{\alpha_j}^n$ to be positive is a sufficient condition for the $f^{(n)}(x)$ to be positive. Similarly, restricting $\Delta_{\alpha_j}^n$ to be negative is a sufficient condition for the $f^{(n)}(x)$ to be negative. In addition, for $q-n=0$ and $q-n=1$, these sufficient conditions are necessary as well since in that case $f^{(n)}(x)$ is piecewise constant or piecewise linear, respectively.

Hence, a penalty reflecting the constraint of a positive n -th order derivative within a range as defined by indicator variable v_j is

$$\sum_{j=n+1}^r v_j w(\alpha)_j (\Delta^n \alpha_j)^2 \quad (15)$$

with

$$v_j = \begin{cases} 1, & \text{if the constraint on } \partial^n f(x)/\partial x^n \text{ is to hold on at least part of the support of } B_j, \\ 0, & \text{otherwise,} \end{cases}$$

and with

$$w(\alpha)_j = \begin{cases} 0, & \text{if } \Delta^n \alpha_j \geq 0, \\ 1, & \text{otherwise,} \end{cases}$$

being asymmetric weights. As can easily be seen, negative values for $\Delta^n \alpha_j$ are penalized whereas non-negative are not. Then, with κ being a user-defined constraint parameter, the overall loss function is:

$$S = \sum_{i=1}^m (y_i - \hat{y}(\alpha)_i)^2 + \lambda \sum_{j=3}^r (\Delta^2 \alpha_j)^2 + \kappa \sum_{j=n+1}^r v_j w(\alpha)_j (\Delta^n \alpha_j)^2, \quad (16)$$

which is convex in α (for a proof, see Appendix A). The corresponding system of equations that follows from minimization of S equals

$$\mathbf{B}'\mathbf{y} = (\mathbf{B}^T \mathbf{B} + \lambda (\mathbf{D}^2)^T \mathbf{D}^2 + \kappa (\mathbf{D}^1)^T \mathbf{V} \mathbf{W} \mathbf{D}^1) \hat{\alpha}, \quad (17)$$

with \mathbf{B} and \mathbf{D}^2 defined as in (5), with \mathbf{D}^1 being the matrix representation of the difference operator Δ^1 and with \mathbf{V} and \mathbf{W} diagonal matrices with elements v_j and $w(\alpha)_j$, respectively. A similar reasoning holds for a constraint of a negative n th-order derivative. In this case, positive values for $\Delta^n \alpha_j$ are penalized whereas non-positive are not.

2.2.2. Multiple regression with two predictor variables

Again, it is straightforward to extend constrained P-splines regression with one predictor variable to constrained P-splines regression with two predictor variables. In the latter case, two constraint penalties are used, such that the corresponding loss function equals

$$S = \sum_{i=1}^m (\mathbf{y}_i - \hat{\mathbf{y}}(\boldsymbol{\alpha})_i)^2 + \lambda_1 \sum_{j=3}^r \sum_{k=1}^{r'} (\Delta_1^2 \alpha_{jk})^2 + \lambda_2 \sum_{j=1}^r \sum_{k=3}^{r'} (\Delta_2^2 \alpha_{jk})^2 \quad (18)$$

$$+ \kappa_1 \sum_{j=n+1}^r \sum_{k=1}^{r'} v_{1,jk} w(\boldsymbol{\alpha})_{1,jk} (\Delta_1^n \alpha_{jk})^2 + \kappa_2 \sum_{j=1}^r \sum_{k=n+1}^{r'} v_{2,jk} w(\boldsymbol{\alpha})_{2,jk} (\Delta_2^n \alpha_{jk})^2$$

with $\Delta_1^n \alpha_{jk}$ being a column-wise and $\Delta_2^n \alpha_{jk}$ a row-wise constraint penalty on the matrix $\mathbf{A} = [\alpha_{jk}]$, $v_{1,jk}$ and $v_{2,jk}$ being indicator variables defining the range for which the constraints should hold, with

$$v_{1,jk} = \begin{cases} 1, & \text{if the constraint on } \partial^n f(x_1, x_2) / \partial x_1^n \text{ is to hold on at least part of the support of } B_{jk}, \\ 0, & \text{otherwise,} \end{cases}$$

$$v_{2,jk} = \begin{cases} 1, & \text{if the constraint on } \partial^n f(x_1, x_2) / \partial x_2^n \text{ is to hold on at least part of the support of } B_{jk}, \\ 0, & \text{otherwise,} \end{cases}$$

$$w(\boldsymbol{\alpha})_{1,jk} = \begin{cases} 0, & \text{if } \Delta_1^n \alpha_j \geq 0, \\ 1, & \text{otherwise,} \end{cases}$$

$$w(\boldsymbol{\alpha})_{2,jk} = \begin{cases} 0, & \text{if } \Delta_2^n \alpha_j \geq 0, \\ 1, & \text{otherwise,} \end{cases}$$

for the constraint of a positive n th-order derivative with respect to x_1 and x_2 , respectively, and with κ_1 and κ_2 being user-defined constraint parameters. The loss function in (15) is convex, of which the proof is a straightforward extension of the one given in Appendix A. Again, a similar reasoning holds for a constraint of a negative n th-order derivative. In this case, positive values for $\Delta^n \alpha_j$ are penalized whereas non-positive are not.

2.3. Algorithm

We will make use of a Newton-Raphson procedure in order to find an optimal solution of the loss functions described in (16) and (18). An iteration of this procedure comes down to calculating $w(\boldsymbol{\alpha})$ on $\boldsymbol{\alpha}$ as estimated in the previous iteration and calculating the new estimates, $\boldsymbol{\alpha}'$, conditional on $w(\boldsymbol{\alpha})$. A schematic presentation of the algorithm reads

as follows:

1. $l \leftarrow 0$
2. set initial weights $\mathbf{VW}^{(l)} = [0]$
3. $l \leftarrow l + 1$
4. estimate $\alpha^{(l)}$ on $\mathbf{VW}^{(l-1)}$
5. calculate $\mathbf{VW}^{(l)}$ on $\alpha^{(l)}$
6. repeat step 3, 4 and 5 until $\mathbf{VW}^{(l)} = \mathbf{VW}^{(l+1)}$
7. if $\mathbf{VW}^{(l)} = \mathbf{VW}^{(l+1)}$, $\alpha^{(l)}$ is the optimal solution sought

For a more in-depth discussion, the reader is referred to Appendix B. The corresponding MATLAB software is available upon request.

3. Application

In this section, we discuss an application on monotonicity-constrained P-splines regression, which is P-splines regression with an additional discrete penalty, forcing the first-order differences to be positive. The data we use come from a study on cognitive development of children (van der Maas, 1993; van der Maas & Molenaar, 1992). In this study the understanding that an amount of liquid remains the same when you pour it into another container (i.e. conservation of liquid) is investigated. Conservation of liquid was introduced by Piaget (1960), the well-known pioneer of stagewise developmental evolution. Piaget distinguished between three acquisition stages: a non-conserving equilibrium stage, a transitional disequilibrium stage and a conserving equilibrium stage. In the non-conserving equilibrium stage, children believe that the amount of liquid may increase or decrease when it is poured from one container to another. In the transitional disequilibrium stage, children start to realize that pouring liquid from one container to another does not change quantity; however, this insight is not yet consolidated. From the conserving equilibrium stage only, children truly understand conservation of liquid. Based on this theory, van der Maas (1993) developed a computer test to measure conservation understanding. We will now give a description of this test.

3.1. Computer test of liquid conservation

The computer test of liquid conservation consists of three different parts (van der Maas, 1993). Since we only use data with respect to the first part, we restrict our description of the test to this part. The latter consists of eight items. Each item contains: (1) an initial situation consisting of two identical containers filled with liquid, (2) a transformation which involves pouring the liquid of one of the two containers into an empty container with a different shape, and (3) the resulting situation. Both the initial and the resulting situation are to be judged by the respondent using three response alternatives: more liquid in the left container, the same amount of liquid in both containers, more liquid in the right container. Furthermore, three different types of items are included in this part of the test: three standard equality items, three standard inequality items and two guess items. An example of each type is shown in Figure 4. In a standard equality item, though the containers of the initial situation are both filled with a same quantity of liquid, the height of the liquid differs in the two containers of the resulting situation. In a standard

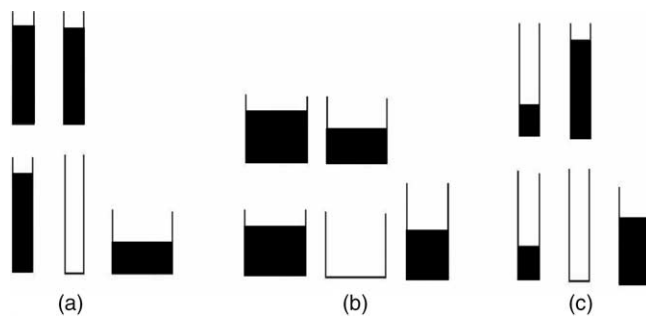


Figure 4. Three different types of items of the liquid conservation test: (a) standard equality item, (b) standard inequality item, (c) guess item. The top two containers constitute the initial situation, the bottom three the resulting situation.

inequality item, the containers of the initial situation are not filled with a same quantity of liquid; however, the transformation results in equal heights of the liquid in both containers. Finally, in guess items, the containers of the initial situation are not filled with a same quantity of liquid nor are the heights equal in the final situation.

With the help of these three types of items, differentiation between conservers, non-conservers and guessers is possible. Conservers, who understand that transformation does not change quantity, are expected to answer correctly that standard equality items contain the same quantity of liquid, whereas standard inequality and guess items do not. Non-conservers are assumed to focus on height and to conclude that two containers are filled with a same quantity of liquid if the heights are equal; this should result in correct answers for guess items and incorrect answers for standard items. Finally, guessers are expected to score at chance level irrespective of the item type.

3.2. Data

Participants were 101 children with ages ranging from 6.2 to 10.6 (van der Maas, 1993). The computer test was administered on 11 consecutive occasions (i.e. 11 test sessions). The time between two successive sessions varied from 2 to 4 weeks. Many test sessions (265 out of 1,111) were missing, however.

3.3. Monotonicity-constrained simple regression

According to Piaget's theory, the performance of children on the liquid conservation test is expected to be a monotone non-decreasing function of time except during the transitional disequilibrium stage, when relapses are possible. With respect to the latter, however, sparse empirical results show that relapses primarily occur when counter-suggestions or completely unfamiliar items are given (Inhelder, Sinclair, Bovet, & Wedgwood, 1974), which is not the case in the study of van der Maas (1993). Therefore, monotonicity can generally be assumed to hold for this study's data.

We check this assumption for four children with different overall performance levels on the liquid conservation test. Overall performance levels, which range from 1.67 to 8.00, are simply computed by averaging the child's scores across the sessions. The levels of the children selected are 3.1, 4.8, 5.9 and 7.5, respectively. For this analysis, we opt for a regression with 12 B-splines of the third degree and a second-order smoothness penalty. Regarding the smoothness weight λ , a value of 0.28 is chosen by making use of

Akaike's information criterion assuming independence of the four children and with the variance $\hat{\sigma}^2$ of the residuals in (8) being estimated on the basis of generalized cross validation. Regarding the monotonicity weight, κ , we choose a value as high as 10^6 to ensure that violations of the monotonicity assumption are negligible. The results, for both unconstrained and constrained regression, are graphically represented in Figure 5.

As a goodness-of-fit measure, we compute squared correlations between observed and predicted scores. If the data are approximately monotone increasing, unconstrained and constrained regression are expected to have a comparable fit; if not, constrained regression is expected to yield a much poorer fit than unconstrained regression. The results are summarized in Table 1. Only for child 1 is a significant discrepancy in fitted values observed, which can be explained by the fact that this child seems to be purely guessing; additional support for this explanation is the child's overall performance level of 3.1, close to the chance level of 2.7. For the other three children, who mainly differ with respect to the time period during which the transition from non-conserving to conserving occurs, the assumption of monotonicity seems justified.

3.4. Monotonicity-constrained regression with two predictor variables

Following Piaget's (1960) theory, it can be hypothesized that performance of children on the liquid conservation test is a monotone non-decreasing function of both time and overall level of performance of the children. This assumption of double monotonicity can be divided into two separate assumptions. First, the performance on the liquid

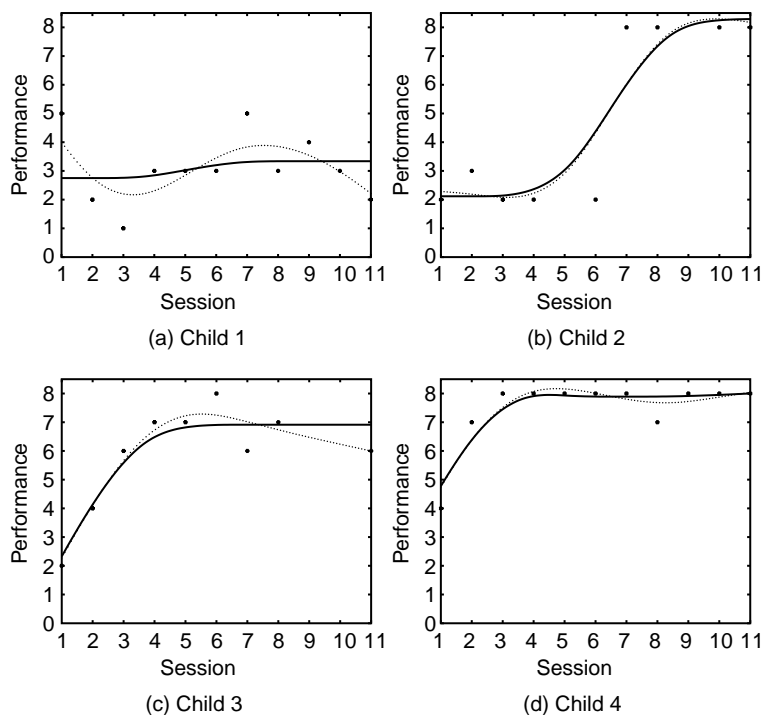


Figure 5. Results of unconstrained (dotted lines) and monotonicity-constrained (full lines) P-splines regression applied to performance on liquid conservation across sessions for four children with different overall performance levels (child 1, 3.1; child 2, 4.8; child 3, 5.9; child 4, 7.5).

Table 1. Squared correlations by child for unconstrained ($\kappa = 0$) and constrained regression ($\kappa = 10^6$)

κ	Child			
	1	2	3	4
0	0.66	0.86	0.93	0.89
10^6	0.055	0.855	0.875	0.87

conservation test is expected to be a monotone non-decreasing function of time conditional on the overall level of performance, in line with the argument and results of the previous section. Second, the performance on the liquid conservation test is expected to be a monotone non-decreasing function of overall level of performance conditional on time. The latter is implied by Piaget's assumption of a stepwise transition from the non-conserving to the conserving stage, with individual differences only occurring with respect to the moment of transition. Hence, at any moment in time, children with a high level of overall performance (i.e. children who are quicker to understand conservation) are expected to perform at least as well as children with a lower level of overall performance (i.e. children who are slower to understand conservation).

The assumption of double monotonicity is checked for all 101 children simultaneously. We opt for a regression with 20 bivariate B-splines of the third degree and second-order smoothness penalties. Regarding the smoothness weights λ_x and λ_y , with X referring to the overall performance level of the children and Y to the 11 different moments in time, values of 0.1 and 44.4, respectively, are chosen by making use of Akaike's information criterion with the variance $\hat{\sigma}$ of the residuals in (8) being estimated on the basis of generalized cross validation. Both monotonicity weights κ_x and κ_y are set at either 0 or 10^6 , resulting in four different analyses: (1) unconstrained regression, (2) monotonicity-constrained regression with respect to overall performance level, (3) monotonicity-constrained regression with respect to time, and (4) double monotonicity-constrained regression. The results are graphically displayed as surface plots in Figure 6 and as contour plots in Figure 7. The contour plots clearly reveal whether monotonicity is imposed. Figure 7a displays the model without monotonicity restrictions. Indeed, it can be seen that in both directions violations of the rank order, as displayed in the colour bar, occur. Figure 7b displays the model with monotonicity imposed on the overall level of performance. In this case, only in the vertical direction do violations of the rank order, as displayed in the colour bar, occur. On the other hand, only violations in the horizontal direction can be seen in Figure 7c, which displays the model with monotonicity imposed on time. Finally, in Figure 7d, which displays the model with monotonicity restrictions in both dimensions, no violations occur at all.

For each of the four regressions, a goodness-of-fit measure – that is, the squared correlation between observed and predicted scores – is computed. The results are summarized in Table 2. Since no significant discrepancies in fit are observed, the assumption of double monotonicity seems justified.

It is interesting to note that monotonicity-constrained models such as those represented in Figures 6 and 7 can be further explored by computing derivatives of the fitted function. The latter can be easily done by making use of expression (14). We

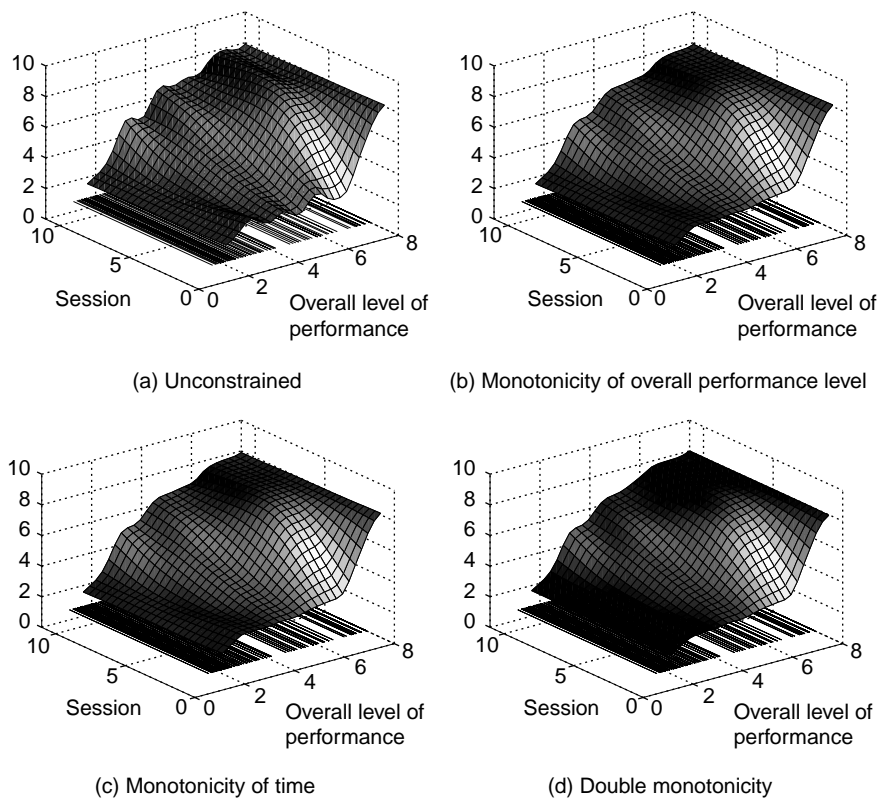


Figure 6. Surface plots of P-splines regression applied to performance of liquid conservation with sessions and overall level of performance as predictor variables: (a) Unconstrained regression, (b) monotonicity-constrained regression with respect to overall level of performance, (c) monotonicity-constrained regression with respect to time and (d) double monotonicity-constrained regression. Each individual child's overall level of performance is indicated with a black horizontal line.

illustrate this by computing the first-order derivative with respect to time of the results of the P-splines regression with monotonicity restrictions on time (see Figures 6c and 7c). Figure 8 displays the surface plot of the computed first-order derivative. Evidently, at any point, the first-order derivative is non-negative, which is as to be expected due to the monotonicity restrictions; values of zero for the first-order derivative indicate stagnation in learning, and positive values indicate learning – the

Table 2. Squared correlations for 2×2 different bivariate regression models

Constraints on Y	Constraints on X	
	No	Yes
No	0.79	0.78
Yes	0.78	0.78

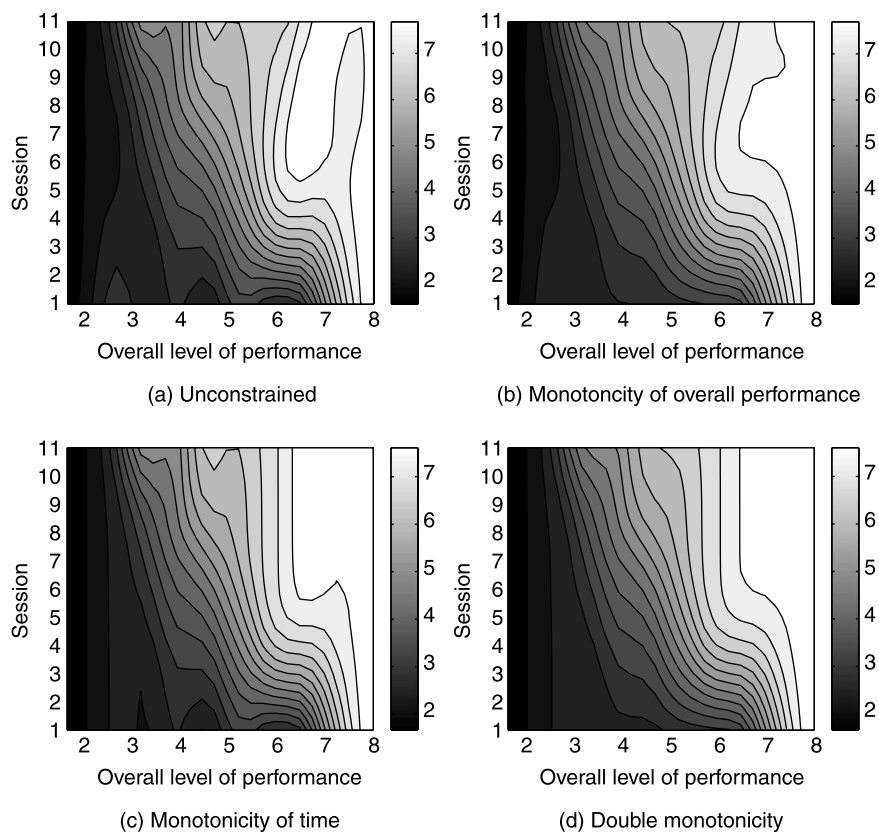


Figure 7. Contour plots of P-splines regression applied to performance of liquid conservation with sessions and overall level of performance as predictor variables: (a) Unconstrained regression, (b) monotonicity-constrained regression with respect to overall level of performance, (c) monotonicity-constrained regression with respect to time and (d) double monotonicity-constrained regression. On the right-hand side of each figure, a colour bar indicating the level of performance of liquid conservation is displayed.

higher these values, the faster the learning. As is clear from Figure 8, children with a lower overall level of performance learn later and more slowly as compared with children with a higher overall level of performance.

4. Concluding remarks

To check non-parametric functional forms that can be formalized as either local or global constraints on the sign of an n th-order derivative, we presented constrained P-splines regression. This is essentially non-parametric regression with additional asymmetric discrete penalties reflecting the assumed functional form. In particular, these penalties restrict the sign of the n th-order differences and, as such, the sign of the n th-order derivative.

For the sake of simplicity, in this paper we used a simple linear model (assuming normality and homoscedasticity) to demonstrate the use of asymmetric discrete penalties enforcing shape constraints. Of course, more complex models can be

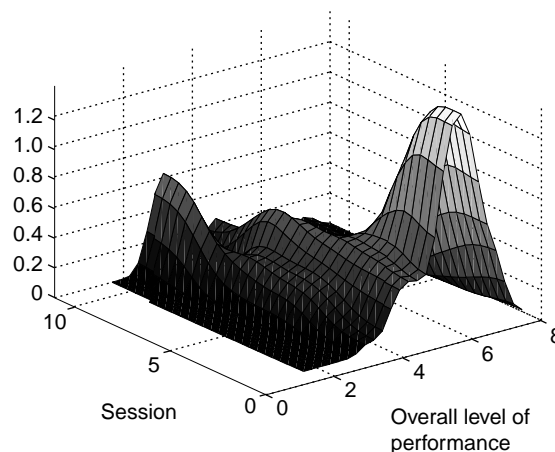


Figure 8. Surface plot of first-order derivative with respect to time of the results of P-splines regression with monotonicity restrictions on time.

considered as well. For instance, our approach can easily be extended to fit models within the framework of *generalized linear models*. In this way, a broad range of non-normal responses can be fitted using almost any monotone transformation of the linear predictor. For instance, Poisson regression using P-splines can be adopted to relax the assumption of homoscedasticity. In addition, the extension to *generalized additive models* (GAMs), in which the response is fitted using a sum of smooth functions, can be made as well. An introduction to GAMs using P-splines is given in Marx and Eilers (1998).

We extended regression with P-splines in order to impose certain shape constraints. Of course, other smoothing techniques can be used as well, some of which we will briefly discuss. In Winsberg and Ramsay (1980) and Ramsay (1988), monotone functions are estimated by taking positive combinations of integrated B-splines (I-splines). More recently, Hall and Huang (2001) suggest a 'tilting' algorithm for monotonizing general kernel-type estimators. The idea is to adjust an unconstrained estimate by tilting the empirical distribution so as to make the least possible change subject to the constraint of monotonicity. An application on monotonizing local polynomials using the *pool adjacent violator* algorithm can be found in Shkedy, Aerts, Molenberghs, Beutels, and Van Damme (2003). Ramsay (1998) estimates smooth monotone functions by solving appropriate differential equations. As can be seen, many different non-parametric approaches exist that can be used to impose shape constraints. However, an additional advantage of the constrained P-splines approach we introduced is its possible extension to the multidimensional setting and the ease with which different types of constraints can be imposed. Because the asymmetric penalties are computed from differences, local constraints, extensions to convex/concave smoothing and arbitrary combinations of monotone and convex/concave constraints are straightforward to implement. As an interesting example, consider the assumption of an ideal point of temperature. This means that a particular temperature (e.g. the ideal point) is judged as most pleasant whereas temperatures deviating from the ideal point are judged as less pleasant, with larger deviations being judged less pleasant. This assumption can be

checked with the help of P-splines regression with two additional asymmetric penalties: a first one that constrains the first-order derivative to be positive up to the ideal point, and a second one that constrains the first-order derivative to be negative after that ideal point. As a consequence, the ranges of the constraints need to be determined, which can be done either on the basis of prior theoretical considerations, or by making use of model selection techniques. The latter can be achieved by fitting several models with different ranges of the constraints and by subsequently selecting the best model. As such, constrained P-splines regression may be useful in checking the assumption of single-peakedness.

In general, constrained P-splines regression is a versatile approach that can easily be modified and used in many applied settings. It can be situated in between an exploratory and a confirmatory data-analytic approach, shifting more towards a confirmatory approach when higher values of the penalty weights are chosen. As such, the penalty weights constitute an interesting source of flexibility by which the method of constrained P-splines regression can easily be fine-tuned according to the researcher's purposes. In this regard, one possible approach that may be of interest consists of determining the minimal values for the constraint weights such that the corresponding assumptions are not violated. These values then indicate the extent to which the constraints fit the data, with lower values indicating a better fit. In this regard, it may also be of interest to investigate reference distributions under the null-model assumption that the assumed functional form holds perfectly, using, for instance, a non-parametric bootstrap type of procedure (Efron & Tibshirani, 1993).

Taken together, constrained P-splines regression constitutes a useful method that can be adapted easily in order to optimally investigate a broad range of substantively guided assumptions on functional forms. Constrained P-splines regression is not to be applied blindly but requires careful choices regarding the weight and the type of constraints.

Acknowledgements

The research reported in this paper was conducted while the first author was at Katholieke Universiteit Leuven. The authors are grateful to Patrick Groenen for his helpful suggestions on convexity and to Han van der Maas for placing his data at our disposal.

References

- Coombs, C. H. (1977). Single-peaked functions and the theory of preference. *Psychological Review*, 84, 216–230.
- De Boor, C. (1978). *A practical guide to splines*. Berlin: Springer.
- Dierckx, P. (1993). *Curve and surface fitting with splines*. Oxford: Clarendon.
- Durban, M., Currie, I. D. & Eilers, P. H. C. (2002). Using P-splines to smooth two-dimensional data. *Proceedings of 17th International Workshop on Statistical Modelling* (pp. 207–214), Crete.
- Efron, B., & Tibshirani, R. (1993). *An introduction to the bootstrap*. New York: Chapman & Hall.
- Eilers, P. H. C. (1994). Sign-constrained, monotone and convex nonparametric regression with asymmetric penalties. *Proceedings of the 9th International Workshop on Statistical Modelling*, Exeter.
- Eilers, P. H. C., Currie, I. D., & Durban, M. (2006). Fast and compact smoothing on large multidimensional grids. *Computational Statistics and Data Analysis*, 50, 61–76.
- Eilers, P. H. C., & Marx, B. D. (1996). Flexible smoothing using B-splines and penalized likelihood (with comments and rejoinder). *Statistical Science*, 11, 89–121.

- Hall, P., & Huang, L.-S. (2001). Nonparametric kernel regression subject to monotonicity constraints. *Annals of Statistics*, 29, 624–647.
- Hastie, T., & Tibshirani, R. (1990). *Generalized additive models*. London: Chapman & Hall.
- Inhelder, B., Sinclair, H., Bovet, M., & Wedgwood, S. (1974). *Learning and development of cognition*. London: Routledge & Kegan Paul.
- Marx, B. D., & Eilers, P. H. C. (1998). Direct generalized additive modeling with penalized likelihood. *Computational Statistics and Data Analysis*, 28, 193–209.
- McCrae, R. R., & Costa, P. T. (1987). Validation of the five-factor model of personality across instruments and observers. *Journal of Personality and Social Psychology*, 52, 81–90.
- O'Sullivan, F. (1988). A statistical perspective on ill-posed inverse problems (with discussion). *Statistical Science*, 1, 505–527.
- Piaget, J. (1960). *The psychology of intelligence*. Paterson, NJ: Littlefield, Adams.
- Ramsay, J. O. (1988). Monotone regression splines in action (with discussion). *Statistical Science*, 3, 425–461.
- Ramsay, J. O. (1998). Estimating smooth monotone functions. *Journal of the Royal Statistical Society, Series B*, 60, 365–375.
- Shkedy, Z., Aerts, M., Molenberghs, G., Beutels, P., & Van Damme, P. (2003). Modelling forces of infection by using monotone local polynomials. *Applied Statistics*, 52, 469–485.
- Strauss, S. (1982). *U-shaped behavioral growth*. New York: Academic Press.
- van der Maas, H. L. J. (1993). *Catastrophe analysis of stagewise cognitive development: Model, method and applications*. Unpublished doctoral dissertation, University of Amsterdam.
- van der Maas, H. L. J., & Molenaar, P. C. M. (1992). Stagewise cognitive development: An application of catastrophe theory. *Psychological Review*, 99, 395–417.
- Winsberg, S., & Ramsay, J. O. (1980). Monotonic transformations to additively using splines. *Biometrika*, 67, 669–674.

Appendix A

To show that the loss function given in (16) is convex, we first rewrite (16) in matrix notation:

$$L(\alpha) = (\mathbf{y} - \mathbf{B}\alpha)^T(\mathbf{y} - \mathbf{B}\alpha) + \lambda(\mathbf{D}_2\alpha)^T(\mathbf{D}_2\alpha) + \kappa(\mathbf{D}_n\alpha)^T\mathbf{I}\mathbf{W}_{(\alpha)}(\mathbf{D}_n\alpha) \quad (19)$$

or, equivalently,

$$L(\alpha) = \mathbf{y}^T\mathbf{y} - 2\mathbf{y}^T\mathbf{B}\alpha + \alpha^T\mathbf{B}^T\mathbf{B}\alpha + \lambda\alpha^T\mathbf{D}_2^T\mathbf{D}_2\alpha + \kappa\alpha^T\mathbf{D}_n^T\mathbf{I}\mathbf{W}_{(\alpha)}\mathbf{D}_n\alpha, \quad (20)$$

with $b_{ij} = B_j(x_i)$ being the elements of \mathbf{B} , \mathbf{D}_2 and \mathbf{D}_n being the matrix representations of the difference operators Δ^2 and Δ^n respectively, \mathbf{I} an $r \times r$ diagonal matrix containing the indicator variables v_j , and $\mathbf{W}_{(\alpha)}$ an $r \times r$ diagonal matrix containing the weights w_{α_j} .

The Hessian of the first term on the right-hand side of (19) equals $2\mathbf{B}^T\mathbf{B}$, which is positive semi-definite; hence, this term is convex in α . Similarly, the second term on the right-hand side of (19) is convex in α since its Hessian equals $2\lambda\mathbf{D}_2^T\mathbf{D}_2$, which is positive semi-definite as well. With respect to the third term, some additional explanation is needed. We first note that:

$$\kappa\alpha^T\mathbf{D}_n^T\mathbf{I}\mathbf{W}_{(\alpha)}\mathbf{D}_n\alpha = \sum_{j=n+1}^r \kappa(\mathbf{D}_n\alpha)_j^T v_j w_{\alpha_j} (\mathbf{D}_n\alpha)_j. \quad (21)$$

Then, for $v_j = 0$, $\kappa(\alpha^T\mathbf{D}_n^T)_j v_j w_{\alpha_j} (\mathbf{D}_n\alpha)_j = 0$, which is, of course, convex in α . For $v_j = 1$, the assignment of the weights w_{α_j} as in (13) implies that $\kappa(\alpha^T\mathbf{D}_n^T)_j v_j w_{\alpha_j} (\mathbf{D}_n\alpha)_j$ is

a truncated power function of second degree with argument $A = \Delta_{\alpha_j}^n$ of the form

$$\kappa(\boldsymbol{\alpha}^T \mathbf{D}_n^T)_j v_j w_{\alpha_j} (\mathbf{D}_n \boldsymbol{\alpha})_j = \begin{cases} 0, & \text{if } A \geq 0, \\ \kappa A^2, & \text{otherwise,} \end{cases} \quad (22)$$

which is convex in A . The latter expression is convex in $\boldsymbol{\alpha}$ as well since a convex function of a linear combination of elements of a vector is convex in that vector too; that is, if $f(\mathbf{x})$ is a convex function of \mathbf{x} , then $g(\mathbf{y}) = f(\mathbf{A}\mathbf{y})$ is a convex function of \mathbf{y} .

To prove the latter, assume $0 \leq \theta \leq 1$; then

$$\begin{aligned} g[\theta \mathbf{y}_1 + (1 - \theta) \mathbf{y}_2] &= f[\mathbf{A}[\theta \mathbf{y}_1 + (1 - \theta) \mathbf{y}_2]] \\ &= f[\theta \mathbf{A} \mathbf{y}_1 + (1 - \theta) \mathbf{A} \mathbf{y}_2] \\ &\leq \theta f(\mathbf{A} \mathbf{y}_1) + (1 - \theta) f(\mathbf{A} \mathbf{y}_2) \\ &= \theta g(\mathbf{y}_1) + (1 - \theta) g(\mathbf{y}_2) \end{aligned}$$

making use, in the third step, of the convexity of f .

Furthermore, since a sum of convex functions is convex, (20) is also convex in $\boldsymbol{\alpha}$. Finally, as we have shown that all three terms of (19) are convex in $\boldsymbol{\alpha}$, (19) is convex in $\boldsymbol{\alpha}$ as well.

Appendix B

To find an optimal solution of (16), we will make use of a Newton-Raphson procedure. At each iteration l , $\boldsymbol{\alpha}^{(l+1)}$ is computed such that

$$\mathbf{g}(\boldsymbol{\alpha}^{(l)}) + \mathbf{H}(\boldsymbol{\alpha}^{(l)})(\boldsymbol{\alpha}^{(l+1)} - \boldsymbol{\alpha}^{(l)}) = 0. \quad (23)$$

with \mathbf{g} being the gradient and \mathbf{H} the Hessian of $L(\boldsymbol{\alpha})$. The gradient of $L(\boldsymbol{\alpha})$ equals

$$\mathbf{g}(\boldsymbol{\alpha}) = -2\mathbf{B}^T \mathbf{y} + 2[\mathbf{B}^T \mathbf{B} + \lambda \mathbf{D}_2^T \mathbf{D}_2 + \kappa \mathbf{D}_n^T \mathbf{I} \mathbf{W}_{(\alpha)} \mathbf{D}_n] \boldsymbol{\alpha} + \kappa \boldsymbol{\alpha}^T \mathbf{D}_n^T \mathbf{I} \frac{\partial \mathbf{W}(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}} \mathbf{D}_n \boldsymbol{\alpha} \quad (24)$$

which can be simplified to

$$-2\mathbf{B}^T \mathbf{y} + 2[\mathbf{B}^T \mathbf{B} + \lambda \mathbf{D}_2^T \mathbf{D}_2 + \kappa \mathbf{D}_n^T \mathbf{I} \mathbf{W}_{(\alpha)} \mathbf{D}_n] \boldsymbol{\alpha}. \quad (25)$$

To see this, consider the following two cases for every j : (1) α_j^* , being values for α_j such that $\Delta_{\alpha_j}^n = 0$; and (2) values for $\alpha_j \neq \alpha_j^*$. Regarding case (1), both the left and the right limit of $g(\boldsymbol{\alpha})$ for α_j going to α_j^* equals (25). Regarding case (2), an infinite small change in α_j will not change the sign of $\Delta_{\alpha_j}^n$, and as such, will not change the corresponding weight. Hence, $\frac{\partial \mathbf{W}(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}}$ in (24) equals zero and as such, in this case, $g(\boldsymbol{\alpha})$ can be simplified to (25) as well.

Expression (25) shows that the gradient is a piecewise linear function of $\boldsymbol{\alpha}$. This implies that the Hessian is a step function of $\boldsymbol{\alpha}$ with discontinuities at α_j^* , $\alpha_j \neq \alpha_j^*$ the

Hessian equals

$$H(\alpha) = 2\mathbf{B}^T\mathbf{B} + 2\lambda\mathbf{D}_2^T\mathbf{D}_2 + 2\kappa\mathbf{D}_n^T\mathbf{I}\mathbf{W}_{(\alpha)}\mathbf{D}_n + 2\kappa\mathbf{D}_n^T\mathbf{I}\frac{\partial\mathbf{W}_{(\alpha)}}{\partial\alpha}\mathbf{D}_n. \quad (26)$$

As argued earlier, for $\alpha_j \neq \alpha_j^*$, $\frac{\partial\mathbf{W}_{(\alpha)}}{\partial\alpha}$ equals zero and as such, in this case, the Hessian $H(\alpha)$ can be simplified to

$$H(\alpha) = 2\mathbf{B}^T\mathbf{B} + 2\lambda\mathbf{D}_2^T\mathbf{D}_2 + 2\kappa\mathbf{D}_n^T\mathbf{I}\mathbf{W}_{(\alpha)}\mathbf{D}_n. \quad (27)$$

Furthermore, the function values of $H(\alpha)$ at α_j^* are uniquely defined as in (24) due to the allocation of the weights.

Then substituting (22) and (24) in (20) yields

$$\begin{aligned} & -2\mathbf{B}^T\mathbf{y} + 2[\mathbf{B}^T\mathbf{B} + \lambda\mathbf{D}_2^T\mathbf{D}_2 + \kappa\mathbf{D}_n^T\mathbf{I}\mathbf{W}_{(\alpha)}^{(l)}\mathbf{D}_n]\alpha^{(l)} \\ & + 2[\mathbf{B}^T\mathbf{B} + \lambda\mathbf{D}_2^T\mathbf{D}_2 + \kappa\mathbf{D}_n^T\mathbf{I}\mathbf{W}_{(\alpha)}^{(l)}\mathbf{D}_n](\alpha^{(l+1)} - \alpha^{(l)}) = 0. \end{aligned} \quad (28)$$

With some algebra, (25) can be simplified to

$$-\mathbf{B}^T\mathbf{y} + [\mathbf{B}^T\mathbf{B} + \lambda\mathbf{D}_2^T\mathbf{D}_2 + \kappa\mathbf{D}_n^T\mathbf{I}\mathbf{W}_{(\alpha)}^{(l)}\mathbf{D}_n]\alpha^{(l+1)} = 0. \quad (29)$$

and, hence,

$$\alpha^{(l+1)} = (\mathbf{B}^T\mathbf{B} + \lambda\mathbf{D}_2^T\mathbf{D}_2 + \kappa\mathbf{D}_n^T\mathbf{I}\mathbf{W}_{(\alpha)}^{(l)}\mathbf{D}_n)^{-1}\mathbf{B}^T\mathbf{y}; \quad (30)$$

the first two terms on the right-hand side have the same form as the normal equations for a least squares linear model, whereas the third term has the same form as the normal equations for a weighted least squares linear model, except that it has to be solved iteratively since \mathbf{W} depends on α .